# Springer Complexity

Springer Complexity is a publication program, cutting across all traditional disciplines of sciences as well as engineering, economics, medicine, psychology and computer sciences, which is aimed at researchers, students and practitioners working in the field of complex systems. Complex Systems are systems that comprise many interacting parts with the ability to generate a new quality of macroscopic collective behavior through self-organization, e.g., the spontaneous formation of temporal, spatial or functional structures. This recognition, that the collective behavior of the whole system cannot be simply inferred from the understanding of the behavior of the individual components, has led to various new concepts and sophisticated tools of complexity. The main concepts and tools – with sometimes overlapping contents and methodologies – are the theories of self-organization, complex systems, synergetics, dynamical systems, turbulence, catastrophes, instabilities, nonlinearity, stochastic processes, chaos, neural networks, cellular automata, adaptive systems, and genetic algorithms.

The topics treated within Springer Complexity are as diverse as lasers or fluids in physics, machine cutting phenomena of workpieces or electric circuits with feedback in engineering, growth of crystals or pattern formation in chemistry, morphogenesis in biology, brain function in neurology, behavior of stock exchange rates in economics, or the formation of public opinion in sociology. All these seemingly quite different kinds of structure formation have a number of important features and underlying structures in common. These deep structural similarities can be exploited to transfer analytical methods and understanding from one field to another. The Springer Complexity program therefore seeks to foster cross-fertilization between the disciplines and a dialogue between theoreticians and experimentalists for a deeper understanding of the general structure and behavior of complex systems.

The program consists of individual books, books series such as "Springer Series in Synergetics", "Institute of Nonlinear Science", "Physics of Neural Networks", and "Understanding Complex Systems", as well as various journals.

# Understanding Complex Systems

## Understanding Complex Systems

Future scientific and technological developments in many fields will necessarily depend upon coming to grips with complex systems. Such systems are complex in both their composition (typically many different kinds of components interacting with each other and their environments on multiple levels) and in the rich diversity of behavior of which they are capable. The Springer Series in Understanding Complex Systems series (UCS) promotes new strategies and paradigms for understanding and realizing applications of complex systems research in a wide variety of fields and endeavors. UCS is explicitly transdisciplinary. It has three main goals: First, to elaborate the concepts, methods and tools of self-organizing dynamical systems at all levels of description and in all scientific fields, especially newly emerging areas within the Life, Social, Behavioral, Economic, Neuro- and Cognitive Sciences (and derivatives thereof); second, to encourage novel applications of these ideas in various fields of Engineering and Computation such as robotics, nano-technology and informatics; third, to provide a single forum within which commonalities and differences in the workings of complex systems may be discerned, hence leading to deeper insight and understanding. UCS will publish monographs and selected edited contributions from specialized conferences and workshops aimed at communicating new findings to a large multidisciplinary audience.

S. Baglio   A. Bulsara   (Eds.)

# Device Applications
# of Nonlinear Dynamics

With 124 Figures

Springer

Professor Salvatore Baglio
Associate Professor
Microsystems Group
Dipartimento di Ingegneria Elettrica
Elettronica e dei Sistemi
University of Catania
V.le A. Doria 6
95125, Catania, Italy
E-mail: salvatore.baglio@diees.unict.it
Web: www.diees.unict.it/baglio

Dr. Adi Bulsara
Space and Naval Warfare Systems Center
Code D-363
49590 Lassing Road
San Diego, California 92152-6147
USA
E-mail: bulsara@spawar.navy.mil

# Preface

The past two decades have witnessed an explosion of ideas in the general field of nonlinear dynamics.

In fact, it has become increasingly clear that areas as diverse as signal processing, lasers, molecular motors, and biomedical anomalies have a common underlying thread: the dynamics that underpin these systems are inherently nonlinear. Yet, while there has been significant progress in the theory of nonlinear phenomena under an assortment of system boundary conditions and preparations, there exist comparatively few devices that actually take this rich behavior into account.

In the presence of background noise (a given, for most practical applications), the underlying dynamic phenomena become even richer, with the noise actually mediating cooperative behavior that, when properly understood, can lead to significant performance enhancements; a striking example of this behavior occurs, for example, when the underlying dynamics undergoes a bifurcation from static to oscillating behavior when a control parameter is swept through a critical value. If properly understood, theoretically, the (suitably quantified) system response can be significantly enhanced near the onset of the bifurcation. Examples of this behavior have been observed in a large number of laboratory experiments on systems ranging from solid state lasers, to SQUIDs, and such behavior has been hypothesized to account for some of the more striking information processing properties of biological neurons. In many cases, background noise can precipitate this behavior, thereby playing a significant role in the optimization of the response of these systems to small external perturbations.

A series of meetings on topics such as Stochastic Resonance, Experimental Chaos, and Neural Coding, have attempted to bring together researchers in this field (of applied nonlinear dynamics), whose ramifications cut across "party lines", however, there has not been, to date, a meeting that brings together researchers who are actually attempting to apply and exploit this knowledge to make devices which operate more efficiently (e.g. without complicated nulling circuits to essentially block the effects of the nonlinearity) and

cheaply, while affording the promise of much better performance. Given the current explosion of ideas in areas as diverse as molecular motors, nonlinear filtering theory, noise-enhanced propagation, stochastic resonance and (perhaps most important) networked systems the time was deemed to be right, particularly in this dawn of the era of nanotechnology, to have a meeting of researchers who are actually attempting to integrate some of these ideas into real devices.

The plenary and invited speakers at this meeting were drawn from a rarified mix. They included a few well-established researchers in the theory of noisy nonlinear dynamic systems, as well as a "new breed" of pioneers (applied physicists, engineers, and biologists) who are attempting to apply these ideas in actual laboratory (and in some cases industrial) applications.

Again, micro- and nano-technology was an area of emphasis; however, the idea was to present developments in a wide spectrum of practical areas. Preference was given to researchers who are actually involved in fabricating and experimenting with devices that are left to operate in their free-running (i.e. inherently nonlinear dynamical) configuration, rather than being "quasi-linearized" as mentioned above.

In addition, the organizers attempted to give some exposure to much younger researchers (advanced graduate students and postdocs) in the form of short contributed talks or posters. Plenty of time for discussions and sightseeing was available.

The goal faced in the organization of this conference was quite ambitious due to the virtually limitless set of areas that can be covered with the nonlinear dynamics umbrella, however the scientific coordination efforts were focused to identified a number of themes, yet still quite large, as can be derived from the technical program of the DANOLD conference that corresponds almost completely to the table of contents of this proceeding book.

The meeting may be regarded as the logical successor to the 1997 ANDM (Applied Nonlinear Dynamics Near the Millennium) meeting, held in San Diego. That meeting brought together researchers from physics, engineering, biology and the social sciences, who were involved in the analysis (and in a few cases experiments with) of applied nonlinear dynamical systems.

This meeting was intended to lead to the next level; it concentrated on the implementation of these ideas into actual devices and systems. However, realizing that theoretical ideas and discoveries march to the beat of their own drum, the meeting also featured, as already stated, some novel theoretical ideas that have not yet made it to the drawing board, but show great promise for the future.

The meeting was held in beautiful Catania, Sicily. In addition to lying at the base of an active volcano, and being home to a number of active and prolific research groups (at the Universita degli Studi di Catania, as well as numerous local high-tech companies), Catania is close to historical and architectural gems such as Taormina, Siracusa, and Agrigento. It provided a

truly beautiful atmosphere and setting to conduct a meeting of this kind, on a subject that is poised on the threshold of an explosion.

The organizers extend their sincerest thanks to the principal sponsors of the meeting: SPAWAR Systems Center (San Diego), Office of Naval Research-Global (London), and the Engineering Faculty of the Università degli Studi di Catania, who, through their support, enabled the organization of a quality meeting without too much of a financial burden placed on individual participants. In addition, many thanks are due to the local sponsors that have allowed to enrich the social program of the conference and, last but not least, our grateful thanks go to the collegues that have chaired the sessions and to the numerous individuals (postdocs, graduate students, and secretaries) who donated long hours of labor to the success of the meeting.

A special thank is for Dr. Bruno Andò and Dr. Visarath In who have helped us in a precious and priceless manner throughout this all adventure.

Finally, we thank Springer-Verlag for their production of an elegant proceedings book.

Catania, Italy                                                            *S. Baglio*
San Diego, USA                                                    *A.R. Bulsara*
March 2006

# Contents

# Use of Chaos to Improve Equipments

L. Fortuna and M. Frasca

Dipartimento di Ingegneria Elettrica Elettronica e dei Sistemi, Università degli Studi di Catania, viale A. Doria 6, 95125 Catania, Italy
`lfortuna@diees.unict.it`

In this communication three applications where the use of chaos improves the device are dealt with. The first application concerns the use of chaos to drive sonar sensors in multi-user scenarios. The second application deals with the use of chaos to enhance motion control of a microrobot. The third application deals with a new synchronization scheme for chaotic systems.

## 1 Application of Chaos to Sonars

Ultrasonic devices are widely used in robotics as exteroceptive sensors for ranging measurements. These applications involve a large number of sonars operating concurrently, giving rise to the phenomenon of crosstalk. The first application presented in this work aims to exploit the peculiarity of chaos to enhance the performance of sonar systems in terms of crosstalk and noise rejection. In particular, since all chaotic systems share properties [1] such as sharp autocorrelation functions and uncorrelation between signals coming from different systems as well as signals coming from different attractors of the same system, chaos is a suitable paradigm in any application in which a unique signature for a source of information is needed. The main idea underlying this application [4] is therefore to drive a sonar with suitable chaotic signals and apply a matched filter technique for a robust rejection of crosstalk and noise. The great advantage of chaotic signals is that they can be easily generated with a low cost analog circuitry. Therefore, the sonar system presented in this work does not require digital units to drive the sensor. Moreover, the uniqueness of sequences is guaranteed by the properties of the circuitry, and no coordination or supervision units are required to avoid crosstalk.

Ranging measurements based on sonars are usually performed by measuring the *TOF* of an ultrasound wave propagating in air. In low cost, autonomous robot applications, the most used sonar sensor is perhaps the *Polaroid Series 600,* driven by the *Polaroid 6500 Ranging Module* [3].

The approach adopted in this work conjugates the ideas of exploiting both the time intervals between the pulses emitted by the sonar and the peculiar characteristics of chaotic signals to build a unique signature belonging to each sensor. The main idea underlying the chaotic modulation used is to generate a sequence of pulses in which the duration of the time interval between a pulse and the next one is provided by a chaotic law. As the information is contained on the temporal distance between pulses, additive noise on the channel does not affect the integrity of the information. Moreover, pulses with a small duty cycle are used, thus involving low power consumption.

In particular, the chaotic sequence is generated on the basis of a continuous chaotic attractor via a voltage-to-time conversion implemented in a modulator circuit designed for the purpose of continuous chaotic pulse position modulation (CPPM). In our case, the continuous circuit generating chaos is the well known *Chua's Circuit* [2].

The experimental setup has been built with a Polaroid series 600 sensor. The power circuitry has been obtained by using only the output stage of the Polaroid series 6500 ranging module and inhibiting the remainder of the board. The train of pulses has been emitted through continuous CPPM driven by a Chua's circuit evolving according to a *double scroll Chua attractor* [2].

The sensor has been characterized for measurements ranging from 5 cm to 145 cm, in 5 cm steps. Three sets of experiments have been carried out. The first one refers to measurements performed by using the sonar driven by CPPM, in a single-user scenario. The second one refers to the same experiment, in presence of another CPPM sensor located close to the sensor to be characterized (two-user scenario). The third one refers to measurements performed by a sensor driven by an original Polaroid 6500-Series Sonar Ranging Module. It is worth remarking that in this case, measurements of distance under 40 cm require specific techniques to damp the echo of the transmitted signal, which would lead to incorrect measurement. This is prevented by the constructor by introducing a blank interval of 238 ms, when the sensor is inhibited. Table 1 reports the average measurement error committed in the three cases. To make a comparison, the table has been worked out by considering the range 40 cm–145 cm. It is worth noticing that the CPPM approach allows to perform measurements under 40 cm without adopting any particular technique to damp the echo of the transmitted signal. The signal transmitted is in fact entirely received by the sensor and gives rise to a large peak of correlation at the origin of time, which can be ignored. In conclusion, the CPPM approach allows us to perform ranging measurements with an error comparable with that committed by the Polaroid Ranging module, despite of the presence of crosstalk and noise, thus obtaining a better overall performance.

**Table 1.** Average error committed in the range 40 cm–145 cm

|      | Polaroid Ranging Module | CPPM  | Multi-user CPPM |
|------|:-----------------------:|:-----:|:---------------:|
| Mean |         2.29%           | 1.87% |     1.84%       |

## 2 Chaos to Improve Motion Control of Microrobots

The second application deals with the use of chaos for motion control in microrobotics. In particular, a microrobot actuated by piezoelectric elements, named PLIF (Piezo Light Intelligent Flea) [5], designed to be fast, small, light and cheap, is taken into account and chaos is used to enhance the motion capabilities on irregular surfaces.

Usually, the actuation of robot legs is controlled by square wave signals characterized by a fixed amplitude and a variable switching frequency. In this application these signals are generated performing a frequency modulation driven by the chaotic evolution of Chua's circuit state variables. The smooth changes of the actuation signal frequency, performed by our chaotic system, enhances the microrobot walking capabilities especially when walking on irregular surfaces. Indeed, when driven with a constant frequency control signal, the microrobot is able to walk on regular surfaces if the frequency is appropriately tuned, but very small irregularities (such as grazes) can be a serious problem for the microrobot. By exploiting the widespread spectrum of a chaotic signal, a control signal with erratically varying frequency is provided to the robot making it able to deal with asperities in the surface and adaptable to different surfaces. In fact, in our microrobot chaos is directly used in the actuation system to modulate the signals devoted to the robot control.

The actuation of the microrobot used in this application is based on piezoelectric ceramic actuators. Piezoelectric materials are particular structures able to produce a voltage when deformed and, viceversa, an excitation voltage induces a deformation that can generate a force. Hence, it is possible to use piezoelectric materials as deformation sensors as well as actuators. The piezoelectric actuator is made up of two piezoceramics joined and isolated through a resin coverage. The two elements are excited alternatively: one of the elements, excited, shortens while the other one stretches making the entire structure bending toward the short side. To recover the original position it is sufficient to reverse the excitation voltage. The piezoelectric actuators are used to build the legs of the robot; each leg is therefore actuated by a flexor-extensor-like pair. The whole structure of the PLIF robot, designed to be light and as small as possible, is shown in Fig. 1(a).

The motion control system generates and controls the locomotion pattern of the microrobot which preliminary experimental tests have been revealed to be the most effective for the adopted structure. The motion pattern is characterized by the simultaneous actuation of the two legs. Each robot leg

(a)                                                 (b)

**Fig. 1.** (**a**) The PLIF microrobot structure. (**b**) Block scheme of the electronic board

follows this movement sequence: femur raising; tibia moving forward; femur going down; tibia moving backward.

In [5] and related works, the locomotion pattern is realized by an oscillator which generates a square wave signal with constant frequency and by a power circuitry (driver) providing the voltage supply needed by the piezoelectric actuator. In this application, in order to provide the robot with adaptive capabilities, this control scheme has been modified as shown in Fig. 1(b), where chaotic modulation of the control signal is included. In this way, the generation of control signals with time-variant frequency can be accomplished. The frequency of the control signal changes as function of the state variables of a chaotic circuit. Thus, the unpredictable behavior of the chaotic modulating signal is exploited to obtain a control system able to explore at each step new solutions to the motion control problem. In particular, one of the state variables of a Chua's circuit is used as modulation driving signal.

In order to evaluate the performances of the microrobot, three kind of tests were performed. In the first set of tests the microrobot walks on different smooth surfaces like an iron or wooden layer. In the second set of tests the surfaces are grazed in order to compare the performances in terms of speed obtained with or without chaotic modulation. The last set of tests concerns the overloading of the microrobot structure in order to verify if the introduction of chaotically modulated control signals is able to improve the motion also in presence of heavy structures.

Comparative results show that driving the actuation by using chaotic modulation leads to consistent improvements in terms of two factors: robot speed and motion on irregular surfaces. In particular on grazed surfaces, the robot, driven by chaotically modulated signals, is able to pass over the scratches while in the case of constant frequency actuation signal the robot often stops or decreases its velocity. To graphically show the improvements obtained using chaotically modulated frequency signal, the robot has been equipped with a led that lights up when the robot is actuated. A camera with a long exposure time has been used in order to take pictures which traces the robot trajectory. As shown in Fig. 2, improvements are clearly visible simply comparing the trajectory of the red led.

**Fig. 2.** Performances on an scratched wooden surface. On the left side the constant frequency case is reported, on the right side the chaotically modulated frequency case

# 3 Separation and Synchronization of Chaotic Circuits

The third application described in this communication deals with a new synchronization scheme for chaotic systems. In the classical scheme based on negative feedback [7], starting from the difference of two corresponding state variables (which are assumed measurable), an error signal is built and fed back into the slave system. In this application the synchronization of two pairs of chaotic systems instead of two chaotic systems by using a negative feedback scheme is investigated. In our case, thus, the master system is formed by two independent chaotic systems (i.e. two different systems which do not interact each other). In general, the synchronization of two pairs of such chaotic systems requires two independent feedback signals. In our case, instead, the question if and under which conditions synchronization can be achieved by using only a feedback signal which depends on both the two chaotic systems of the master (i.e. it is for instance a linear combination of the state variables of the two master chaotic systems) has been investigated. We refer to this problem as separation and synchronization of chaotic signals.

In particular, the problem of separation and synchronization for a class of chaotic systems, namely those with piece-wise linear (PWL) nonlinearities, is investigated with an approach based on linear matrix inequalities (LMI) [6].

The main idea underlying the application is the following. Chaotic systems characterized by PWL nonlinearities are considered: in each region of the PWL, the systems of this class assume different linear behavior switching through the PWL regions. Therefore, a PWL system is characterized by the set of its possible linearizations. Since, in each region, each linear system can be observed using the classical linear control techniques, our idea is to design an observer which simultaneously guarantees asymptotically stable error dynamics in each of these regions. Therefore, to solve the problem of separation and synchronization, the observer should be designed by solving a simultaneous stability problem. This can be done by formulating an LMI problem. If this problem is feasible, the corresponding problem of separation

(a)  (b)

**Fig. 3.** Separation and synchronization of a pair composed by a Chua's circuit and a Kennedy oscillator. (**a**) Synchronization plot $x_{1m}$ vs. $x_{1s}$. (**b**) Synchronization plot $x_{2m}$ vs. $x_{2s}$

and synchronization may admit a solution. Moreover, numerical results often show that this condition is also sufficient for the existence of a solution to the separation and synchronization problem.

An example related to the synchronization of a pair of chaotic systems made of a Chua's circuit [2] and a Kennedy's oscillator [8] is reported. In Fig. 3 the synchronization plots related to corresponding variables of the master and slave Chua's circuit (named $x_{1m}$ and $x_{1s}$, respectively) and to the corresponding variables of the master and slave Kennedy's oscillator (named as $x_{2m}$ and $x_{2s}$, respectively) are shown emphasizing that the two systems are perfectly synchronized.

Experimental results confirm the suitability of the approach even in the real case, when nonidentical systems are necessarily considered.

The solution to the problem of separation and synchronization, introduced in this application, can be adopted in chaotic communication systems. For instance, in order to enlarge the bandwidth of the communication channel it might be possible to use two different chaotic carriers transmitting two different information at the same time. In this case, the two chaotic systems of the slave have to be synchronized to the two chaotic systems of the master, starting from only one signal containing the carriers and the data mapped on them. This can be achieved applying the proposed separation and synchronization scheme.

## 4 Conclusions

Since the discovery of chaos in physical systems, chaotic behavior has been also observed in many engineering fields. For instance, many electrical and electronic systems such as DC-DC converters may show chaotic behavior. In

such systems the chaotic behavior is often an undesirable behavior which the designer should avoid. More recently, another point of view has gained interest and possible applications of chaos (and in particular of chaotic circuits) to engineering problems are searched for. In this communication three applications from different application fields in which the use of chaos improves the performance of the device have been presented. In particular, the use of chaos to drive sonar sensors in multi-user scenarios, the use of chaos to enhance motion control of a microrobot and a new synchronization scheme for chaotic systems have been discussed.

# References

1. Strogatz S H (1994) *Nonlinear Dynamics and Chaos.* Perseus Book, Oxford
2. Madan R N (1993) *Chua's circuit: a paradigm for chaos*, World Scientific Series on Nonlinear Sciences, Series B, Vol. 1 (World Scientific, Singapore)
3. Everett H R (1995) *Sensors for Mobile Robots – Theory and Application.* A. K. Peters Ltd., Natick, MA
4. Fortuna L, Frasca M, Rizzo A (2003) *IEEE Trans. Instrumentation and Measurement* 52: 1809–1814
5. De Ambroggi F, Fortuna, L, Muscato G (1997) PLIF: Piezo Light Intelligent Flea. New micro-robots controlled by self-learning techniques. In *Proc. of the 1997 IEEE Int. Conf. Robotics and Automation*
6. Boyd S, El Ghaoui L, Feron E, Balakrishnan V (1994) Linear Matrix Inequalities in System and Control Theory. SIAM Books
7. Kapitaniak T (1994) *Phys. Rev. E* 50: 1642–1644
8. Elwakil A S, Salama K N, Kennedy M P (2000) A system for chaos generation and its implementation in monolithic form. In Proceedings of IEEE Int. Conf. of Circuits and Systems, ISCAS 2000

# Invited Papers

# Noise Induced Switching Between Oscillation States in a Nonlinear Micromechanical Oscillator

H.B. Chan and C. Stambaugh

Department of Physics, University of Florida, Gainesville, FL 32608, USA

The interplay of noise and nonlinearity often leads to novel phenomena in micro- and nano-systems. Such phenomena are of fundamental and practical interest since they have the potential to offer new functionalities and to improve the performance of sensors. For instance, nonlinear systems often develops bistability when the periodic driving is sufficiently strong. In the presence of fluctuations, the nonlinear system can be induced to escape from one metastable dynamical state into the other. Unlike equilibrium systems where the escape rate can be calculated from the height of the free-energy barrier [1], driven systems are, generally, far from thermal equilibrium and cannot be characterized by free energy [2–5]. Consequently, the escape rate in such non-equilibrium systems must be derived from system dynamics [6–9]. Experiments on noise induced switching has been performed in a number of driven nonlinear systems, including parametrically driven electrons in a Penning trap [10], doubly clamped nanomechanical beams [11, 12] and radio frequency driven Josephson junctions [13]. Calculation of the escape rate in such non-equilibrium systems is a non-trivial problem and has attracted much interest [7, 8, 14, 15].

Here we describe our investigation of noise-activated switching in an underdamped micromechanical torsional oscillator periodically driven into nonlinear oscillations. Within a certain range of driving frequencies, the oscillator has two stable dynamical states with different oscillation amplitudes. We induce the oscillator to escape from one state into the other by injecting noise in the driving force. By measuring the frequency of random transitions as a function of noise intensity, we extract the transition rate and demonstrate the activated behavior for switching.

In our experiment, the micromechanical oscillator is fabricated using a silicon surface micromachining process. The oscillator consists of a polysilicon plate (500 µm by 500 µm by 3.5 µm) that is supported by two torsional rods. After a 2-µm-thick sacrificial silicon oxide layer beneath the plate is etched away, the plate becomes free to rotate about the torsional springs. Figures 1a and 1b are scanning electron micrographs of a typical oscillator. The moment

**Fig. 1.** (**a**) Scanning electron micrograph of the micromechanical torsional oscillator. The *large square* in the middle is a movable polysilicon top plate. The *small squares* are bond pads that provide electrical connections to the top plate and the two underlying electrodes. (**b**) Close up on one of the torsional springs. (**c**) A cross sectional schematic of the torsional oscillator with electrical connections and measurement circuitry

of inertia $I$ for the plate is $4.3 \times 10^{-17}$ kg m$^2$, while the spring constants are each $9.2 \times 10^{-9}$ N m rad$^{-1}$. The other ends of the springs are anchored to the silicon substrate. Beneath the top plate, there are two fixed electrodes (500 μm by 250 μm) on each side of the torsional springs. One of the electrodes is used to excite the torsional oscillations electrostatically and the other electrode is used to measure the response.

Figure 1c shows a cross sectional schematic of the oscillator with electrical connections and measurement circuitry. The application of a periodic voltage with dc bias $V_{dc1}$ to one of the electrodes leads to an electrostatic attraction between the grounded top plate and the electrode. Torsional oscillations of the top plate are excited by the periodic component of the electrostatic torque. The detection electrode is connected to a dc voltage $V_{dc2}$ through a resistor $R$. As the plate oscillates, the capacitance between the plate and the detection electrode changes. The detection electrode is connected to a charge sensitive preamplifier followed by a lock-in amplifier that measures the signal at the excitation frequency. Measurements were performed at liquid helium temperature and at pressure of less than $2 \times 10^{-7}$ torr. The quality factor $Q$ of the oscillator is about 16,000.

The nonlinear behavior in our torsional oscillator originates mainly from the strongly distance dependent electrostatic interaction between the top plate

and the excitation electrode. The equation of motion of the oscillator is given by [16]:

$$\ddot{\theta} + 2\gamma\dot{\theta} + \omega_o^2\theta = \tau/I \tag{1}$$

where $\theta$ is the angular rotation of the top plate, $\gamma$ is the damping coefficient, $\omega_o$ is the natural frequency of the oscillator and $\tau$ is the driving torque. $\tau$ arises from the electrostatic interaction between the top plate and the driving electrode. When the rotation angle $\theta$ is small, $\tau$ can be written as:

$$\tau = bF \tag{2}$$

where $b$ is half the distance from the axis of rotation to the edge of the top plate and the electrostatic force $F$ is evaluated with separation $d - b\theta$:

$$F = \frac{\varepsilon_o A V^2}{2(d - b\theta)^2} \tag{3}$$

As shown in Fig. 1c, $d$ is the separation between the top plate and the electrode when no electrostatic force is applied, $A$ is the area of the electrode and $\varepsilon_o$ is the permittivity of free space. The excitation voltage $V$ is a sum of three components:

$$V = V_{dc1} + V_{ac}\sin(\omega t) + V_{noise}(t) \tag{4}$$

The three terms on the right side of (4) represent the dc voltage, periodic ac voltage and random noise voltage respectively. $V_{dc1}$ is chosen to be much larger than $V_{ac}$ and $V_{noise}$ to linearize the dependence of $F$ on $V_{ac}$ and $V_{noise}$. The strong spatial dependence of the electrostatic force leads to nonlinear contributions to the restoring torque. A Taylor expansion of the electrostatic force $F$ about $d$ gives:

$$F(d - b\theta) = F(d) - F'(d)b\theta + \frac{F''(d)(b\theta)^2}{2} - \frac{F'''(d)(b\theta)^3}{6} \tag{5}$$

where $F'$, $F''$ and $F'''$ denotes the first, second and third spatial derivative of $F$ respectively. Substituting $F(d - b\theta)$ in (1) and (2) leads to:

$$\ddot{\theta} + 2\gamma\dot{\theta} + [\omega_0^2 - \eta]\theta + \beta\theta^3 + C = E\sin(\omega t) + N(t) \tag{6}$$

where $\eta = \left(b^2\varepsilon_o A/Id^3\right)V_{dc1}^2, \beta = -\left(2b^4\varepsilon_o A/Id^5\right)V_{dc1}^2, C = -\left(b\varepsilon_o A/2Id^2\right)V_{dc1}^2,$ $E = \left(b\varepsilon_o A V_{dc1}/2Id^2\right)V_{ac}$ is the effective amplitude of the periodic excitation, and $N(t) = \left(b\varepsilon_o A V_{dc1}/2Id^2\right)V_{noise}(t)$ is the effective noise in the excitation. The contribution of the quadratic term to the nonlinearity proves to be negligible compared to the cubic term, while the linear term produces a shift in the natural frequency.

We first focus on the response of the oscillator with no injected noise in the excitation. Figure 2 shows the frequency response of the oscillator at two different excitation amplitudes. For both cases the peak oscillation amplitudes are normalized to unity. The squares represent the response of the oscillator

**Fig. 2.** Normalized frequency response of the oscillator for excitation voltages amplitudes of 23 μV (*solid squares*) and 91 μV (*hollow circles*). The *dotted line* represents a fit to the data at smaller driving force using the response of a damped harmonic oscillator. For the large driving force, two dynamical states coexist from 3296.45 Hz to 3297.40 Hz. The *dashed line* fits the data to a damped oscillator with cubic non-linearity [16], yielding $\beta = -1.3 \times 10^{14}$ rad$^{-2}$ s$^{-2}$

when the excitation is small. The resonance peak is fitted well by the dotted line that corresponds to the response of a damped harmonic oscillator. As the periodic excitation is increased, contributions from the cubic term in (6) become more important and lead to nonlinear behavior in the oscillations. The resonance curve becomes asymmetric, with the peak shifting to lower frequencies. At a high enough excitation, the frequency response becomes hysteretic, as shown by the circles in Fig. 2. Within the certain range of driving frequencies from 3296.4 Hz to 3297.4 Hz, there are two stable dynamical states with different oscillation amplitudes. The system resides in the high-amplitude state or the low-amplitude state depending on the history of the oscillator. In the absence of fluctuations, the oscillator remains in one of the stable states indefinitely.

In the presence of noise in the excitation, the oscillator could be induced to escape from one state into the other. Since this driven, bistable system is far from thermal equilibrium and cannot be characterized by free energy, calculation of the escape rate is a non-trivial problem. Theoretical analysis [6, 7] suggests that the rate of escape $\Gamma$ at a particular driving frequency depends exponentially on the ratio of an activation energy $E_a$ to the noise intensity I$_N$:

$$\Gamma = \Gamma_0 \exp(-E_a/\mathrm{I}_N) \tag{7}$$

Close to the bifurcation frequency where the high-amplitude state disappears, the activation energy is expected to display system-independent scaling. For cubic nonlinearity, the activation energy is given by:

$$E_a \propto \Delta\omega^\alpha \tag{8}$$

where the frequency detuning $\Delta\omega$ is the difference between the driving frequency and the bifurcation frequency. The activation energy is predicted [6,7] to increase with frequency detuning with critical exponent $\alpha = 3/2$ for all dynamical systems with cubic nonlinearity. We will describe below our comprehensive experimental investigation of activated switching from the high-amplitude to the low-amplitude state of the micromechanical oscillator.

In our experiment, transitions from the high-amplitude state to the low-amplitude state are induced by injecting noise in the excitation with a bandwidth of 100 Hz centered about the resonant frequency. The bandwidth of the noise is much larger than the width of the resonance peak. Figure 3a shows typical switching events where the oscillator resides in the high amplitude state for various durations before escaping to the low amplitude state. Due to the random nature of the transitions, a large number of switching events must be recorded to determine the transition rate accurately. During the time interval between switching events in Fig. 3a, the oscillator is reset to the high amplitude state using the following procedure. First, the noise is turned off and the driving frequency is increased beyond the range of frequencies where bistability occurs (>3297.4 Hz as shown in Fig. 2). The driving frequency is then decreased slowly towards the target frequency so that the oscillator remains in the high-amplitude state. Once the target frequency is reached, the noise is turned back on and the time for the oscillator to escape from the high-amplitude state is recorded. This process is then repeated multiple times to accumulate the statistics for the switching. This procedure is necessary because the energy barrier for transitions from the low-amplitude state back to the high-amplitude state is much larger than the barrier for transition in the opposite direction. Thus, noise induced transitions from the low-amplitude to



**Fig. 3.** (**a**) In the presence of noise in the excitation, the oscillator switches from the high amplitude state to the low amplitude state at different time intervals. The system is reset to the upper amplitude state between switching events. (**b**) Histogram of the residence time in the upper state before switching occurs, at detuning frequency of 0.15 Hz. The *dotted line* is an exponential fit

the high-amplitude state will fail to occur in the duration of the experiment and the oscillator must be reset to the high-amplitude state using the steps described above. Figure 3b shows a histogram of the residence time $T_R$ in the high-amplitude state before a transition occurs. The exponential dependence on the residence time indicates that the transitions are random and follow Poisson statistics as expected.

The activation energy at a particular detuning frequency is determined by recording a large number of transitions for multiple noise intensities ($I_N$). The average residence time at each noise intensity is extracted from the exponential fit to the corresponding histograms. Figure 4 plots the logarithm of the transition rate as a function of inverse noise intensity. The transition rate varies exponentially with inverse noise intensity, demonstrating that escape from the high-amplitude state is activated in nature. According to (7), the slope in Fig. 4 yields the activation barrier for escaping from the high-amplitude state at the particular detuning frequency.



**Fig. 4.** Logarithm of the transition rate from the high amplitude state to the low amplitude state as a function of inverse noise intensity at detuning frequency of 0.15 Hz. The slope of the linear fit yields the activation energy

Measurement of the activation barrier at different detuning frequencies is currently in progress. Such data will yield the critical exponent for noise activated switching between dynamical states in a driven, nonlinear system and allow experimental verification of the system-independent scaling of activation energy near the bifurcation point.

Apart from the scaling behavior of the activation barrier near the bifurcation point, novel phenomenon arises when the transition rates out of the two states are comparable. As we have discussed, the transition rate depends exponentially on the activation barrier (7). The ratio of the populations of the two dynamical states is given by:

$$w_1/w_2 \propto e^{(E_2 - E_1)/I_N} \tag{9}$$

**Fig. 5.** (**a**) Occupation of the high-amplitude state (*upright triangle*) and the low amplitude state (*inverted triangle*) as a function of the driving frequency. (**b**) The oscillation amplitude switches between two values as a function of time. (**c**) Histogram of the oscillation amplitude

where $E_1$ and $E_2$ are the activation barriers to escape out of states 1 and 2 respectively. Due to the exponential dependence of the population ratio on $E_2 - E_1$, the oscillator remains in one of the states at most driving frequencies. Only over a very narrow range of driving frequencies will the occupation of the two states be comparable [7,17]. This behavior resembles equilibrium systems with multiple phases such as vapor and liquid. In thermal equilibrium, these systems are usually in one of the phases and coexistence of the phases only occurs at the phase transition. Albeit our micromechanical oscillator is driven far from equilibrium, theoretical analysis predicts that a similar kinetic phase transition occurs at certain driving frequencies [7].

Figure 5a shows the occupation of the two states as a function of driving frequency. On the low frequency side of the hysteresis loop, the occupation of the high-amplitude state is negligible while the occupation of the low-amplitude state is almost unity. As the frequency increases, the activation energy for escaping from the high-amplitude state increases. On the high frequency side of the hysteresis loop, the occupation of the two states are reversed, with the probability of finding the oscillator in the upper state close to unity. While the oscillator is predominantly in one of the states at most driving frequencies, a small range of frequencies exists where the populations of the two states are comparable. The oscillator undergoes a kinetic phase transition at these driving frequencies [7,18,19]. Figure 5b plots the oscillation

amplitude as a function of time, illustrating the oscillator switching between two states. The relative occupation of the two states at this driving frequency is obtained by calculating the area under the two peaks in the histogram of the oscillation amplitude in Fig. 5c.

We are currently investigating a range of fluctuation phenomena near the kinetic phase transition. Micromechanical torsional oscillators provide well-characterized systems with tunable characteristics, allowing quantitative comparisons of experimental results with theoretical predictions of activated escape.

# References

1. H. A. Kramers, Physica (Utrecht) **7**, 284 (1940).
2. R. Graham and T. Tel, *Existence of a Potential for Dissipative Dynamical-Systems*, Physical Review Letters **52**, 9 (1984).
3. J. Lehmann, P. Reimann and P. Hanggi, *Surmounting oscillating barriers: Path-integral approach for weak noise*, Physical Review E **62**, 6282 (2000).
4. J. Lehmann, P. Reimann and P. Hanggi, *Surmounting oscillating barriers*, Physical Review Letters **84**, 1639 (2000).
5. R. S. Maier and D. L. Stein, *Noise-activated escape from a sloshing potential well*, Physical Review Letters **86**, 3942 (2001).
6. A. P. Dmitriev and M. I. Dyakaonov, *Activated and tunneling transitions between two forced-oscillation regimes of an anharmonic oscillator*, Sov. Phys. JETP **63**, 838 (1986).
7. M. I. Dykman and M. A. Krivoglaz, *Theory of fluctuational transitions between stable states of a nonlinear oscillator*, Sov. Phys. JETP **50**, 30 (1979).
8. M. I. Dykman and M. A. Krivoglaz, *Fluctuations in Non-Linear Systems near Bifurcations Corresponding to the Appearance of New Stable States*, Physica A **104**, 480 (1980).
9. M. I. Dykman, C. M. Maloney, V. N. Smelyanskiy and M. Silverstein, *Fluctuational phase-flip transitions in parametrically driven oscillators*, Physical Review E **57**, 5202 (1998).
10. L. J. Lapidus, D. Enzer and G. Gabrielse, *Stochastic phase switching of a parametrically driven electron in a Penning trap*, Physical Review Letters **83**, 899 (1999).
11. J. S. Aldridge and A. N. Cleland, *Noise-enabled precision measurements of a Duffing nanomechanical resonator*, Physical Review Letters **94**, 156403 (2004).
12. R. L. Badzey, G. Zolfagharkhani, A. Gaidarzhy and P. Mohanty, *Temperature dependence of a nanomechanical switch*, Applied Physics Letters **86**, 3587 (2005).
13. I. Siddiqi, R. Vijay, F. Pierre, C. M. Wilson, L. Frunzio, M. Metcalfe, C. Rigetti, R. J. Schoelkopf, M. H. Devoret, D. Vion and D. Esteve, *Direct observation of dynamical bifurcation between two driven oscillation states of a Josephson junction*, Physical Review Letters **94**, 027005 (2005).
14. M. I. Dykman, R. Mannella, P. V. E. McClintock, S. M. Soskin and N. G. Stocks, *Noise-Induced Narrowing of Peaks in the Power Spectra of Underdamped Nonlinear Oscillators*, Physical Review A **42**, 7041 (1990).

15. D. Ryvkine, M. I. Dykman and B. Golding, *Scaling and crossovers in activated escape near a bifurcation point*, Physical Review E **69**, 061102 (2004).
16. D. Landau and E. M. Lifshitz. *Mechanics* (Pergamon, London, 1976).
17. P. Hanggi and H. Thomas, *Stochastic-Processes – Time Evolution, Symmetries and Linear Response*, Physics Reports-Review Section of Physics Letters **88**, 207 (1982).
18. M. I. Dykman, D. G. Luchinsky, R. Mannella, P. V. E. McClintock, N. D. Stein and N. G. Stocks, *Supernarrow Spectral Peaks and High-Frequency Stochastic Resonance in Systems with Coexisting Periodic Attractors*, Physical Review E **49**, 1198 (1994).
19. M. I. Dykman, R. Mannella, P. V. E. McClintock and N. G. Stocks, *Fluctuation-Induced Transitions between Periodic Attractors – Observation of Supernarrow Spectral Peaks near a Kinetic Phase-Transition*, Physical Review Letters **65**, 48 (1990).

# Nonadiabaticity in Modulated Optical Traps

J.R. Kruse, D. Ryvkine, M.I. Dykman, and B. Golding

Department of Physics and Astronomy, Michigan State University, E. Lansing, MI 48824-2320 (USA)

**Abstract.** Experiments on noise-induced escape of a mesoscopic particle in a double-well potential are described. The potential is created by the interaction of two focused laser beams with a single sub-micrometer dielectric particle. By mapping the 3-dimensional trapping potential, the eigenfrequencies of the trapped particle are found. Over-barrier transitions are directly measured as a function of the rate and amplitude of a modulation that periodically tilts the potential. At low modulation rates and amplitudes the particle follows the potential adiabatically. As the system approaches its saddle-node bifurcation, different scaling regions emerge, each characterized by distinctive power-laws as predicted by recent theories. Of particular interest is the presence of a weakly non-adiabatic region with novel critical behavior.

## 1 Introduction

Optical trapping of sub-micrometer colloidal particles has proven to be a useful technique for non-invasive manipulation of physical and biological objects [1–3]. The method is most often applied to transparent dielectric objects in aqueous solution. The trapping forces exerted on the colloidal particles are of order 1–100 pN under typical conditions [4]. Since these forces are greater than the gravitational force on the particle in solution, the traps can be employed as optical "tweezers", capable of grabbing the particle and moving it in three dimensions. The particle can be visualized as being confined by a potential created by the interaction of the dielectric particle with the time-averaged radiation intensity of the optical beam. In water, significant viscous forces are brought into play when the colloidal particle is manipulated by translating the beam. This leads to an overdamped response controlled by a Stokes viscous damping rate $\Gamma$.

In determining the response of the colloidal particle to a time-dependent change of the potential, it is generally assumed that the particle remains in equilibrium with its environment, a statement of adiabaticity. However, when the drive frequency becomes sufficiently high, adiabaticity is violated since the particle is no longer able to follow the changing potential [5]. In this paper, we

explore the conditions under which adiabatic conditions exist in a periodically modulated double optical trap. Quite surprisingly, we have found that non-adiabatic effects become significant at very low modulation frequencies, of order 1 Hz, for a 0.6 μm colloidal glass sphere in water.

In the experiments reported here, a bistable optical potential is created by two independent optical traps with sub-μm separation [6,7]. The potential is tilted periodically at a constant frequency and amplitude by intensity modulating one, or both, optical beams. The single observable is the position of the particle as it oscillates about one of the two stable points of the potential or makes a transition over the single potential barrier separating the two wells. The optical potential is determined by applying a statistical analysis to experimental data. Over-barrier transition rates are also measured. As the modulation amplitude is increased, the transition rate increases in a way that is well-described by recent theories. Near the bifurcation point, the transition rates follow a universal form, with critical exponents that depend on modulation frequency, i.e., degree of adiabaticity. Agreement between experiment and theory [8,9] in the previously unexplored non-adiabatic region is quite good.

## 2 Static Optical Potential

A single sub-micrometer diameter silica sphere is suspended in water and trapped in the bistable potential formed by the particle's interaction with two 0.633 nm wavelength beams. Figure 1a shows a rendering of the two focused laser beams and the trapped sphere near the focal point in the sample cell. The beams propagate in the $+z$ direction with separation along the $x$-axis by 0.3 to 0.4 μm. The intent of this geometry is the creation of a single saddle-point that separates two stable trapping points. A one-dimensional schematic of the potential along the $x$-direction is shown in Fig. 1b. The particle undergoes Brownian motion near the bottom of one of the potential minima until a large fluctuation causes the particle to overcome the barrier, resulting in an interwell transition. By controlling the relative beam separation and the beam intensities, the potential barrier is tunable from zero to 10 $k_BT$.

The two HeNe lasers, with nominally 20 mW maximum power, are stabilized by Pockels cell electro-optic modulators. The two beams are combined at a beam splitter, and imaged into a sample cell by a $100 \times 1.4\,\mathrm{NA}$ PlanApo objective lens. As they enter the objective, the beams are mutually incoherent and circularly polarized. A portion of the beams is extracted before entering the objective so that beam steering fluctuations can be suppressed. The beams are separated and allowed to impinge on quadrant photodiodes. Feedback voltages are generated and used to control galvo positioners located before the beam splitter/combiner. This active feedback circuit has proven necessary to reduce relative motion of the two beams in the focal plane of the objective.

**Fig. 1.** Optical trap formed by two focused laser beams. Representation of the (**a**) two beam optical trap with 0.6 μm sphere and (**b**) potential energy $U(x)$ cross-section. The beams propagate in the $+z$ direction and are separated in the $x$ direction. Interaction of the beams and particle creates a potential with barrier $b$ and stable points $a$ and $c$

In general, system stability to better 10 nm for periods of an hour or more is needed to insure reliable transition rate measurements.

To measure the transition statistics (and the three dimensional potential energy) the particle is imaged onto a CCD camera whose output is transferred to a digital computer for analysis. The computer extracts the three coordinates of the center of mass $\mathbf{r}(t)$ of the sphere in less than 5 ms, the reciprocal frame rate of the present camera. The resulting time series is then analyzed to obtain statistical measures of the transition dynamics, such as instantaneous and average transition rates.

A quantitative determination of the three-dimensional potential is made by use of the Boltzmann expression

$$\rho(\mathbf{r}) = Z^{-1} \exp\left[-U(\mathbf{r})/k_B T\right], \tag{1}$$

where $\rho(\mathbf{r})$ is the probability density for finding the particle in a volume element $\delta\mathbf{r}$. The equation is inverted to find $U(\mathbf{r})$ where $U(\mathbf{r})$ is the position-dependent potential energy and $Z$ is a normalization constant [6].

The particle position is sampled for at least $10^6$ frames, a duration much longer than the mean interwell transition time. Slices through the potential in the three orthogonal directions yield the energy contours shown in Fig. 2.

The potential is parameterized in the vicinity of the extrema by assuming a quadratic coordinate dependence and performing least-squares fits to the data. Figure 3 shows the one-dimensional profiles along orthogonal $y, z,$ and $x$ directions at $r_1$, $r_2$, and $r_b$, respectively. Table 1 summarizes the results of the fits.

**Fig. 2.** Three dimensional potential energy slices for a two-beam trap. The contours represent 1 $k_BT$ energy intervals with local minima at the stable points $r_{1,2}$ given by black diamonds. *Upper left*: $x-z$ cross-section through the center of the trap. *Upper right*: $y-z$ cross-section. *Lower*: $x-y$ cross-section. To make a transition the particle passes through the vicinity of the saddle point $r_b$, requiring $x$ and $z$ displacements. Note that symmetry is broken about the focal plane owing to radiation pressure displacement of the particle along $z+$

**Fig. 3.** One dimensional profiles along $y, z$, and $x$ directions with quadratic fits. The plots show, from *left* to *right*, fitted profiles for the left well, right well, and saddle point. Results of quadratic fits (*solid lines*) are shown by the solid lines. The horizontal scales have arbitrary origins

## 3 Periodically-Modulated Optical Potential

Of principal interest here is the escape rate dependence on modulation amplitude when the amplitude is large enough that the distance between stable and unstable states becomes small. At frequencies small compared to $t_r^{-1}$, these states merge at an *adiabatic* critical amplitude $A_c^{\mathrm{ad}}$. For a portion of the period, the potential near the saddle point becomes flat. At the bifurcation point, the restoring force disappears; the particle is free to diffuse. Slow motion close to a bifurcation point is controlled by a soft mode. It is generically one-dimensional.

The motion of an overdamped one-dimensional Brownian particle is described by the Langevin equation $\dot{q} = K(q; A, t) + \xi(t)$, where $K(q; A, t) =$

**Table 1.** Double-well parameters from quadratic fits to one-dimensional profiles. The subscript 1(2) refers to the left(right) well, b refers to the barrier, $\omega$ is the oscillation frequency, the relaxation rate $t_r^{-1} = \omega^2/2\pi\Gamma$, $\Delta U$ is the potential energy relative to $r_1$, and $r_o$ is the position of the potential energy extremum taken from the data in Fig. 3. Note the large anisotropy in the potential, leading to slow relaxation in the z-direction

|        | $\omega(10^4 \text{ s}^{-1})$ | $t_r^{-1}$ (Hz) | $\Delta U(k_B T)$ | $r_0(\mu\text{m})$ |
|--------|------|------|------|------|
| $x_1$ | 14   | 125  | 0    | 0.16  |
| $x_2$ | 16   | 175  | 0.5  | 0.54  |
| $x_b$ | 5.9  | 22   | 3.8  | 0.36  |
| $y_1$ | 30   | 570  | 0    | 0.125 |
| $y_2$ | 29   | 540  | 0.6  | 0.125 |
| $y_b$ | 27   | 460  | 3.4  | 0.127 |
| $z_1$ | 2.3  | 3.4  | 0.1  | 0.95  |
| $z_2$ | 2.2  | 3.1  | 0.6  | 0.81  |
| $z_b$ | 2.0  | 2.5  | 3.3  | 1.26  |

$K(q; A, t + \tau_F)$ includes the potential which is periodically modulated with amplitude $A$ and period $\tau_F$. The adiabatic stable and unstable states $q_{a,b}^{\text{ad}}(t)$ are the solutions to $K(q_{a,b}^{\text{ad}}; A, t) = 0$. The interstate barrier is assumed to be much larger than the noise intensity $D$ ($k_B T$ in this experiment), and the mean interwell transition rate $W$ satisfies $W \ll \omega_F, t_r^{-1}$.

The motion of the particle is adiabatic when $\tau_F \gg t_r$, the interwell relaxation time. As $q_a^{\text{ad}}$ and $q_b^{\text{ad}}$ approach each other, the potential curvature decreases, changing the restoring force. Near the critical amplitude $\partial K/\partial q \to 0$, the instantaneous relaxation time $t_r = -(\partial K/\partial q)$ diverges, and the adiabatic approximation is violated. This leads to a crossover to a non-adiabatic asymptotic behavior [7–9].

The transition rate is affected exponentially strongly by the change in barrier height. Assuming $\Delta U \gg k_B T$ (or $W \ll \omega_F$), a transition is most probable in a small time interval about $t = n\tau_F$ when the barrier is at its lowest. Expanding $K(q; A, t)$ around $t = 0$ and $A = A_c^{\text{ad}}$ and keeping the lowest order terms,

$$K \approx \alpha q^2 + \beta \delta A^{\text{ad}} - \alpha\gamma^2 (\omega_F t)^2 . \tag{2}$$

where $\delta A^{\text{ad}} = A - A^{\text{ad}}$, $\alpha = \frac{1}{2}\frac{\partial^2 K}{\partial q^2}$, $\beta = \frac{\partial K}{\partial A}$, and $\gamma^2 = \frac{-1}{2\alpha\omega_F^2}\frac{\partial^2 K}{\partial t^2}$, are evaluated at $A = A_c^{\text{ad}}$, $t = 0$, and $q = q_a^{\text{ad}}(0)$.

The adiabatic relaxation time becomes

$$t_r^{\text{ad}} = \frac{1}{2}[(\alpha\gamma\omega_F t)^2 - \alpha\beta\delta A^{\text{ad}}]^{-1/2} \tag{3}$$

around the stable and unstable points $q_{a,b}^{\text{ad}} = \mp 1/(2\alpha t_r^{\text{ad}})$. For $t = 0$ and $\delta A^{ad} = 0$, $t_r^{\text{ad}}$ diverges. Using (2) the condition for adiabaticity is $\omega_F t_r^{\text{ad}} \ll 1$.

However, a second stronger condition is necessary for adiabaticity: the time-dependence of the relaxation time must satisfy $\left|\partial t_r^{\mathrm{ad}}/\partial t\right| \ll 1$. This condition is also expressed as $t_r^{\mathrm{ad}} \ll t_l$ for $t_l = (\alpha \gamma \omega_F)^{-1/2}$, or $\omega_F \ll \left|\beta \delta A^{\mathrm{ad}}\right|/\gamma$. The time scale $t_l$ limits the adiabatic approximation.

As the stable and unstable points approach each other their trajectories become distorted. Solving (2) for $K = 0$ as $\delta A^{\mathrm{ad}} \to 0$, it is found that the extrema merge along the line $q_{a,b}(t) = \gamma \omega_F t$. The critical amplitude becomes frequency dependent, $A_c^{\mathrm{sl}} = A_c^{\mathrm{ad}} + \gamma \omega_F/\beta$, where $A_c^{\mathrm{sl}}$ is the critical amplitude in the *nonadiabatic* slow driving regime. The condition $\omega_F t_r^{\mathrm{ad}} \ll 1$ still holds in this region.

As the system approaches criticality, asymptotic behavior emerges for the transition rate. Here, the time-averaged escape rate can be expressed as $\ln W = B + C(A_c - A)^\xi$. In the adiabatic regime, $A \to A_c^{\mathrm{ad}}$ and $\xi = 3/2$. The exponent $3/2$ has been derived previously in several contexts: Josephson junctions [10], magnetic materials [11], and other slowly varying processes [12].

As noted above, as $A_c^{\mathrm{ad}}$ is approached, the system must invariably fall out of equilibrium before $A_c^{\mathrm{ad}}$ is reached, whereupon new scaling appears. In the weakly nonadiabatic case, a crossover to the critical exponent $\xi = 2$ occurs. In this regime, the asymptotic $A_c^{\mathrm{sl}}$, given above, depends on the modulation frequency and is always greater than $A_c^{\mathrm{ad}}$.

## 4 Experimental Results

The adiabatic regime was established by studying transition probabilities vs. modulation amplitude as a function of modulation frequency. A frequency-independent region was found only for $f < 0.5\,\mathrm{Hz}$. This result is surprising as it implies that, for the trap parameters used here, the adiabatic regime appears only at extremely low frequencies and is consequently difficult to access. Therefore, at 20 Hz one expects a weakly nonadiabatic regime to exist over a substantial range of modulation amplitudes. Figure 4 shows the results of an experiment carried out for two different values of the static barrier height. A fitting procedure using the functional form $\ln W = B + C(A_c^{sl} - A)^\xi$ was used. The parameters $B, C, A_c^{\mathrm{sl}}$, and $\xi$ are treated as unknown parameters. At the critical point is approached, the system becomes phase-locked to the modulation field with transitions occurring every half-cycle. Therefore we can set $B = \omega_F/\pi$ and $\ln(W) = B$ for $A \geq A_c^{\mathrm{sl}}$. The results of the procedure, Fig. 4, show that the critical exponent $\xi$ is consistent with the value of 2 in the weakly nonadiabatic regime. The results clearly exclude the adiabatic exponent $3/2$, previously observed in this system at modulation frequencies below 1 Hz [7].

**Fig. 4.** Plot of the time-averaged escape rate $W$ vs. modulation amplitude (arbitrary units). The modulation frequency is 20 Hz. The static potential barrier $\Delta U/k_B T$ for the data represented by the *open circles* is 2.9; for the *solid squares*, 4.1. The *solid lines* represent a non-linear least squares fit to the data yielding the critical exponent $\xi$

## 5 Discussion

In view of the extremely low modulation frequency capable of inducing nonadiabatic dynamics in this bistable system, it is worthwhile to understand its origin. From the characteristics of the 3-dimensional potential shown in Fig. 2 and summarized in Table 1, we can calculate the crossovers in the amplitude-frequency plane based on the criteria outlined in the theoretical development. Figure 5 shows the plane with several boundaries inscribed. The boundary $ft_{r,x} = 1$ is the standard criterion for separating adiabatic from nonadiabatic regions. If this criterion were applicable, at 20 Hz one should expect adiabatic dynamics except in a narrow region close to $A_c^{\mathrm{ad}}$. This contrasts also with the observation at 5 Hz of a weakly nonadiabatic exponent $\xi = 2$. The weakly nonadiabatic criterion, $dt_{r,x}/dt \ll 1$ is more restrictive and, for a given modulation amplitude, predicts a much lower crossover frequency. As noted above, it is in this region that the system falls out of equilibrium with the driving field for only a portion of the modulation cycle.

It is interesting to speculate that there may be another contributor to low frequency nonadiabaticity. This may arise from the large anisotropy of the 3-dimensional potential, in particular, the "softness" in the $z$-direction. When

**Fig. 5.** Boundaries of adiabatic and nonadiabatic regimes in the $(f, A)$ plane. Black lines are the adiabatic critical amplitude $A_c^{\mathrm{ad}}$ (*solid*) and $A_c^{\mathrm{sl}}$ (*dashed*). Red lines are the weakly nonadiabatic crossover, $|\partial t_{r,x}/\partial t| = 0.1$ (*solid*) and $|\partial t_{r,x}/\partial t| = 1$ (*dashed*). The blue line is the strongly nonadiabatic crossover, $ft_{r,x} = 1$, and the cyan line is $ft_{r,x} = 1$, the crossover due to $t_{r,z}$. The modulation frequency f is plotted on a $\log_{10}$ scale

a transition occurs, the dominant motion between stable points is along the $x$-axis, but the particle must also move in the $z$-direction to pass through the saddle point, as seen in Fig. 2. Consequently, if the criterion $dt_{r,z}/dt \ll 1$ is enforced, the boundary is pushed to much lower frequencies. One would then expect to find evidence for nonadiabaticity setting in at frequencies near 1 Hz. The data are not sufficiently precise to establish experimental crossovers with great certainty. They do suggest, however, that it may be necessary to take into account the full 3-dimensional dynamics of the particle in the bistable trap to account for the nonequilibrium behavior with slow driving.

## 6 Conclusions

In this paper we have shown how an object, subjected to periodic forcing in the presence of noise, can easily fall out of equilibrium even with low frequency perturbations. The rate at which interstate transitions occur is greatly enhanced as the amplitude of the modulation is increased. This can be advantageous in situations in which control over thermally activated transition rates is desired without changing the system temperature. The use of periodic modulation of a barrier height for studying escape processes has a number

advantages over the use, for example, of linear ramps. One often uses *ad hoc* measures for choosing a ramping rate window that guarantees adiabaticity and rate-independence of escape probabilities.

Interest in optical trapping has been stimulated by advances in biological applications and in the development of large arrays of traps [13,14]. Dynamics is imposed on the ensemble of trapped objects by modulating the array potentials or by introducing flow in the trapping medium. Arrays of traps are generated by simultaneously rastering and modulating a single beam [15,16] or by the use of a spatial light modulator [13,14]. With rastering, the particles see a time-averaged potential, since the time the light interacts with the particle is much shorter than any relaxation time. With spatial light modulators, modulation rates are limited by the computational overhead of calculating a Fourier transform or addressing the modulator.

With the ability to control large numbers of particles by optical methods, new applications are emerging. Interacting dynamical systems underlie many existing and future industrial and military platforms where one needs to control interacting vehicles of different types moving in a coordinated way. Examples include mobile sensing arrays or arrays of different objects in random environments, such as miniature submarines in the ocean and insect robots on land or air. Such systems typically have their own dynamics, but interact with others via a communication link. In many of these applications it is necessary to control the dynamics of the entire collection while operating in the presence of random environments and stochastic communication. In the most realistic cases, the dynamical systems are subjected to noise. Noise largely complicates the dynamics and inter-platform communication, thereby influencing methods for control. Optically trapped particle arrays studied with the methods outlined in this investigation should prove to be a useful model system to test the development of control algorithms.

## Acknowledgments

## References

1. A. Ashkin, Phys. Rev. Lett. **24**, 156–159 (1970).
2. A. Ashkin, J.M. Dziedzic, J.E.. Bjorkholm, and S. Chu, Optics Letters **11**, 288–290 (1986).
3. R.M. Simmons, J.T. Finer, S. Chu, and J.A. Spudich, Biophys. J. **70**, 1813–1822 (1996).
4. L.P. Ghislain, N.A. Switz, and W.W. Webb, Rev. Sci. Inst. **65**, 2762–2768 (1994).
5. V.N. Smelyanskiy, M.I. Dykman, and B. Golding, Phys. Rev. Lett. **82**, 3193–3197 (1999).

6. L.I. McCann, M.I. Dykman, and B. Golding, Nature **402**, 785–787 (1999).
7. M.I. Dykman, B. Golding, J.R. Kruse, L.I. McCann, and D. Ryvkine in *Unsolved Problems of Noise and Fluctuations*, AIP Conf. Proc. Vol. **665**, edited by S.M. Bezrukov, pp. 428–434 (AIP Publishing, Melville, NY) 2003.
8. M.I. Dykman, B. Golding, and D. Ryvkine, Phys. Rev. Lett. **92** (2004).
9. D. Ryvkine, M.I. Dykman, and B. Golding, Phys. Rev. E **69**, 080602 (2004).
10. J. Kurkijärvi, Phys. Rev. B **6**, 832 (1972).
11. R.H. Victora, Phys. Rev. Lett. **63**, 457 (1989).
12. M.I. Dykman and M.A. Krivoglaz, Physica A **104**, 480 (1980).
13. P.C. Mogensen and J. Glückstad, J. Optics Commun. **175**, 75–81 (2000).
14. D.G. Grier, Naturae **424**, 810–816 (2003)
15. D.L.J. Vossen, A. van der Horst, M. Dogterom, and A. van Blaaderen, Rev. Sci. Inst. **75**, 2960–2970 (2004).
16. R. Nambiar, A. Gajraj, and J.C. Meiners, Biophys. J. **87**, 1972–1980 (2004).

# Signal Processing and Control in Nonlinear Nanomechanical Systems

R.L. Badzey, G. Zolfagharkhani, S.-B. Shim, A. Gaidarzhy and P. Mohanty

Department of Physics, Boston University, Boston, MA 02215, USA

## 1 Introduction

Bestriding the realms of classical and quantum mechanics, nanomechanical structures offer great promise for a huge variety of applications, from computer memory elements [1] and ultra-fast sensors to quantum computing. Intriguing as these possibilities are, there still remain many important hurdles to overcome before nanomechanical structures approach anything close to their full potential. With their high surface-to-volume ratios and sub-micron dimensions, nanomechanical structures are strongly affected by processing irregularities and susceptible to nonlinear effects. There are several ways of dealing with nonlinearity: exceptional fabrication process control in order to minimize the onset of nonlinear effects or taking advantage of the interesting and oftentimes counterintuitive consequences of nonlinearity.

A fundamental control mechanism popular for its counter-intuitive ability to amplify coherent behavior via the addition of noise, stochastic resonance (SR) has waxed and waned in popularity since its inception over 20 years ago. Originally, it was postulated [2] as an ad hoc explanation for the periodic onset of ice ages over the Earth's climate history. Given a nonlinear system with an energy threshold subject to a sub-threshold periodic modulation and white noise, there is a certain regime of added noise power in which the system oscillates between its states. This oscillation is synchronized with the modulation. Although this theory has since fallen out of favor in the climate-modeling community, the appeal of the concept encouraged its extension to a large variety of systems. These include (but are by no means limited to) bistable ring lasers [3], neurophysiological systems (mechanoreceptors in crayfish [4] and crickets [5]), SQUIDs [6], mechanical systems [7,8], electronic systems [9] such as amplifiers, and geophysical systems [10]. However, there have not been any studies demonstrating the effect of stochastic resonance in nanomechanical systems. Aside from being simply one more system in which the phenomenon has been demonstrated, nanoscale systems [11] are interesting because of their proximity to the realm of quantum mechanics.

The combination of stochastic resonance and quantum mechanics has been the subject of intense theoretical activities [12–15] for many years; nanomechanical systems present a fertile ground for the study of a broad variety of novel phenomena in quantum stochastic resonance. Additionally, the physical realization of such nonlinear nanomechanical strings offer the possibility of studying a whole class of phase transition phenomena, particularly those modeled by a Landau-Ginzburg quantum string [16, 17] in cosmology.

## 2 Nonlinearity in Nanomechanical Structures

The nanomechanical structure under consideration is a simple suspended beam, clamped at both ends and subjected to a transverse driving force. Following continuum mechanics, it is well known that such a structure has a resonance frequency for the fundamental flexural mode which is given by

$$f \sim \sqrt{\frac{E}{\rho}} \frac{t}{L^2} \tag{1}$$

where E is the Young's modulus, $\rho$ is the material density, t is the thickness (dimension parallel to the transverse forcing), and L is the length of the beam. Our doubly-clamped beams were fabricated from silicon, with dimensions of 7–8 $\mu$m $\times$ 300 nm $\times$ 200 nm; these yielded resonance frequencies around 25 MHz. Clamping losses and electrode mass loading reduced the resonance frequency from that which could be assumed from the simple continuum mechanics derivation. Fabrication of these structures is a well-established sequence of techniques, including PMMA spinning, e-beam exposure, development, metallization, and both dry and wet etching techniques. The details of fabrication are given in several locations [1, 18]. The magnetomotive technique has been in use for many decades, providing an efficient and low-noise method for exciting the resonant modes of a suspended structure. A sinusoidal current is pushed through a metallic electrode in the presence of a magnetic field. The resulting Lorentz force is then:

$$\overrightarrow{\mathbf{F}} \cos \omega t = \overrightarrow{\mathbf{I}} l \cos \omega t \times \overrightarrow{\mathbf{B}} \tag{2}$$

It is easy to see that the Lorentz force will be perpendicular to both the field and the current; a field perpendicular to the plane of the substrate will produce an in-plane transverse vibration. This in turn produces a motional EMF on the two clamping electrodes. This voltage is proportional to the velocity of the beam:

$$V_{avg}(t) = \xi l B \frac{dx(t)}{dt}, \tag{3}$$

where the factor $\xi$ is a proportionality constant given by the mode shape. For a fundamental mode, $\xi = 0.53$. In the harmonic approximation, this induced

voltage is directly proportional to the displacement of the beam. Under this force, the beam will behave as a damped, driven harmonic oscillator:

$$m\ddot{x} + \gamma\dot{x} + kx = F\cos\omega t. \tag{4}$$

The amplitude of motion describes a Lorentzian lineshape as a function of frequency, centered at the resonance frequency and with a quality factor $Q = \omega_0/2\gamma$. In the limit of small dissipation $Q = \omega_0/\Delta\omega$. There are two main roads into the nonlinear regime: inherent nonlinearities in the suspended bridge, or increasing the driving force until the beam response becomes nonlinear. Both methods result in the equation of motion described by the Duffing equation,

$$m\ddot{x} + \gamma\dot{x} + kx \pm k_3 x^3 = F\cos\omega t, \tag{5}$$

where the addition of the cubic restoring force term yields a dramatic change in the resonance behavior. It is a well-known phenomenon that a Duffing oscillator exhibits hysteresis and bistability near the linear resonance frequency. The evolution of an oscillator from linear response into nonlinear behavior is well-studied and covered by a number of excellent texts – a simple qualitative explanation will suffice. Simply put, as the beam is subjected to ever-increasing driving forces, it stretches in response, thereby increasing its length. The clamping points being a fixed distance apart, the situation is reached whereupon the length of the bridge exceeds that between the clamping points. Therefore, when the bridge passes through equilibrium, it feels a compressive force from the pads, giving rise to a quartic term in the beam potential and therefore a cubic term in the equation of motion. This is the well-known Euler instability, which occurs under the influence of any compressive or tensile force. For a doubly-clamped beam, there is a critical force which describes the boundary between the linear and nonlinear regimes [19]:

$$F_c \sim \sqrt{\frac{\gamma^3}{k_3\omega_0}}. \tag{6}$$

Depending on the elasticity of the material, it is possible to switch between the linear and nonlinear response regimes at will by simply tuning the driving force. Eventually, however, the elastic response gives way to a plastic deformation of the bridge and the nonlinear effects become more and more prevalent and permanent. When excited nonlinearly, the response of the bridge evolves from the simple Lorentzian into curved peak with a dramatic drop as seen in the following figure.

   Once the beam is driven into a nonlinear regime it demonstrates hysteresis as a function of frequency, as predicted by mechanics. As the frequency is increased, the amplitude of the response function follows the pseudo-Lorentzian lineshape until it reaches the sharp drop as seen in Fig. 1. In reality, the full analytic solution of the frequency-response function shows a rather dramatic bend to lower frequencies, creating a frequency domain in which the amplitude

**Fig. 1.** Beam response with increasing drive force, showing the transition from the linear to nonlinear regime. The labeled curves represent the linear response regime of the oscillator in response to different applied forces, from 0.5 pN (curve A) to 3 pN (curve D). Beyond this critical force, the beam response is nonlinear and bistable

function is multi-valued. Therefore sweeps in one direction in frequency follow the upper curve of this region, while opposite sweeps follow the lower curve. Eventually, each sweep reaches a frequency in which the system is unstable and therefore causes a transition to the other state. This is clearly illustrated in Fig. 2.



**Fig. 2.** Hysteresis caused by the nonlinear bistable response of a silicon nanomechanical bridge under strong driving

As stated previously, the Euler instability occurs under the influence of either a compressive or tensile strain on the beam. These two situations lead to qualitatively different but dynamically similar situations. In the Duffing equation, this is taken into account by the alternative sign on the cubic term, with compression yielding a positive sign and tension a negative sign. These two situations also present slightly different frequency response spectra as well.

The compressive case is revealed by a sharp drop on the right-hand side; the tensile spring has a drop on the left. For illustrative purposes, Fig. 1 shows a compressive case, while Fig. 2 shows what appears to be a tensile case. Regardless of whether there is compression or tension, the oscillator is still bistable and hysteretic. The fact that bistability arises in both nonlinear cases is important, because the effect the electrical measurement circuit can have on the apparent behavior of the oscillator is worth noting. It is understood that a mechanical harmonic oscillator system can be written as an equivalent LCR circuit [20] by the same token an electrical system can be constructed as a mechanical one. Therefore, when a mechanical system is connected to an electrical LCR measurement circuit, the signals measured are due to both constituent parts. The electrical resistance, capacitance, and inductance can be related to the mechanical spring constant and dissipation

$$R \sim \frac{1}{\gamma k}, \quad C \sim k \quad L \sim \frac{1}{k}. \tag{7}$$

These relationships also hold as the oscillator evolves into nonlinear response – the electric circuit need not be nonlinear for its signal to combine with that of the oscillator. Therefore, the relationship between the effective linear ($k$) and nonlinear ($k_3$) spring constants of the combined electromechanical system can be either compressive (same sign) or tensile (opposite sign). Regardless, the physical state of the mechanical structure itself is determined by the physical geometry of the system – in all of the cases described here, the compressive response is the physically realistic one even though the frequency response might suggest otherwise. This supposition was borne out by an experimental check – changing the cables (to ones with different electrical parameters) resulted in a flip of the nonlinear response from a left-hand drop to a right-hand drop.

   In addition, since both tensile and compressive cases result in bistability, the distinction from the standpoint of device applications is a largely academic one. From the derivation of the frequency response, assuming that the solution can be written as a harmonic function, the states are of the vibrational amplitude. In actuality, the most accurate statement is that they are different vibrational velocity states. This is borne out by the response seen by the magnetomotive technique, which is at its heart a velocity measurement. Only should the harmonic approximation prove valid can it be said that the two states are indeed amplitude states. For all of the illustrative abilities of the previous analyses, they all neglect the effect of temperature on the oscillatory characteristics of the bridge. At room temperature, slight fluctuations will have a negligible effect, and large variations in temperature are neither realistic nor easy to implement. Low temperatures, however, are easily attained through a variety of techniques. We inserted a nanomechanical silicon bridge into the cryostat of a $^3$He refrigerator and varied the temperature from 300 mK up to 5 K.

As is clear from the Fig. 3, both the spring constant and the dissipation were affected by changes in temperature. The first panel is a series of frequency scans taken at a variety of temperatures. Changes in temperature have two separate but related effects. First is a shift in the resonance frequency largely due to a reduction of the spring constant. The second is a gradual broadening of the peak and decrease in the amplitude at resonance, stemming from an increase in the dissipation factor. However, these two results are not completely independent. The frequency shift is also partly responsible for the change in measured Q, making a true separation of the effect of temperature on $\gamma$ and $k$ difficult. However, to a large degree, measuring Q gives a good approximation of $\gamma$, and measuring the shift in frequency sheds light on the change in $k$.

From 300 mK to 5 K, the frequency shifts by 9.5 kHz; taken as a fraction of the initial resonance frequency of approximately 22.0945 MHz, this shift is $4 \times 10^{-4}$. This means that the spring constant also changes by approximately four parts in $10^4$, as $\Delta f \sim \Delta k$. Interestingly, what is also clear from this temperature sweep is the effect on the dissipation, measured by tracking the changes in the Q of the oscillator. While the frequency shift is negligible, the change in dissipation most certainly is not. The dissipation (1/Q) at 300 mK has a value of $25 \times 10^{-5}$; heating to 5 K increases this value to $45 \times 10^{-5}$, an increase by almost a factor of two.

It is important to note that the preceding data comes from a nanomechanical oscillator in linear response. In nonlinear response, the notions of resonance frequency and Q become very difficult to formulate. However, the characteristics of the hysteresis region itself contain information about the behavior of $\gamma$, $k$, and $k_3$. These three parameters are even more closely linked than in the linear regime. As the concept of resonance frequency is somewhat misleading for the nonlinear case, the better measure is the central frequency of the nonlinear hysteresis region. By the same token, the concept of Q as the ratio between the resonance and the FWHM of the peak loses merit in the nonlinear case; a measurement of the change in the width of the hysteresis is a better indicator. Again, extracting the exact contribution due to each particular parameter is difficult when looking at the temperature alone. When the temperature is increased as shown in Fig. 4, there are two main results – a slight shift to lower frequencies and a reduction in the width of the hysteresis. From a low of 300 mK to a high of 4 K, the central frequency of the hysteresis region shifts by approximately 1.25 kHz, and the width decreases by approximately 500 Hz. It is important to note, however, that transitions within the hysteresis region are probabilistic and affected by the ambient temperature. Many sweeps through the hysteresis region are required in order to fully characterize the frequency shift and width change. In general, though, this result agrees with earlier experiments [21] that examined the change in the nonlinear characteristics with the addition of white noise.

**Fig. 3.** Temperature effects on the dissipation and frequency shift of a nanomechanical silicon bridge. (**a**). Variations in the resonance peak with temperature. As the temperature is increased the peak gets broader and also undergoes shift to lower frequencies. (**b**). Graph of the frequency shift as a function of temperature. The total shift is approximately 9.5 kHz, or $4 \times 10^{-4}$ of the initial resonance at 300 mK. (**c**). The dissipation (inverse Q) as a function of temperature. This quantity changes by almost a factor of 2 over the temperature sweep

## 3 Control by Stochastic Resonance

One of the more interesting nonlinear effects to come forth in the past few decades is the phenomenon of stochastic resonance. Compelling because of its ability to draw out coherence from a background of noise, stochastic resonance has suffered a bit in recent years because of a dearth of new and interesting experimental observations. Most have been received as simply one more example of stochastic resonance in a system, with little impact beyond adding to a pile of curiosities. However, the presence of nonlinear behavior in

**Fig. 4.** Alteration of the hysteresis region by an increase in the temperature from 300 mK to 4 K

a nanomechanical system [22] opens up an interesting realm of inquiry, a new avenue into the exploration of the phenomenon.

These nanomechanical systems, by virtue of their small sizes, high frequencies, and low temperatures venture quite close to the regime of quantum mechanics. Additionally, these nanomechanical bridges are closely related to a system which has become even closer to reality – a truly macroscopic, mechanical quantum harmonic oscillator [23, 24]. One of the very basic and necessary conditions for any oscillator to exhibit quantum behavior is that its modal energy must be greater than the thermal spreading of the energy states. That is, $hf > k_B T$, where $k_B$ is the Boltzmann constant and $h$ is Planck's constant. For an oscillator with a resonance frequency of 1 GHz, this condition is achieved when the temperature is below 48 mK. Both of these are experimentally realizable conditions, and several interesting experiments have been put forth which show the deviation of such oscillators from the classical norm. Therefore, exhibiting stochastic resonance in these MHz-range nanomechanical oscillators can lead to exploring the phenomenon in the closely related oscillators beyond the quantum limit. At the quantum level, prying coherent behavior out of a noisy environment could prove to be a powerful means of achieving quantum control.

To look more universally, these systems closely approximate quantum strings described by a Landau-Ginzburg equation.

$$\frac{\partial \Phi(x,t)}{\partial t} = m\Phi(x,t) - \Phi^3(x,t) + \kappa \frac{\partial^2 \Phi(x,t)}{\partial t^2} + A\cos(\Omega t) + \eta(x,t) \qquad (8)$$

Here, the field variable $\Phi$ replaces the position variable $x$, and $\eta$ is the noise term. Therefore, these and related nanomechanical structures provide a powerful tool for examining fundamental properties of a whole host of systems which are governed by similar equations and the resultant phase transition effects. The central requirement for the observation of stochastic resonance is a nonlinear two-state system subjected to a subthreshold modulation in the presence of tunable white noise. With the well-behaved bistable behavior seen

in these nanomechanical oscillators, they present a very clean system in which to study the effect. Also, as the temperature and the driving power affect the hysteresis, the beam characteristics can be adjusted so as to study an entire parameter space.

For this experiment, we used two silicon bridges, cooling them in the $^3$He cryostat and exciting in-plane vibrational modes with the magnetomotive technique. At 300 mK bridge 1 had a linear resonance frequency of approximately 23.57 MHz (Q = 3700), while bridge 2 had a linear resonance frequency of 20.835 MHz (Q = 1000). Both exhibited nonlinear bistable behavior under the influence of suitably strong forcing (4.0 dBm for bridge 1, 1.0 dBm for bridge 2). It was found during previous studies of bridge 1 that it was possible to excite the oscillator nonlinearly at a single frequency within the bistable region. With the addition of a square-wave modulation it was then possible to control the state of the oscillator, forcing it back and forth between its two states with the modulation frequency $\Omega$.

Under the influence of the driving force and modulation, then, the equation of motion for either of these bridges is a modified Duffing equation of the form:

$$m\ddot{x} + \gamma\dot{x} + kx \pm k_3 x^3 = F_d \cos\omega_d t + f_m(\Omega)$$

$$f_m(\Omega) = \frac{F_m \Theta(t)}{2}; \quad \Theta(t) = \begin{cases} -1 & (n-1)T < n < n-1/2)T \\ 1 & (n-1/2)T < t < nT \end{cases} \qquad (9)$$

This shows the equation with square wave forcing, although in principle any modulation term can be used. The response of the bridge was highly dependent on the amplitude of the modulation, with low powers leading to a loss of switch fidelity but not in the amplitude of the switches. This is consistent with expectations, as the voltage difference of the hysteresis is fixed by the driving amplitude and temperature – the modulation only serves to overcome (or fail to overcome) the existing barrier between the two states. Figure 5 illustrates both the circuit and sample layout and an example of the switching behavior as a function of modulation power. Once controllable switching was established, the modulation was reduced to sub-threshold, and a broadband electrical noise source was introduced into the circuit, and the noise power was swept from $-71$ dBm ($\sim 79$ pW) to $-41$ dBm ($\sim 79$ nW). At each noise power, a time-resolved scan was performed in order to establish whether or not the addition of the noise would result in the re-emergence of switching. Low noise powers resulted in no response from the bridge – it remained in the initially prepared state. As the noise was increased, however, first sporadic, and then more synchronized switch events were seen. Each switch event was synchronized with the modulation, and occupied exactly one-half period of the modulation frequency of 0.05 Hz. If a switch was skipped, no new switch would appear until the modulation passed through at least one-half period and again entered the positive-going part of the signal. The initial increase in switching was dramatic and quite rapid, establishing full switching within only a few fractions of a dBm increase in noise power. The signature of stochastic

**Fig. 5.** (**a**) Schematic of the electrical circuit used to excite a nanomechanical silicon oscillator into nonlinear response, modulate switching between the two resultant bistable states, and introduce electrical broadband noise. (**b**) Sample switching events, showing the effect of reducing modulation. The bottom panel shows the synch out from the modulation source

resonance is an increase in the signal-to-noise ratio (SNR) as a function of applied noise power. SNR is defined as:

$$SNR = \frac{S(\Omega)}{N(\Omega)} \tag{10}$$

here, $S(\Omega)$ is the height of the peak in the power spectral density at the modulation frequency, and $N(\Omega)$ is the background noise. With the SNR data accumulated, the following graph was compiled (see Fig. 6), showing a quite dramatic increase in the SNR with increasing noise power. An interesting feature of this data is the fact that the SNR does not exhibit the typical sharp peak as a function of the noise, but is instead rather broad, extending for several dBm before declining sharply down to near zero again. It should be noted that our system was distinct for several reasons. One, it was subjected to a square wave modulation as opposed to the canonical sine wave. The effect of a square wave on stochastic resonance has been the subject of considerable theoretical interest [25,26]. Secondly, the electronic noise source was of broad but still finite bandwidth $-15$ MHz, to be exact. Canonical theory quite explicitly states that the noise introduced is Gaussian and broadband.

Additionally, because of the very clean nature of the system switching, the barrier between the two states was quite strong. There is a considerable amount of rigidity to the system – it is difficult to induce switching between the two and the switching, once begun, is robust. By the same token, dynamic changes to the conditions of the oscillator were not felt instantaneously. It is possible that a longer time series at each noise power would reveal subtle differences in the switching behavior that would lead to changes in the SNR. Finally, the nature of the oscillator also lends itself unusual switching characteristics

**Fig. 6.** SNR as a function of applied noise power. (**a**) Three representative time series, showing the reemergence of switching as the noise power is increased. The applied noise increases from the uppermost panel downward. (**b**) Taking a power spectrum of each time series yields the SNR for each noise power  these were combined to create a graph of the SNR versus the applied noise power. The signature increase of stochastic resonance is clear, even though the overall shape does not follow the expected canonical model

once the noise has increased above the effective power for stochastic resonance. In the canonical formulation and most experiments in stochastic resonance, the decrease in the SNR at higher noise powers is due to a swamping of the switching signal by a large number of incoherent switches. However, this is not the case here – the oscillator simply stops switching and instead remains in the upper state. As this is a symmetric system, there is no preference for the eventual settling in the upper state – subsequent runs would leave the system in the lower state with just as much probability. And lastly, the effect of the fact that the noise source bandwidth did not intersect or overlap with the system resonance frequency should be explored in more detail, as it could contribute to the rigidity. This matter of system rigidity is certainly an open question, and would benefit from additional experimentation and theoretical consideration.

As temperature has already proven itself to be a powerful noise source [18] in the behavior of nanomechanical oscillators, a natural check was to see if stochastic resonance occurred with the source of noise being temperature, not external and electrical.

The second oscillator was disconnected from the noise source and brought into nonlinear response. The modulation required for switching was exceptionally strong (19 dBm), to the extent that the electrical signal appeared in the time series on top of the normally-static beam response voltage. Again reducing the modulation to below threshold (16.5 dBm), the temperature was increased from 300 mK to 4 K and the SNR was extracted from the time series of the beam response. It is clear from Fig. 7, there is another peak in the SNR as the temperature increases beyond 2 K. Although it is tempting

**Fig. 7.** SNR as a function of ambient temperature. (**a**) A sequence of time series for this asymmetric bridge shows that even when the coherence is regained, it is not as clean as for the first bridge. Also, the first panel, which shows the sharp spikes and immediate decays, demonstrates the asymmetric character of this bridge. (**b**) The SNR graph demonstrates a strong resonance. The data is simply line-connected – there is no fit implied, and is affected by the coarseness of the sweep

to apply the same analysis to the SNR response of the second bridge as the first, the results should be scrutinized carefully. It is obvious from the time series and the power spectra that this bridge had a non-negligible asymmetry. The presence of asymmetry has been broached in a number of theoretical investigations [27], and it can be used to understand this system, as well.

Because the temperature is changing, its effect on the dissipation and spring constant of the oscillator should not be overlooked. As shown in the first section, the dissipation is quite strongly changed as a function of the temperature. Therefore it is necessary to look at the equation of motion with new eyes. Typically, the noise considered in the formulation of the stochastic resonance problem is purely additive. However, when the dissipation is a function of temperature, the damping coefficient $\gamma$ can change the dynamics of the situation. In this case, it would be best to consider the noise as being partially additive (from the direct effect of heating brought on by the increase in temperature) and partially multiplicative (from the changes to dissipation). Over the temperature range where the stochastic resonance is seen ($2\,K$ to $3\,K$), the dissipation (based on the linear data) changes by approximately 10 per cent – certainly it is an effect worthy of further investigation. The effect on the spring constant, however, is not large.

In conclusion, we have introduced stochastic resonance into an entirely new class of systems, ones which by virtue of their small size and high frequencies touch upon an area where this and other nonlinear phenomena may have a deep and lasting impact. Nanoscale systems, and nanomechanical systems in particular, are much more than another example of stochastic resonance. Their highly controllable nature makes them perfect candidates for the study of a host of exciting phenomena, to an extent rarely seen before.

## Acknowledgements

## References

1. Badzey R, Zolfagharkhani G, Gaidarzhy A, Mohanty P (2004) Appl. Phys. Lett. 85:3587–3589.
2. Benzi R, Sutera A, Vulpiani A (1981) J. Phys. A 14:L453–L457.
3. McNamara B, Wiesenfeld K, Roy R (1988) Phys. Rev. Lett. 60:2626–2629.
4. Douglass J, Wilkens L, Pantazelou E, Moss F (1993) Nature 365:337–340.
5. Levin J, Miller J (1996) Nature 380:165–168.
6. Rouse R, Han S, Lukens J (1995) Appl. Phys. Lett. 66:108–110.
7. Spano M, Wun-Fogle M, Ditto W (1992) Phys. Rev. A 46:R5253–R5256.
8. Stambaugh C, Chan H-B (2004) cond-mat 0405791.
9. Fauve S, Heslot S (1983) Phys. Lett. 97A:5.
10. Alley R, Anandakrishnan S, Jung P (2001) Paleoceanography 16:190–198.
11. Lee I Y, Liu X, Kosko B, Zhou C (2003) NanoLetters 3:1683–1686.
12. Wellens T, Buchleitner A (2000) Phys. Rev. Lett. 84:5118–5121.
13. Goychuk I, Hanggi P (1999) Phys. Rev. E 59:5137–5141.
14. Grifoni M, Hanggi P (1996) Phys. Rev. Lett. 76:1611–1614.
15. Lofstedt R, Coppersmith S (1994) Phys. Rev. Lett. 72:1947–1950.
16. Benzi R, Sutera A, Vulpiani A (1985) Journ. Phys. A 18:2239–2245.
17. Stephens G J, Calzetta E, Hu B.-L, Ramsey S (1999) Phys. Rev. D 59:045009.
18. Badzey R, Zolfagharkhani G, Gaidarzhy A, Mohanty P (2005) Appl. Phys. Lett. 86:023106.
19. Feynman R, Leighton R, Sands M (1989) The Feynman Lectures on Physics Vol. 1 Ch. 38 Addison-Wesley, Redwood, CA.
20. Mohanty P, Harrington D, Ekinci K, Yang Y, Murphy M, Roukes M (2002) Phys. Rev. B 66:085416.
21. Aldridge J, Cleland A, (2005) Phys. Rev. Lett. 94:156403.
22. Badzey R, Mohanty P (2005) Nature 437:995–998.
23. Gaidarzhy A, Zolfagharkhani G, Badzey R, Mohanty P (2005) Phys. Rev. Lett. 94:030402.
24. Gaidarzhy A, Zolfagharkhani G, Badzey R, Mohanty P (2005) Appl. Phys. Lett. 86:254103.
25. Dykman M, McClintock PVE, Mannella R, Stocks N (1990) JETP 52:780–782.
26. Morillo M, Gomez-Ordonez J (1995) Phys. Rev. E 51:999–1003.
27. Inchiosa M, Bulsara A, Gammaitoni L (1997) Phys. Rev. E 55:4049–4056.

# Signal Modulation by Martensitic Control of Shape Memory Alloy Thin Film Actuator Architectures[1]

C.M. Craciunescu, I. Mihalca, and V. Budau[2]

**Abstract.** The paper reports on theoretical and experimental results related to the thermal control of microactuators based on shape memory alloy thin films. The behaviour of actuators that have one or more phase transforming films deposited on a non-transforming substrate is influenced by the thermal stress that grows in the bimorph or trimorph architecture on cooling from the deposition or annealing temperature. When a phase transition occurs or is induced in the film it leads to a corresponding change in the stress state in the film/substrate architecture and can be reflected accordingly in the actuation of a cantilever. The hysteretic characteristics of the actuation by shape memory alloy films can be controlled by appropriately selecting the chemical composition of the film, the substrate material, the film/substrate thickness ratio and in some cases the external stimuli. The deposition temperature was a factor considered for modulating the output signal of the bimorph and trimorph cantilevers, as well as the sequence(s) of deposition in case of multilayers and trimorphs. Bimorphs with bilayer and structurally graded films and trimorph architectures have been characterized based on known results for bimophs with single layer deposited on the same type of Si cantilever-type substrates. The results show how the martensitic transformation occurring in the films or in the layers or microlayers is affecting the response of the actuator to thermal stimuli. The models proposed could allow the selection of appropriate parameters in order to generate a specific type of actuation or a modulated sensorial response to thermal (for shape memory alloy) or thermal and magnetic (for ferromagnetic shape memory alloy) stimuli.

[2] C.M. Craciunescu is Associate Professor and Victor Budau is Professor in Materials Science Department at Politehnica University of Timisoara, Romania, craciunescucm@yahoo.com, budau@mec.utt.ro; I. Mihalca is Professor in Physics Department at Politehnica University of Timisoara, Romania, mihalca@etv.utt.ro

# 1 Introduction

Shape memory alloys have frequently been considered for sensing and actuating applications including for micro applications [1–3] due to their capacity to change the shape in response to thermal stimuli. The basic idea for their control is related to the martensitic transformation and involves two phases with different properties i.e austenite (A) and martensite (M) [e.g. 4]. Since the phase transition in shape memory alloys (usually thermally controlled) is reversible, it is possible to control the actuation by one-way or two way shape memory effect [e.g. 5]. When the shape memory alloy is heated (by conduction, convection, radiation, by external or direct electrical heating, etc.), the low temperature phase (martensite) transforms into the high temperature phase (austenite), while on cooling the reverse transformation occurs. The main disadvantage of the one way shape memory effect relies to the fact that an external force is always needed to deform the alloy on cooling, because the shape recovery occurs spontaneously only on heating. On the other hand, the two-way shape memory effect obtained as a result of special thermomechanical treatments (intrinsic effect) or when the shape memory alloy is elastically attached to a counterpart (extrinsic effect) allows a spontaneous shape change on heating and on cooling [e.g. 6]. Such a thermally controlled actuation shows a swifter actuation compared to bimetals, and occurs when the shape memory alloy undergoes the martensitic phase transition. The temperatures for the phase transition (usually between $-50$ and $+150°$C) are influenced by the alloy composition [4–6].

The austenite $\leftrightarrow$ martensite phase transition associated with the actuation in shape memory alloys is reversible and occurs gradually in a temperature range ($M_s$–$M_f$) on cooling, and in the ($A_s$–$A_f$) temperature range on heating ($s$ and $f$ subscripts denote the *start* and *finish* temperatures for the formation of the corresponding phase – martensite (M) and austenite (A)). During the phase transition, an (austenite + martensite) mixture exists and the mechanical properties in the corresponding temperature range can be approximated using the rule of mixtures. The martensitic phase transition occurs not only controlled by temperature, but also by the combined effect of stress and temperature. The "stress-induced martensite" appears at temperatures higher than the ones for which the martensite would occur if only thermally controlled and may be unstable if the energetic conditions are not satisfied by the complex of thermomecanical factors.

Architectures based on shape memory alloys are excellent solution to overpass shape memory alloys disadvantages and are recently under intense scrutiny for applications in micro-electro-mechanical (mems) systems [e.g. 7]. In fact, they combine bimetal and shape memory alloy advantages leading to interesting actuators, with the output potentially controlled by the sequence of phase transformation in the transforming films or microlayers. Figure 1 synthesizes some of the observation concerning shape memory actuation, compared to bimetals.

| Actuator | | | Cantilever-type model | Actuation (δ) | Observations |
|---|---|---|---|---|---|
| Bimetal | | | Metal A / Metal B | $\delta$ vs $T$, linear | Linear dependence |
| Shape memory alloy based | | One-way | Shape memory alloy $+$ F | $\delta$ vs $T$ | Steeper actuation Need for an external force Shape recovery on heating only |
| | Architectures | Two-way | Shape memory alloy | $\delta$ vs $T$ | Spontaneous shape recovery on heating and cooling |
| | | Bimorph | Shape memory film / Substrate | $\delta$ vs $T$ | Controlled by thermoelastic stresses |
| | | Trimorph | Shape memory film / Substrate / Shape memory film | ? | Influenced by the nature of the films and the deposition sequence |

**Fig. 1.** Actuation principles based on thermally controlled shape memory alloys, compared to bimetals

By comparison to bulk shape memory alloys, architectures are becoming extremely important for microactuation because the shape change can be better controlled and the two-way memory effect can be obtained in the composite system. Thin film processing techniques can be used to fabricate microactuators based on shape memory alloys composites. With the film deposited on the substrate, a bimorph (when the film is deposited only on one side of the substrate) [8] or trimorph [9] (when the film is deposited on booth sides of the substrate) composite architectures can be obtained. Bimorphs have the actuation controlled by the phase transition the shape memory alloy film, while for trimorphs, the actuation could be controlled by the combination of the phase transition in each of the composing films.

This presentation aims to explore at theoretical and experimental levels, the ways the phase transformations in the films and the thermoelastic stresses developing in bimorph [10] and trimorph composite architectures can be used to modulate the deflection of shape memory alloys composite actuators under thermal control. In fact, starting from the analysis of the actuation in shape memory alloy thin film composites – based on previous experimental results – it will be shown that the deposition in trimorph architectures could allow the design of particular actuations. This can be done by appropriately selecting the phase transformations (and the corresponding transformation temperatures or expansion coefficients) in the films, deposited on booth sides

of the architecture. A model is designed in order to explain how the actuation
signal can be modulated in trimorph architectures, taking into account the
thermoelastic stresses and the phase transition in the films. The experiments
designed for two cases with NiTi shape memory alloy films – in bimorph and
trimorph architectures – demonstrate the influence of the deposition sequences
on the resulting actuation.

## 2 Actuation Principles in Shape Memory Alloys Thin Film Composites

The main particularity in composites architectures based on thin shape mem-
ory films deposited on substrates is related to the need to have at least some
degree of crystallinity in the film, in order to observe an actuation. When de-
posited on the substrate, the film and the substrate behave like a bimetal un-
less a phase transition occurs in the film (providing that a good film-substrate
adherence can be obtained). On the other hand, the film crystallizes either
as a result of annealing (usually when the film is deposited on unheated sub-
strates) or during the deposition (if the substrate is heated).

### 2.1 Typical Behavior in Bimorph Architectures

When the film $(f)$ has a higher expansion coefficient $(\alpha)$ than the substrate
$(s)$ $(\alpha_f > \alpha_s)$, both in the martensitic and austenitic state, a tensile stress
grows in the film on cooling from the deposition $(T_D)$, or crystallization $(T_C)$
temperature – slope (1) in Fig. 2.



**Fig. 2.** Schematic model for the stress vs. temperature relationship in case of a
shape memory alloy film-based bimorph microactuator (heating 1–3, cooling 4–6).
The $\sigma$ subscripts denote the fact that the martensitic transition occurs under stress

On reaching the $M_{s\sigma}$ temperature, at which the stressed-induced martensitic transformation starts in the film ($\sigma$ denotes the stress contribution), the stress is gradually decreasing, proportional with the fraction of martensite formed – slope (2). The ideal case would lead to a total cancellation of stress. Such a case has not been observed and one of the suppositions relates to the fact that the film that is already martensitic when the $M_{f\sigma}$ temperature is reached still has a fraction in contact with the substrate. This fraction cannot transform because of the constraints coming from the substrate [7]. On further cooling, the martensite reaches its plastic deformation limit and a new slope (3) proportional to the tensile stress develops this time due to the mismatch between the expansion coefficients of the martensitic film ($\alpha_{fM}$) and substrate ($\alpha_{fM} > \alpha_s$).

During the heating from martensitic state a similar sequence of slopes develops but in reverse order, i.e.: a stress decrease when the martensitic film is heated (4), followed by a steep increase in the stress as the martensite $\rightarrow$ austentite transformation develops (5), and finally, a decrease of the thermoelastic stress until it is fully cancelled at the deposition or crystallisation temperature (6).

## 2.2 Optimisation Potential of Shape Memory Alloys Thin Film Composites

The particularities of the shape memory alloy film based architectures and of the phase transformation allow additional ways to control the behaviour of such composites, derived from the possibilities to control the composition of the film, the deposition or annealing temperature, the development of a graded phase transition over the thickness of the film, the disposition of the films with respect to the substrate and additional layers. All these features can be used to modulate the actuation vs. temperature behaviour, in various combinations and architectures. Some of these possibilities are described below.

### Control of the Stress at the Onset of the Phase Transition in the Film

The films deposited at room temperature are usually not crystalline in the as-deposited state and a crystallisation annealing around $600°C$ is required in order to observe the shape memory effect [11]. Such crystalline films can reveal the phase transition in the film during the film-substrate temperature change. For example, Figs. 3a and b describe the experiment used to determine the change in the deflection of the free-end of a $100\,\mu m$ Si cantilever with (100) orientation, with a NiTi film deposited by sputtering.

The film/substrate curvature reflects the stress in the composite architecture that can be calculated using the Stoney equation [12] for $d_f \ll d_s$:

$$\sigma = \frac{E_s}{1 - \nu_s} \cdot \frac{d_s^2}{6 \cdot R \cdot d_f} \tag{1}$$

**Fig. 3.** Example of the influence of the stress at the onset of the phase transformation in the film on bimorph architectures behavior (cooling stage): (**a**) cantilever's deflection ($\delta$); (**b**) typical deflection vs. temperature dependence, emphasizing the phase transition in the film; (**c**) stress temperature dependence calculated using the Stoney formula, for NiTi films deposited on Si substrates at 385°C ($\Delta$) and 250°C (o) respectively. Films deposited at higher temperatures transform under higher stress

where $d_s$ and $d_f$ are the substrate ($s$), and film ($f$) thickness, $E_s$ şi $\nu_s$ the Young modulus and Poisson ratio of the substrate, and R – the bimorph curvature.

Figure 3c shows the stress change calculated based on the Stoney equation, recorded on cooling two bimorphs with NiTi films deposited at 250 and 385°C respectively. It can be observed that the films deposited at higher temperatures reach higher stress as at the onset of the phase transformation, compared to those deposited at lower temperatures. The main reason for such a difference resides in the thermoelastic stresses that are built in the film on cooling from the deposition temperature as a result of the mismatch between the expansion coefficient of the film and the substrate. As long as there is a difference between the thermal expansion coefficients of the film and the substrate thermoelastic stresses occur by changing the temperature of the bimorph, and can be expressed by the following relationship:

$$\sigma_f = E_f \cdot \Delta\alpha \cdot \Delta T \tag{2}$$

where $\Delta\alpha$ is the difference between the thermal expansion coefficients of the film and the substrate, $E_f$ – the elastic modulus of the film, and $\Delta T$ – the temperature factor reflects the difference between $T_C$ or $T_D$ and the particular temperature $T_X$ at which the stress is calculated. This relationship is valid only in the temperature range corresponding to the austenitic phase (from $T_D$ or $T_C$ until the maximum stress before the phase transition is achieved).

**Fig. 4.** Schematic representation of the influence of the substrate's expansion coefficient (**a**) and the transformation temperature of the shape memory alloy (**b**) on the stress at the onset of the transformation in the film. Distinct cases have been considered for comparison in each figure: same film composition and different substrates leading to different $\Delta\alpha_i$ – Fig. 4a; and same substrate but different transformation temperatures of the films, leading to different temperature factors $\Delta T_i = T_D(T_C)$-$M_{Si}$ – Fig. 4b. In each case, the thermoelastic stress before or at the onset of the transformation is affected. See text for details

Figure 4 describes schematically how the amount of the stress relief (and the corresponding actuation) can be controlled by appropriately selecting the deposition or crystallization temperature (see also Fig. 2), the substrate's material (considering the contribution of film and substrate expansion coefficients ($\Delta\alpha$) – Fig. 4a) and the shape memory alloy transformation temperature, respectively. The later case is schematically illustrated in Fig. 4b ($\Delta T_i$ reflects the difference between $T_D$ or $T_C$ and the temperature corresponding to the onset of the transformation in the film – $M_{s\sigma}$) and the observation is in line with the results of Winzeck et al. [13] who studied the phase transition in Ni-TiPd films with different compositions deposited on Mo substrates. The stress relief during the phase transition is likely to occur in a range delimited by the thermoelastic stresses of the austenite and martensite, respectively (see also Fig. 3b).

According to Fig. 4, the stress at the onset of the phase transformation can be influenced by appropriately selecting the substrate (and its expansion coefficient $\Delta\alpha_s$), the temperature for deposition ($T_D$) or crystallization ($T_C$), or the composition of the film. It is has also been previously described how the stress relief is increasing (and so is the actuation capacity) when the onset of the transformation occurs at higher temperatures. To a certain extent, the slope of the phase transition can also been influenced by the occurrence of different types of martensitic or premartensitic transformations. For example, in TiNiCu shape memory alloys, depending on the actual composition transformation sequences may occur and since each phase has its own elastic characteristics more than one slope can be obtained [14].

## Graded Transition in the Film

In composite architectures with film deposited on the substrate different behaviour of the film layer that is on top is expected with respect to the layer that is in contact with the non-transforming substrate. If the film is homogenous, i.e. the transformation temperature is identical over the thickness of the film, it follows that as the temperature is changed during the thermal cycling the film has the possibility to transform over its entire thickness (Fig. 4-left), and the transformation may progress in plane and out of plane based on the best energetic conditions.

However, if the film is graded in way that makes it composed out of layers ($\delta$) with different transformation temperatures the transformation will develop in plane and over the thickness of one layer at a time, rather that over the entire thickness of the film, since the neighbouring layer is either already transformed or its transformation temperature has not been reached (Fig. 4-right). In this case, the transition develops gradually, from the top of the film to its part attached to the substrate, leading to a complete and ordered transition over the thickness of the film and correspondingly to a higher actuation, compared to the prior case where the transition occurs randomly over the thickness of the film (see also [15]).

## Films or Layers Deposition Sequence

The behaviour of shape memory bimorphs with layered films and trimorph architectures can also be influenced by the sequence of deposition of the layers or the films, respectively. The main reason for such influences is related to the fact that a transforming layer sandwiched between two non-transforming parts of the architecture may not be able to relive the stress since it finds itself stacked. Such a case is possible to occur in graded films also (Fig. 6), if the layers are deposited in an opposite sequence compared to the one previously described in Fig. 4b., i.e the highest transformation temperature in the $\delta$ microlayers.

In Fig. 6, as the temperature is lowered, the first microlayer to transform is the one in the vicinity of the substrate, because its Ms temperature is the highest of all microlayers. However the effect of the transformation is reduced by the fact that the rest of the microlayers composing the film have lower Ms temperatures an have not transformed. By contrast to the previous case where the transformation of the microlayers started from the top layer, the actuation in this case is expected to be reduced.

Trimorph architectures based on shape memory alloys can be generated by simultaneous or successive deposition of the films on the sides of the substrates. In order to generate an actuation in the trimorphs with the films deposited simultaneously, different compositions of the films are needed, otherwise the "stress-relief" effect is compensated and no actuation can be observed. If the films are deposited successively on each side of the substrate,

**Fig. 5.** Comparison between the phase transformation in homogenous (*left*) and in a graded film with the transformation temperature $M_s$ increasing from the substrate interface to the top free surface (*right*). See text for details



**Fig. 6.** The influence of the deposition sequence in graded films. The $M_s$ is decreasing from the substrate interface to the top free surface, contrasting with the case described in Fig. 5. See text for details

the temperature of the substrate can be adjusted for each film so that the transitions occur distinctively (see also Fig. 3, where the temperature of the stress relief for each film is slightly different, in spite of the same composition).

Figures 7a and b show the stress in the bimorph and trimorph cantilevers deposited at two different temperatures ($T_1$ and $T_2$ respectively) on heated substrates. The Stoney equation can be used to calculate the stress based on the displacement of the cantilever. This is possible because the cantilever bends as a result of the interaction between one film and the substrate. By comparison, this equation cannot be used to calculate directly the stress of the trimorphs. If the expansion coefficients of the films deposited in trimorph architectures are equal, there will be no bending data to be used and the Stoney equation does not apply (it is only valid for bimorphs). In this case, the stress in the films results from adding the stresses that would have occurred in the corresponding bimorphs.

The following different regions can be considered for the films deposited at two different temperatures $T_1$ and $T_2$:

**Fig. 7.** Model for stress-temperature dependence for trimorph shape memory alloy architectures, compared to bimorph ones. In NiTi/Si bimorphs, the film is under tensile stress when cooled from the deposition temperature and leads to the bending of the cantilever toward the film. In trimorph architectures the bending of the two bimorphs in opposite directions is balanced by the initiation of the phase transition in each film ($a_2 \rightarrow t_2$ and $b_2 \rightarrow t_4$). (For clarity, the trimorph model has been shifted on the stress scale)

-$a_1$, $b_1$ – growth of thermoelastic stresses in the corresponding bimorphs due to the difference in the thermal expansion between the NiTi films in austenitic state and the Si substrate; The corresponding slope $\alpha$ and $\beta$ are influenced by the thermal expansion of the films.

-$a_2$, $b_2$ – stress relief in the bimorphs as result of the martensitic transformation;

-$a_3$, $b_3$ – growth of thermoelastic stresses in the bimorphs due to the difference in the thermal expansion between the NiTi film in martensitic state and the substrate;

-$t_1$ – the deflection of the free-end of the cantilever at a given temperature depends on the difference between the $a_1$ and $b_1$ stresses. The ($a_1$) stress is usually higher than ($b_1$) leading to a slight slope oriented toward bimorph I;

-$t_2$ – the onset of the transformation in bimorph I ($a_1$) is associated with a change in the curvature. The trimorph tends to adopt a curvature similar to the one of bimorph II

-$t_3$ – the transition in bimorph I is complete and its stress tends to increase ($a_3$), leading to an equilibrium with the stress generated in bimorph II ($b_1$). The trimorph actuation shows an intermediary plateau.

-$t_4$ – the onset of the transformation in bimorph II ($b_2$) leads to a similar result as the one observed during ($t_1$) stage, but in opposite direction.

-$t_5$ – the transitions in both films is complete and the stress in the corresponding bimorphs is increasing ($a_3$, $b_3$). The trimorph tends to bend toward the higher stressed film.

Different actuation curves can be designed using the model described before. The actuation can be modulated through the composition of the films, the transformation temperatures or the width and relative position of the films hysteresis. The order of deposition is basically influencing the state of stress in the films, however further research needs to be carried out in order to fully understand the influence of the deposition sequence.

## 3 Experimental Details

In order to relate the theoretical observations regarding the possibilities to control the actuation of shape memory alloy bimorph and trimorph architectures a dc magnetron system has been used to sputter-deposit films out of a $Ni_{50}Ti_{50}$ and Ni targets respectively on one or both sides of (100) Si cantilever-type substrates (Fig. 8) heated with a plate attached to the substrate.

The cantilevers were manufactured using the lithography technique and the surface of the cantilevers was oxidized at $1100°C$ in air. The depositions were made using the following parameters: 100 W power, $10^{-6}$ preliminary vacuum, 10 mTorr Ar pressure. The films were deposited on one side of the substrate at a time. Four sets of actuators have been made, with specific features described in Table 1. Bimorph-type cantilevers ($B$) have the film deposited on only one side of the cantilever, while trimorph-type cantilevers ($T$) have films deposited on both sides of the cantilever.

**Table 1.** Parameters used for the deposition of bimorph and trimorph actuators

| Architecture | Type | Target/ Deposition Temperature / Film Thickness | |
|---|---|---|---|
| | | First deposition /(side) | Second deposition /(side) |
| Bimorph (with one transforming film) | $B_1$ | Ni-Ti /385 °C/2 μm /(a)- | |
| | $B_2$ | Ni-Ti /250°C /2 μm /(a)- | |
| Bimorph deposited on heating (H) and cooling (C) | $B_H$ | Ni-Ti /250°C $\rightarrow$ 385°C /2 μm /(a) | |
| | $B_C$ | Ni-Ti /385°C $\rightarrow$ 250°C /2 μm /(a) | |
| Trimorph (with one transforming films) | $T_{n1}$ | Ni/385°C/2 μm /(b) | Ni-Ti /385°C/2 μm /(a) |
| | $T_{n2}$ | Ni/385°C/1 μm /(b) | Ni-Ti /385°C/2 μm /(a) |
| Trimorph (with two transforming films) | $T_{t1}$ | Ni-Ti /385°C /1 μm /(a) | Ni-Ti /250°C /1 μm /(b) |
| | $T_{t2}$ | Ni-Ti /250°C /1 μm /(b) | Ni-Ti /385°C /1 μm /(a) |

Bimorphs deposited at different temperatures ($B_1$ and $B_2$) have been used as reference for the analysis of the influence of the stress at the onset of the transformation and the results have been referred to in Fig. 3. Bimorphs with graded films deposited on heating ($B_H$) and on cooling ($B_C$) have been designed to analyze the influence of the sequence of deposition in bimorph architectures.

The influence of a non-transforming film in trimorph architectures has been investigated. It has been observed – as expected – that the slope of the stress in austenite and martensite can be influenced. Trimorphs with two transforming films allow an expanded control on the actuation profile. The main difference between ($T_{t1}$) and ($T_{t2}$) trimorph cantilevers is related to the sequence of the deposition temperature of the corresponding films. For example, the manufacturing process for $T_{t1}$ trimorphs was initiated by heating the cantilever at 385°C (the limit of the deposition system), followed by deposition of the NiTi film. Before the deposition of the second film, the deposition chamber was opened and the cantilever was turned with the un-deposited side toward the sputtering gun and reattached to the heating plate. The trimorph $T_{t1}$ architecture resulted after heating at 250°C and the deposition of the second film.

The composition of the films was determined by wavelength dispersive spectroscopy (WDS) and showed, as expected [11], the depletion of Ti, thus making the Ni-Ti films Ni-rich. X-ray spectra of the films deposited on the substrates were recorded at room temperature using a Rigaku X-ray diffractometer and CuK$\alpha$ radiation.

The actuation – reflected by the bending of the bimorph and trimorph cantilevers – was determined as a function of temperature, using the clamped free-reed vibration method], which also allowed the measurement of the damping and the modulus defect. A heterodyne circuit was used to excite the fundamental mechanical resonance of the bimorph and trimorph cantilevers. The bending of the cantilever was quantified by measuring the variation of the carrier frequency of the vibrating-reed apparatus. The measurement method is described in Fig. 3. First the cantilever was gradually forced to bend toward the electrode. This caused a variation of the capacitance between the cantilever and the electrode and therefore a change in the carrier frequency. The bending was measured and the corresponding carrier frequency was recorded. Then the cantilever was thermally cycled in the temperature range corresponding to the transformation of the film(s). The heating/cooling rate was 1°C/min. The carrier frequency data recorded during the thermal processing was converted into the displacement of the free end of the cantilever using the relationship previously established.

Figure 8 synthesizes the main aspects related to the experiments made in order to confirm the observations concerning the actuation control in bimorph and trimorph architectures.

Bimorphs deposited on heating ($B_H$) show an exceptionally large actuation considered to be the positive consequence of a complete and ordered

**Fig. 8.** Influence of bimorph (*left*) [15] and trimorph (*right*) architectures on the displacement vs. temperature signal for NiTi films deposited on one or both sides of Si cantilever-type substrates (*top*). $B_H$, $B_C$ bimorphs and $T_{t1}$ and $T_{t2}$ trimorphs respectively reflect the sequence of deposition. Se also Table 1 for details

transition that occurs gradually as the temperature is changed. By contrast, the actuation of bimorphs deposited on cooling ($B_C$) is limited and over a larger temperature range. The fact that a typical slope starts to shape at lower temperature seems to confirm the suggestions that only the last transforming layers effectively contribute to the actuation (see also the 2.2.2. and 2.2.3 sections).

Trimorph architectures with shape memory alloy films deposited on booth sides of the substrate show a behaviour that is dependent on the sequence of deposition. Both films have the same composition but the slope for the ($T_{t1}$) architecture (only one transformation type slope) significantly differs from the one of the ($T_{t2}$) architecture (one large slope and a smaller one in the opposite direction), also suggesting that the actuation signal can be modulated by appropriately selecting the composition of the films and the deposition sequence. (see also section 2.2.3).

The results show additional ways to control the actuation and can be used to further expand the limits of development in microrobotic applications based on shape memory alloy films [16].

# 4 Conclusions

Shape memory alloys have an outstanding potential to be used in microactuation systems due to the large actuation and the possibilities to control the actuation thermally or magnetically. The composite architectures based on shape memory alloy films are an excellent solution that can be integrated in

advanced fabrication techniques. The actuation is induced by the phase transformation in the shape memory alloy film and the thermoelastic stress that develops in bimorph and trimorph composites on cooling from the deposition or annealing temperature.

For a given composition of the film, the stress at the onset of the transformation can be adjusted by appropriately selecting the substrate and the deposition or annealing temperature respectively. Alternatively, the composition of the films and the corresponding phase transition temperatures can also be used to increase or decrease the stress at the onset of the phase transition.

Advanced variants with graded shape memory alloy films obtained by continuously changing the deposition temperature lead to increased actuation compared to regular films as long as the deposition sequence provides a transformation temperature gradient that allows an ordered transformation of the film from the free surface to the non-transforming substrate. Such a case has been observed for the depositions on heating. By contrast, very small actuation has been experimentally obtained for the depositions on cooling where the phase transformation in the film is expected to start in the microlayer sandwiched between the substrate and the not-yet-transformed rest of the film.

Trimorph architectures are capable to expand even more the ways to modulate the actuation signal by combining shape memory film properties on both sides of the substrate. In addition to the typical bimorph-type actuation it has been described the way the effect of the films can be combined for films with the same composition but deposited at different temperatures in the composite architecture.

In summary, the theoretical and experimental results lead to the conclusion that a significant potential to improve the actuation of shape memory alloy architectures exists and various new possibilities arising from the particularities of shape memory alloys can be used to further improve the actuation capacity.

# References

1. Fu Y, Du H, Huang W, Zhang S, Hu M (2004) TiNi-based thin films in MEMS applications: a review. Sensors and Actuators A 112: 395–408.
2. Ishida A, Martynov V (2002) Sputter-deposited shape-memory alloy thin, films: properties and applications, Materials Research Society Bulletin 27: 111–114.
3. Krulevitch P, Lee AP, Ramsey PB, Trevino JC, Hamilton J, Northrup MA (1996) Thin Film Shape Memory Alloy Microactuators. Journal of Microelectromechanical System 5(4):270–281.
4. Otsuka K, Kakeshita T (eds) (2002) Science and Technology of Shape Memory Alloys: New Developments. Materials Research Society Bulletin 27(2):91–100.
5. Eucken S (ed) (1992) Progress in Shape Memory Alloys. DGM Informationsgesellschaft mbH, Oberursel.

6. Otsuka K, Wayman CM (1998) Shape Memory Materials. Cambridge University Press.
7. Gill JJ, Chang DT, Momoda LA, Carman GP (2001) Manufacturing issues of thin film NiTi microwrapper. Sensors and Actuators A 93:148–156.
8. Roytburd AL, Kim TS, Su Q, Slutsker J, Wutig M (1998) Martensitic transformation in constrained films. Acta Materialia 46(14): 5095–5107.
9. Craciunescu CM, Mihalca I, Budau V (2003) Trimorph Actuation based on Shape Memory Alloys. Journal of Optoelectronics and Advanced Materials 7(2): 315–321.
10. Craciunescu CM, Li J, Wuttig M (2003) Thermoelastic Stress-Induced Thin Film Martensites. Scripta Materiallia 48(1): 65–70.
11. Miyazaki S, Ishida A (1999) Martensitic transformation and shape memory behavior in sputter-deposited TiNi-base thin films. Materials Science and Engineering A 273–275: 106–133.
12. Stoney GG (1909) The Tension of Metallic Films Deposited by Electrolysis. Proceedings of the Royal Society London A 82: 172–175.
13. Winzek B, Quandt E (1999) Shape-memory Ti-Ni-X-Films (X=Cu, Pd) under constraint. Zeitschrift fuer Metallkunde 90: 796–802.
14. Craciunescu CM, Li J, Wuttig M (2003) Constrained Martensitic Transformations in TiNiCu Films. Thin Solid Films 434: 271–275.
15. Craciunescu CM, Wuttig M (2003) New Ferromagnetic and Functionally Graded Shape Memory Alloys. Journal of Optoelectronics and Advanced Materials 5(1): 139–146.
16. Winzeck B, Schmitz S, Rumpf H, Sterzl T, Hassdorf R, Tienhaus S, Feydt J, Moske M, Quandt E (2004) Recent developments in shape memory thin film technology. Materials Science and Engineering A 378: 40–46.

# Exploiting Dynamic Cooperative Behavior
# in a Coupled-Core Fluxgate Magnetometer

V. In[1], A.R. Bulsara[1], A. Kho[1], A. Palacios[2], P. Longhini[2], S. Baglio[3],
B. Ando[3], V. Sacco[3], and J.D. Neff[1]

[1]  Space and Naval Warfare Systems Center, Code 2373, 49590 Lassing Road A341,
    San Diego, CA 92152-5001, USA
[2]  Mathematics Dept., San Diego State University, San Diego, CA 92182, USA
[3]  Dip. di Ingegneria, Univ. degli Studi di Catania, Viale A. Doria 6, 95125
    Catania, Italy

## 1 Introduction

Overdamped bistable dynamics, of the generic form $\dot{x} = -\nabla U(x)$, underpin
the behavior of numerous systems in the physical world. The most-studied
example is the overdamped Duffing system, the dynamics of a particle in
a bistable potential $U(x) = -ax^2 + bx^4$. Absent an external forcing term,
the state-point $x(t)$ will rapidly relax to one of two stable attractors, for
any choice of initial condition. It has been shown [1], however, that coupling
similar elements via a linear uni-directional coupling with cyclic boundary
conditions, can lead to oscillatory behavior past a critical value of the coupling
coefficient. Typically, this behavior is dictated by symmetry conditions [2], and
is generated via Hopf bifurcations; it appears to occur in any coupled system
of overdamped bistable elements, none of which would oscillate when isolated
and undriven, subject to the appropriate choice of parameters and operating
conditions (albeit through different bifurcation mechanisms).

  To better understand this behavior, we focus on a specific nonlinear dy-
namic system in which the state point $x(t)$ represents the (suitably normal-
ized) magnetic induction in a ferromagnetic sample; the dynamical model
is obtained via the continuum limit of a discrete spin model of individ-
ual domain dynamics and has been recently used [3] to characterize the
response of a specific magnetic measurement system, the (single-core) flux-
gate magnetometer (see [4] for good descriptions and references). The single-
core fluxgate magnetometer can be treated [3] as a nonlinear dynamic sys-
tem by assuming the core as approximately single-domain, and writing down
an equation for the evolution of the macroscopic magnetization parameter
$x(t)$: $\dot{x}(t) = -\nabla_x U(x)$ in terms of the potential energy function $U(x,t) =
\frac{x^2(t)}{2} - \frac{1}{c} \ln \cosh c[(x(t) + h(t)]$. $c$ is a temperature-dependent nonlinearity pa-
rameter which controls the topology of the potential function: the system

becomes monostable, or paramagnetic, for $c < 1$ corresponding to an increase in the core temperature past the Curie point. The overdot denotes the time-derivative, and $h(t)$ is an external signal. In single-core fluxgate magnetometers, $h(t) = h_1(t) + \varepsilon$. Here, $\varepsilon$ is the target signal (taken to be dc and quite small compared to the energy barrier height $U_0$ of the potential) to be detected, and $h_1(t) = A\sin\omega t$, a *known* bias signal that is applied to make the device switch readily between its saturation states (the target signal is too small to accomplish this alone), so that the asymmetry introduced by the dc target signal manifests itself as a (measureable) difference in the residence times spent in the "up" and "down" saturation states of the core hysteresis loop.

Single-core fluxgate magnetometers, that operate on the above dynamics, have always been of interest to the technical and scientific communities as practical and convenient sensors for vector magnetic field measurements requiring a resolution up to 100pT at room temperature; they have found applicability [4] in fields such as space and geophysical exploration and mapping, and non-destructive testing, as well as assorted military applications. These magnetometers are, most often, operated via a readout based on the spectral amplitude of the second harmonic of an applied time-dependent reference signal [5]; the presence of a target dc magnetic flux signal leads to the appearance of peaks at even harmonics of the reference signal frequency, with the peak amplitudes a function of the target signal. Recently, the possibilities offered by new technologies and materials in realizing miniaturized devices with improved performance, have lead to renewed interest in a new generation of cheap, compact and low-power fluxgate sensors. Miniaturization of the fluxgate sensors is complicated by the rapid increase of the magnetic noise with the (inverse) device dimensions, and general practical rules for achieving high sensitivity (large number of windings, large cross-sectional sensor area, and large driving current), but these requirements are at odds with the desired characteristics (low cost, power and noise) of the miniaturized sensors. Nonetheless, despite the difficulties manifest in integrated devices with better performance, the literature does contain good examples of fluxgate sensors in PCB [6] and even CMOS [7]. In particular, CMOS affords the possibility of realizing the sensing part (fluxgate) and the read-out circuit on the same chip, resulting in enhanced reliability, and lower costs in batch production.

An alternative to the power-spectral-based readout is the time-domain-based readout method which was first introduced by Strycker and Wulkan [8]. It relies on a compilation of the residence times in the stable steady states; this procedure, well established as the mechanism whereby neural firing events convey information in the nervous system, is the backbone of the readout scheme in our recent single-core fluxgate development work [3,9]. Hence, although the residence times based readout is used in our coupled sensor system, we do not describe it in this work; rather, we refer the reader to a recent publication [3] in which the technique, applied to a single-core fluxgate magnetometer in the presence of a noise floor, is detailed. It is easy to show that, for the single-core

fluxgate operated in the conventional (second-harmonic) mode, the device sensitivity is proportional to the amplitude and frequency of the applied bias signal; hence, the conditions for increasing the sensitivity lead to a greater on-board power consumption as well as an elevated noise floor since a significant fraction of the latter arises from the bias signal generator. By contrast, one can show [3, 9] that the sensitivity of the single-core fluxgate subject to the residence times readout is inversely proportional to the bias amplitude and frequency. Hence, the (theoretical) conditions for increasing the sensitivity in the residence times readout correspond to those for decreasing the onboard power requirement, as well as the noise floor. Of course, one cannot lower the bias amplitude and frequency too far; the frequency must remain high enough to ensure a good data sampling rate as well as strong electromagnetic coupling between the core and the readout circuit elements, while the amplitude must remain somewhat higher than the energy barrier height separating the steady states (these correspond to the saturation states of the hysteresis loop) of the potential energy function, to ensure reliable switching between the magnetization steady states in the absence of the target signal. Further, as the signal amplitude approaches the deterministic switching limit (the point at which the deterministic forcing alone is just sufficient to drive switching events) from above, the switching events become increasingly contaminated by the sensor noise floor, and a large number of crossing events must be gathered in order to obtain a reliable estimate for the difference in mean residence times. Hence, in practical applications, the decrease in onboard power (with the concommitant decrease in the noise contribution from the bias signal generator) must be offset against an increase in the observation time, required to obtain reliable statistics [3].

In two recent papers [11], we have demonstrated that coupling an *odd* number of overdamped bistable elements in a ring, with unidirectional coupling, can lead to oscillatory behavior when the coupling strength exceeds a critical value. Clearly, the practical importance of this effect lies in the potential sensitivity enhancement when the system is "tuned" very close to the oscillation threshold (i.e. in the regime of very low frequency). We now provide an overview of the (deterministic) results [11] of our investigations to date including, where applicable, experimental results that confirm our theoretical predictions.

## 2 Coupled-Core Fluxgate Magnetometers

We start by writing down the dynamics for three *uni*directionally (cyclically) coupled ferromagnetic cores ($i = 1 \ldots 3, x_4 \equiv x_1$):

$$\tau \dot{x}_i = -x_i + \tanh(c(x_i + \lambda x_{i+1} + h(t))), \tag{1}$$

where $x_i(t)$ represents the (suitably normalized) magnetic flux at the output (i.e. in the secondary coil) of unit $i$, and $h(t) \ll U_0$ is an externally applied

"target" magnetic flux, $U_0$ being the energy barrier height (absent the coupling) for each of the elements (assumed identical for theoretical purposes). In the absence of the coupling, each core is seen to be described by the dynamics already introduced in the preceeding section; the coupled system is, however, no longer describable by a (3-body) potential energy function, due to the unidirectional nature of the coupling. $\tau$ represents the device time constant (inverse bandwidth); it is a measure of how rapidly the switching events occur between the two steady states of each core. In well-fabricated "single domain" cores, this switching rate is quite large ($\tau \approx 10^{-6}$) with the hysteresis loop being narrow and sharp. It is important to note [11] that the oscillatory behavior occurs even for $h(t) = 0$, however when $h(t) \neq 0$, the oscillation characteristics change; these changes can be exploited for signal quantification purposes. We reiterate that the oscillations, corresponding to the periodic back-and-forth switching of the state point of each element between the stable (or saturation) states of the input-output (hysteretic) transfer characteristic and occurring when the coupling parameter $\lambda$ exceeds a threshold value $\lambda_c$, do *not* occur for a single (uncoupled) core, due to the overdamped nature of the dynamics (1).

## 2.1 dc Target Signal

We start with a straightforward numerical simulation of the system (1), taking $h(t) = \varepsilon(\ll U_0)$, a constant. The oscillations commence when the coupling coefficient exceeds a threshold value [11]

$$\lambda_c = -\varepsilon - x_{inf} + \frac{1}{c}\tanh^{-1} x_{inf} , \tag{2}$$

with the inflexion point (computed for $\lambda = 0$) $x_{inf} = \sqrt{(c-1)/c}$; note that in our convention, $\lambda < 0$ so that oscillations occur for $|\lambda| > |\lambda_c|$. We confine the system to a small neighbourhood of the bifurcation (i.e. having a small "separation" $\Delta\lambda \equiv \lambda_c - \lambda$); the benefits of this are, immediately, evident from a glance at Fig. 1 which shows the results of the simulations.

   For small separation $\Delta\lambda$, it is clear that the state-points spend the bulk of their transition times reaching the inflexion points $\pm x_{inf}$, after which the passage to the opposite minimum (at $\approx \pm 1$) is very rapid. Put differently, the combination of dc and coupled fluxes in each of the elements of (1), cause that particular potential to skew or tilt so that a minimum and the saddle point approach each other, coalescing into an inflexion point. At this point, an infinitesimal further tilt, causes the state-point to drop into the opposite minimum, all the time providing an input to the next (forward-coupled) element via the coupling, so that a soliton-like periodic disturbance travels around the ring. One also notes that the elements evolve (approximately) individually, with two element always remaining in (or very close to) their steady states (these have opposite signs) while the other evolves; this is particularly

**Fig. 1.** Emergent oscillatory behavior in the coupled core system (1) for $N = 3$. The top panel shows the oscillations near the critical point. Summed response is indicated by *thick black lines*, and individual element responses follow the *gray lines* in all panels. The parameters are set at $\lambda = -0.60, \varepsilon = 0$. The second panel shows the oscillations for a higher coupling strength $\lambda = -0.75$, and $\varepsilon = 0$; the oscillation frequency increases significantly, scaling as the square root of $|\lambda|$ and $\varepsilon$. The third panel shows the individual element oscillations for $\lambda = -0.60, \varepsilon = 0.1$; the frequency decreases (compared to the *top* panel) with increasing $\varepsilon$. The initial conditions for all simulation runs are $(x_1, x_2, x_3) = (1.0, 0.0, -1.0)$, $c = 3$, and the time step size is 0.001. For each panel, the critical coupling $\lambda_c$, at the onset of the oscillations, may be determined from (2)

evident in the slow oscillation regime at very small $\Delta\lambda$. In essence, the particular structure of our system introduces a "dynamic" time-scale separation wherein one element is always evolving much faster than all the others, until it reaches its equilibrium, at which point the next element acquires this rapid (by comparison to the others) time-scale. This happens even though the original dynamics (1) does not, explicitly, show a well-defined separation of slow and fast temporal scales. The above behavior, is reminiscent of what might be expected in a discrete line of magnetic spins, subject to a dc magnetic field. For an odd number of spins, there will always be two spins that have the same alignment, and are therefore "frustrated", with each spin trying to orient itself anti-parallel to the other; the net result is a "ripple" in the spin orientations that propagates through the chain, and continues to propagate as long as the boundary conditions are cyclic.

The considerations of the last paragraph lead to a calculation of the oscillation period in terms of the "up" and "down" state residence times $T_\pm$ of the *summed* output $\sum_i x_i(t)$. These are found as [11]:

$$T_+ = \frac{\pi}{\sqrt{cx_{inf}}} \frac{1}{\sqrt{\Delta\lambda}}; \quad T_- = \frac{\pi}{\sqrt{cx_{inf}}} \frac{1}{\sqrt{\Delta\lambda + 2\varepsilon}} \,, \qquad (3)$$

whence we can readily write down $T_\Sigma = T_+ + T_-$ for the period of the summed output, and $T_i = NT_\Sigma$ for each individual elemental period [11]. The period shows a characteristic dependence on the inverse square root of the bifurcation distance $\Delta\lambda$, as well as the target signal $\varepsilon$; these oscillations can be experimentally produced at frequencies ranging from a few Hz to high kHz. Recalling that the oscillations are, actually, switching events between the stable states of each core, it is clear that they will occur as long as at least one element has an initial condition that is different from the remaining elements (clearly, in any practical system with a noise floor, this condition is easily satisfied). The setup, clearly, eliminates the need to apply the reference bias signal, as would be required for the single fluxgate; by generating the oscillations we are, effectively, forcing each core to switch between its stable steady states (the saturation states of its hysteresis loop). Increasing $N$ changes the frequency of the individual elemental oscillations, but the frequency of the summed response is seen to be independent of $N$.

The residence times difference, i.e. the difference in the times spent by the system in its two stable magnetization states, can be computed (for the summed output signal) as $\text{RTD} = T_+ - T_-$; note that the (dimensionless and independent of $N$) residence times ratio (RTR) can also be computed from these residence times. The RTD vanishes (as expected) for $\varepsilon = 0$, and can be used as a quantifier of the target signal, analogous to the time-domain operation of the single fluxgate. For very small target signals $\varepsilon$, we can obtain an approximation to the coupled core RTD [11]:

$$\text{RTD} \approx \frac{\pi\varepsilon}{\sqrt{cx_{inf}}} \Delta\lambda_0^{-3/2} \,, \qquad (4)$$

where $\Delta\lambda_0 \equiv \lambda_c(\varepsilon = 0) - \lambda$. We readily observe that the sensitivity $\partial(\text{RTD})/\partial\varepsilon$ is significantly enhanced as we get closer to the critical point; this is illustrated in Fig. 2. An experimentally obtained analog of this result is shown in Fig. 3 (see [10, 11] for experimental details). We note that decreasing the temperature-dependent control parameter $c$ close to unity can also lead to enhanced sensitivity to small $\varepsilon$, as is readily apparent in (4). However, for $c < 1$ (corresponding to an increase in the temperature past the Curie point) the system ceases to be bistable and the material becomes paramagnetic. It is worth reiterating that a sensitivity $\partial T_\Sigma/\partial\varepsilon$, defined via the (summed) oscillation period, is actually proportional to $\varepsilon$ (for small $\varepsilon$) as can readily be calculated from (3). This may not be desireable in practical sensors where one would like to develop the optimal sensor configuration independently of the target signal. Hence, from this standpoint, the RTD constitutes the more desireable measure with its sensitivity proportional to the factor $\Delta\lambda_0^{-3/2}$ and independent of $\varepsilon$.



**Fig. 2.** Theoretical response curve of the coupled core fluxgate system to an applied target dc magnetic field $\varepsilon$ vs. coupling strength. As $\lambda$ approaches the critical value $\lambda_c = -0.5345$, the response curve rises almost vertically which suggests that the sensitivity of the device increases dramatically in this regime. $c = 3$, $\varepsilon = 0.1$, and $-1.0 \leq \lambda \leq -0.54$. Both the Residence Times Difference (*dashed curve*) and Residence Time Ratio (*solid curve*) response curves are plotted

The system sensitivity, defined via the derivative $\partial(\text{RTD})/\partial\varepsilon$, is found to increase dramatically as one approaches the critical point in the oscillatory regime; this suggests that careful tuning of the coupling parameter so that the oscillations have very low frequency (i.e. just past the critical point of the bifurcation), could offer significant benefits for the detection of very small target signals. This preceding statement must, however, be qualified by an important caveat: in practical setups the oscillation frequency cannot be set too low, in

**Fig. 3.** Experimentally obtained responsivity curves, using the Residence Time Ratio (RTR) vs. the applied target magnetic field $\varepsilon$ for different coupling strengths. As expected, the coupled core system is less responsive as the coupling strength is increased (*top* 3 curves). The *bottom* curve is the responsivity of an "equivalent" single fluxgate magnetometer with bias signal amplitude selected to be slightly suprathreshold, thereby yielding the maximal sensitivity. Note that the RTR for the single fluxgate is, conveniently, frequency independent

order to ensure good coupling between the cores and the circuit components; in turn, this places an upper bound on the sensitivity. The absence of the bias signal generator should also result in less system noise, however this must be balanced against the presence of the coupling circuits which require power and will, therefore, contribute to the sensor noise floor.

## 2.2 Time-Sinusoidal Target Signal

We now consider the situation wherein the target signal is $h(t) = \varepsilon \sin \omega t$ and the amplitude $\varepsilon$ is very small ($\varepsilon \ll U_0$ as before). Numerically integrating the system 1, with non-identical initial conditions, reveals three distinctive regimes of oscillatory behavior that are clearly separated (Fig. 4) in the parameter space $(\lambda, \varepsilon)$:

(I) The *supercritical* regime wherein the coupling parameter is below the critical value ($\lambda < \lambda_0$, i.e. $|\lambda| > |\lambda_0|$ in our convention). In this regime, the coupled system oscillates with a traveling wave pattern as described above, even for

**Fig. 4.** Theoretical Phase Diagram: Oscillatory behavior of the coupled fluxgate model 1 in parameter space $(\lambda, \varepsilon)$. In the supercritical regime, the oscillations form a traveling wave pattern. In the subcritical regime, with $h(t)$ small, the system oscillates about one of the steady states $\pm 1$, while with $h(t)$ large, the system oscillates between two steady states. In both cases the oscillation form a travelling wave and their frequency is exactly $\omega/3$. For $\varepsilon$ greater than a critical value, all three waves are in-phase with each other and frequency synchronized with the external signal $h(t)$ in region (III). Recall that $\lambda_0 < 0$ in our convention (see text). An experimentally obtained versionof this phase diagram is given in [12]

$h(t) = \varepsilon$, as long as the initial conditions for at least two of the elements are non-identical. In the presence of the target signal $h(t)$, the system responds by oscillating asymmetrically between the two stable magnetization states of each element. The response displays a frequency mixing of the inherent oscillations of the coupled system and the target signal. The bifurcation diagram (not shown) in this regime is quite complex ranging from the simple oscillations (for $h(t) = 0$), to quasi-periodicity and, eventually, to chaos.

(II) The *subcritical* regime wherein the coupling strength exceeds the critical value ($\lambda > \lambda_0$, $|\lambda| < |\lambda_0|$), so that there are no spontaneous oscillations. For small $h(t)$, the system oscillates about the steady states near $\pm 1$. With sufficiently large $h(t)$, the system oscillates between the two steady states in a travelling wave pattern where the amplitude and frequency of each oscillation are the same but a phase shift of $\frac{2\pi}{3}$ exists between the different waveforms. This behavior is quite similar to that already observed for the case of dc (or zero) target signal; however, the onset of the oscillations occurs sooner in parameter space when the applied signal is time-periodic. The oscillation frequency is exactly $\omega/3$.

(III) Frequency matching of the output waveform to that of the target signal. With the control parameter $\lambda$ held constant in the subcritical regime, increasing $\varepsilon$ past a critical value causes the coupled system to switch to another oscillation mode wherein the frequency of the output waveform precisely match with that of the target signal. This behavior occurs solely in the presence of a time-periodic applied signal, within the parameter space indicated in Fig. 4. For signal detection purposes, the subcritical regime is more relevant since it is relatively easy to extract information about the target signal, via

the RTD method, because of the simplicity of the oscillation characteristics e.g. constant amplitudes, frequencies, and phases; the RTD technique is not, however, as easy to use in the supercritical regime.

The important point to emphasize here is that, unlike the preceeding case of the dc target signal one can, for this case, generate oscillatory behavior (with a suitably chosen signal amplitude) even with the sensor set in the subcritical regime (i.e. no oscillations in the absence of the target signal). This affords a novel "power-saver" mode of operation for detecting time-periodic target signals. For a given $\omega$ and moderate values of $(\lambda, \varepsilon)$ above the boundary line for the supercritical regime, each element oscillates at $\frac{1}{3}\omega$, with an inter-element phase difference of $2\pi/3$. When the amplitude is large enough, the oscillations switch to an in-phase pattern with a frequency perfectly matched to the external signal frequency. This out-of-phase region is bounded (Fig. 4) by the *super*critical region (below) and the in-phase region (above). To the right, the region is bounded by the line connecting $\lambda_0$ and the critical signal amplitude $\varepsilon_c$ where the entrainment between the uncoupled ($\lambda = 0$) system and the external signal occurs. So the critical coupling $\lambda_{c_{sub}}$ for the onset of the oscillations, for a given $\varepsilon$, is

$$\lambda_{c_{sub}} = \lambda_0 - \left(\frac{\lambda_0}{\varepsilon_c}\right)\varepsilon\,. \tag{5}$$

Clearly, it is important to obtain an analytical expression for $\lambda_{c_{sub}}$ in terms of the critical value $\varepsilon_c$ for the onset of oscillatory behavior when the system is set to the subcritical regime. The calculation has been carried out in [12]; the result is presented here:

$$\varepsilon_c = \frac{A\sqrt{(b + 2aA)^2 + \omega^2\tau^2}}{\sqrt{1 + \omega^2\tau^2} + \sqrt{(b + 2aA)^2 + \omega^2\tau^2}}\,. \tag{6}$$

For the *super*critical case, one can readily write down (see Fig. 4)

$$\lambda_{c_{sup}} = \lambda_0 + \left(\frac{\lambda_0}{\varepsilon_c}\right)\varepsilon\,, \tag{7}$$

Here, the following definitions are used:

$$\begin{aligned}
a &= \frac{(4c^2 - 4c^2 e^{2c})e^{2c}}{(e^{2c}+1)^3}, \\
b &= \frac{(-3 + 4c - 8c^2)e^{2c} + (-3 + 4c + 8c^2)e^{4c} - e^{6c} - 1}{(e^{2c}+1)^3}, \\
d &= -(a + b) + \left(\frac{e^{2c}-1}{e^{2c}+1} - 1\right), \\
A &= (-b - \sqrt{b^2 - 4ad})/2a\,.
\end{aligned} \tag{8}$$

Figure 5 shows the oscillation characteristics in the different operating regimes of the device; the figure should be examined in the context of the phase diagram Fig. 4. This figure has been generated via our experimental setup [12].

**Fig. 5.** Oscillation waveforms associated with the different regimes in the *experimental* system. (**a**) System (emergent) oscillations (at 44 Hz) in the supercritical regime (see Fig. 4) without an external field. (**b**) characteristic modulation of the oscillations in the supercritical regime with a small applied ac magnetic flux signal (at 150 Hz). (**c**) Oscillations in the subcritical regime where the system oscillates at $\omega/3$ (50 Hz) with no modulation of the waveforms. (**d**) Oscillations in phase to each other and frequency-locked to the external signal. The external signal amplitude $\varepsilon$ increases from panels (a) to (d). Each core response is transduced from a core magnetization to a (easily measured) voltage output (see [9] for details)

In the experimental run, the system is set up with the coupling strength in the supercritical regime so that it is oscillating (44 Hz) without any applied external field (top panel of Fig. 5). The next panel illustrates the modulation of the oscillation waveforms by a small amplitude ac external signal (at 150 Hz) while the system is still in the supercritical regime; note that the the system remains oscillating at the natural frequency (44 Hz). Thereafter, increasing the amplitude of the ac signal pushes the coupled system into the subcritical regime (see Fig. 4), and the resulting oscillations occur (panel (c) of Fig. 5) at $\frac{1}{3}$ the frequency of the ac signal without the amplitude modulation of panel (b). The last panel illustrates the case when the ac signal amplitude is increased sufficiently (into region III) so that the system switches to another behavior in the subcritical regime where all three waveforms are phase-locked to each other and the oscillation frequency exactly matches that of the external signal. All four scenarios illustrated here are predicted by theory (Fig. 4) and verified in

numerical simulations (not shown). Clearly, a good theoretical understanding of the frequency mixing behavior affords an avenue to using these phenomena for the detection and quantification of time-periodic target signals.

# 3 Discussion and Concluding Remarks

We have introduced a new paradigm for operating fluxgate magnetometers. Instead of relying on a time-periodic bias signal to force a single-core device to switch states, we introduced a uni-directional coupling scheme for the sensor array to create switching behavior between the stable steady states of each core; the coupling, essentially, carries out the function of the applied bias signal in the single fluxgate. The coupled configuration affords the possibility of tuning the operational frequencies over a wide range for specific applications, by tuning the coupling strengths between individual sensors in the network. It also appears (although, clearly, more detailed investigations of this point are necessary) that the coupled core system could yield potential enhancements to the system sensitivity when compared to a single fluxgate, particularly if one is able to control the coupling parameter so as to always operate in the low frequency oscillation regime (just past the critical point). Other benefits are also realizeable in the context of signal detection applications: since the characteristic response of the coupled core system to the external field results in changes to both the oscillation frequency and the symmetry of the output waveform, we can also (in addition to the change in the RTD that has been the thrust of this paper) use these changes to quantify the external signal. The fact that the sensitivity, for the RTD readout in the coupled core system, is independent of the target signal is an important advantage to using this approach. Clearly, the relative simplicity of the readout electronics for the RTD method in general [9], as well as the potential to have a smaller noise floor (since the externally applied bias signal is unnecessary in the coupled setup) are also important factors when evaluating the potential usefulness of this setup for a real application. The last statement must, of course be coupled with the caveat of a contribution to the noise floor, arising from the coupling circuit. Hence, while we do expect the noise floor to be reduced from its current value (see preceding section) with the careful implementation of low-noise circuit components in future device realizations, the extent of the reduction is still unclear. The fact that the residence times readout achieves its optimal sensitivity (to a given target signal) under conditions that are opposite to those required to achieve optimal sensitivity in the conventional (spectral-based) readout and, furthermore, that implementing these conditions leads to a lower power consumption on the device, are also important considerations when deciding on the particular readout strategy to be employed, for a given application.

Another important point to note is that the implementation of the coupling requires a certain on-board power in the system described in this work; hence,

the coupled core system, when experimentally implemented, does not violate fundamental conservation laws. In a numerical simulation of the basic (1), one does not include an explicit power term, however, the initial condition for at least one element must be taken to be different from the others, in order for the oscillations to be generated (for an appropriately chosen value of the coupling parameter $\lambda$ past the critical value $\lambda_c$). Experimentally, in the presence of a noise floor, the individual fluxgates will, in fact, always have different initial configurations unless specific constraints are implemented; this also happens because it is impossible to have identical devices in practice.

A few additional comments regarding the sensitivity of the coupled device, particularly in the context of the comparison (Fig. 2) to the single RTD fluxgate, are in order. In the zero target field limit, one readily obtains the (theoretical) expressions for the sensitivities:

$$S_{s0} = \frac{4}{\pi\omega} \frac{1}{\sqrt{A^2 - b_t{}^2}}$$
$$S_{c0} = \frac{\pi}{cx_{inf}} \Delta\lambda_0^{-3/2} \tag{9}$$

for the single and coupled core systems, respectively; here, we recall that $h(t) = A \sin\omega t$ is the bias signal applied to the single-core device. Note that the first of these expressions is readily derivable from a simple, physically intuitive (level-crossing) description [3, 9] of the single device, subject to a time-periodic bias signal of amplitude $\varepsilon$ and frequency $\omega$, with $b_t$ being the location of one of the thresholds (the saturation state of the hysteresis loop in the absence of the dc target signal); the second expression is obtained by differentiation of (4). At first glance, it would appear that the sensitivity (for the residence times readout) can be enhanced by decreasing the bias signal frequency $\omega$ for the single fluxgate or, equivalently, by allowing $\lambda$ to be very close to its critical value (thereby, effectively decreasing the oscillating frequency) in the coupled core system. However, one must weigh this against the constraints of a (usually, in practice) finite observation time, as well as the necessity of maintaining a good magnetic-electric coupling between the core(s) and the circuit elements; this coupling breaks down at extremely low frequencies. Hence, one is limited to frequencies that should be above a threshold value (dictated by the device and circuit properties) for optimality in practical scenarios. It is, therefore, difficult to provide a completely accurate comparison between the performance of single and coupled devices, based on sensitivity alone. Note that using the Residence Times Ratio (RTR) removes the frequency from $S_s$, yielding a dimensionless quantity; while convenient for comparison purposes, this does not, however, remove the above-described constraint on the frquency. These considerations are, of course, complicated by the presence of background noise which can arise from a multitude of sources, as already discussed in the preceding section.

It is appropriate to discuss, briefly, the effect of the sensor noise floor on the response; in a real application, this noise can arise from internal (materials,

electronics, etc.) sources, as well as contamination of the target signal. The voltage output signal from a single core fluxgate (operated via the residence times readout) has been shown [9] to have a noise component that can be well-approximated by a gaussian distribution. However, the individual residence times have noise components which are, in general, non-gaussian; they have noise-dependent tails and, with increasing noise intensity the tails get longer, a feature that is quite common to two-state devices. Our earlier theoretical and experimental work [3] showed these features, and also showed that decreasing the noise intensity (alternatively, increasing the bias amplitude $A$) reduced the tail and made the distributions more gaussian-like; in the small $\sigma^2/A$ limit (where $\sigma^2$ is a theoretical noise variance) the residence times distributions are gaussian [3], which also has been observed (not shown) in experiments. Of course, this comes at the price of reduced sensitivity (since the sensitivity, for the residence times readout, is inversely proportional to the bias amplitude). The above ratio can be reduced (for a given noise floor) by increasing the bias amplitude $A$, but this increases the onboard power reqauirement as well as the contribution to the noise floor arising from the bias generator. A careful optimization of geometrical and other core parameters is also known to lower the noise in the voltage signal (see e.g. [13]).

To better understand the ramifications of background noise, we have introduced [3] the (critical to a practical system) observation time $T_{ob}$, and defined a response signal-to-noise ratio (SNR) which is directly proportional to $\sqrt{T_{ob}}$. A longer observation time leads to an enhanced response (to very weak target signals), however, practical constraints may limit $T_{ob}$. One can increase the bias frequency $\omega$, thereby increasing the number of crossing events and improving the statistics of the measurement process, however this implies a larger power requirement. Hence, in a practical application, one must strike a balance between the physical constraints (e.g. onboard power, noise from the bias signal generator) and the need to carry out a reliable measurement of the mean RTD. The practical configuration is also, of course, heavily dependent on the amplitude (relative to the energy barrier height) of the target signal $\varepsilon$ to be quantified; it helps to have an a priori idea of the range of target signals under consideration for a given application. However, if the target signal is larger than the energy barrier height (or the coercive field) it is, usually, easier to detect it by standard techniques that do not require a finely adjusted sensor such as that described in this work.

The noise floor of the (single core) residence times fluxgate can be obtained via the power spectral density of a time series of the voltage response, taken at 1 Hz (usually); this quantity can be shown to be proportional to the ratio of the standard deviation of the RTD (this quantity has been experimentally measured [9]), and can also be theoretically computed [3, 9]) and the Sensitivity $S$ (slope of the RTD vs target signal response curve). The theoretical results and calculation details will be presented elsewhere. The coupled core system's noise floor (measured via the response PSD [10]) is, currently higher than that of the single device, no surprise given the uni-directional coupling

and the added noise generated in each coupling circuit. We do, however, expect that the noise floor will be lowered in subsequent realizations of the coupled core device, following a careful optimization of the circuit components.

Our experimental results, together with all the caveats ennunciated above, make clear that this coupling scheme holds out the promise for a class of fluxgate sensors that would, potentially, overcome many of the limitations of using conventional (i.e., typically readout via the second harmonic in the PSD) single core fluxgates. However, a sensor based on our coupling scheme has not yet been realized (this review describes the results of the first experiments on a laboratory prototype), and must await the outcome of additional experiments aimed at quantifying the noise floor and responsivity of the coupled core system, together with precise comparison of its signal detection performance (and limitations) to a single core fluxgate.

Finally, it should be clear that our coupling scheme is quite general; it can readily be applied to a vast array of dynamical systems which follow the basic "particle-in-potential" paradigm with $U$ being any bi- or multi-stable potential and $x$ the appropriate state variable. The ability to control the oscillation frequency (our system can be made to oscillate at frequencies ranging from a few Hz to several kHz) dramatically broadens the range of applications that can benefit from this scheme. When the target signal is time-periodic, the sensor can be operated in a "power-saver" mode wherein it is set to the subcritical regime in the absence of the signal. Then, the emergent behavior can be used to quantify time-periodic target signals, since the internal oscillation frequency (44 Hz in our setup), can usually be controlled by an appropriate choice of system parameters (in this case, the coupling coefficients). The results of this work are expected to be applicable to a much larger class of nonlinear dynamic systems (of which our 3-core micro-fluxgate system [10] is only one example) coupled in this manner.

## Acknowledgements

## References

1. In V, Kho A, Neff J, Palacios A, Longhini P, Meadows B (2003) Phys. Rev. Lett. 91:244101.
2. Golubitsky M, Stewart I, Schaeffer D (1988) Singularities and Groups in Bifurcation Theory: Vol. II. In: Appl. Math. Sci. 69. Springer-Verlag, New York. Aronson D, Golubitsky M, Krupa M (1991) Nonlinearity 4:861. Golubitsky M, Stewart I (1999) Symmetry and Pattern Formation in Coupled Cell Networks. Springer-Verlag, New York.

3. Bulsara A, Seberino C, Gammaitoni L, Karlsson M, Lundqvist B, Robinson J (2003) Phys. Rev. E67:016120.
4. Ripka P (2001) Magnetic Sensors and Magnetometers. Artech House, Boston. Kaluza F, Gruger A, Gruger H (2003) Sensors and Actuators A106:48–51.
5. Primdahl F (1970) IEEE Trans. Magn. 6:376–383.
6. Tipek A, Ripka P, O'Donnell T, Kubik J (2004) Sensors and Actuators A115:286–292. Dezuari O, Belloy E, Gilbert S, Gijs M (2000) Sensors and Actuators A81:200–203. Kejik P, Chiesi L, Janossy B, Popovic R (2000) Sensors and Actuators A81:180–183.
7. Chiesi L, Kejik P, Janossy B, Popovic R (2000) Sensors and Actuators A82:174–180. Gottfried R, Budde W, Ulbricht S (1995) A miniaturized magnetic field sensor system consisting of a planar fluxgate sensor and a CMOS readout circuitry. In: Proceedings of Transducers'95.
8. Strycker S, Wulkan A (1961) AIEE Trans. 80:253–257.
9. Ando B, Baglio S, Bulsara A, Sacco V (2005) IEEE Trans. Instr. Measur. 54:1366 (2005) IEEE Sensors 5:895 (2005). Measurement 38:89 (2005).
10. In V, Sacco V, Kho A, Baglio S, Ando B, Bulsara A, Palacios A (2005) IEEE Sensors, submitted.
11. In V, Bulsara A, Palacios A, Longhini P, Kho A, Neff J (2003) Phys. Rev. E68:045102(R). Bulsara A, In V, Kho A, Longhini P, Palacios A, Rappel W-J, Acebron J, Baglio S, Ando B (2004) Phys. Rev. E70:036103.
12. In V, Bulsara A, Palacios A, Longhini P, Kho A (2005) Phys. Rev. E submitted.
13. Koch R, Rozen J (2001) Appl. Phys. Lett. 78:1897–1899.

# Motion Sensors and Actuators Based on Ionic Polymer-Metal Composites

C. Bonomo, L. Fortuna, P. Giannone, S. Graziani, and S. Strazzeri

Dipartimento di Ingegneria Elettrica, Elettronica e dei Sistemi-DIEES Università degli Studi di Catania – Viale A. Doria, 6-95100 Catania – ITALY

**Abstract.** Ionic Polymer Metal Composites or IPMCs are innovative materials, light and soft, that show very interesting electromechanical properties to be used in several fields of research, such as robotics, measurements, biomedics. In this paper details on IPMCs will be given: from their state of the art to their modeling and characterization; tools and equipments, designed and realized to perform measurements on the IPMCs will be presented.

**Keywords:** Displacement sensors, motion actuators, electromechanical conversion, IPMCs

## 1 Introduction

Several fields of research, such as robotics, measurements, biomedics, are interested in finding new materials. Most of the materials available on the market to realize motion systems are both rigid and high power consumer. Hence, the need to search for innovative materials, in order to emulate the movement of the biological tissues, emerges. Ionic Polymer Metal Composites or IPMCs seem to be a response to these requests. Indeed IPMCs are known also as *Artificial Muscles* [1] thanks to their capability to undergo large bending under the effect of a low voltage stimulus (the order of magnitude being few volts) [2]. Nevertheless the potentialities of IPMCs go beyond. They work also in a reversible way: if mechanically deformed they generate a voltage across their thickness. They can hence be used also as motion sensors [3]. A unique material offers the possibility to realize two important functions: perception and consequent activation. Moreover, being light and soft they result very silent. All these properties could not be fully exploited because of the lack of both: adequate technologies to fabricate the material and complete models able to describe their electromechanical behavior.

In the following sections an overview on IPMCs will be given, then preliminary results on their characterization will be reported and several

instrumentations and tools made to measure in order to collect all the data needed to model IPMC behavior, both as motion actuator and sensor, will be described.

## 2 IPMCs Fundamentals

A typical IPMC sample is a thin ionic polymeric membrane, the usual thickness being $200\,\mu m$, covered on both sides by metallic layers, to form the electrodes; the metal deposition thickness can vary from $5\,\mu m$ to $10\,\mu m$. The polymeric bulk can be Nafion® (perfluorosulfonate made by Dupont®) or Flemion® (perfluorocarboxilate made by Asahi Glass) [4], while noble metals are used to avoid oxidation phenomena and to improve the performance of the electrodes (i.e. platinum or gold). Figure 1 shows a microscopic image of the cross section of an IPMC made of Nafion® and platinum.



**Fig. 1.** An optical microscopy (20x) of the cross section of an IPMC. Platinum distribution after deposition on the Nafion® membrane surface is pointed out in the enlargement (50x)

The Nafion® molecule has the structure given by the following formula [5]:

$$[(CF_2CF_2)_n CFCF_2]_x$$
$$|$$
$$(OCF_2CF)_m OCF_2CF_2SO_3H$$
$$|$$
$$CF_3$$

Ionic polymers like Nafion® have inner ionisable groups; a property of these groups is that they dissociate and move in the molecular net in a variety of solvent media generating a strong electric field [1]. In the polymeric matrix of commercial Nafion® $-SO3^-$ is the fixed group while the cation $H^+$ is free to move.

By applying a voltage to a typical electrolyte, cations and anions move in opposite directions: no energy is transferred from the molecular network to the

solvent, and no solvent molecule is carried [6]. Instead, in the above-mentioned polymeric membrane the solvent molecules can be carried parasitically by the mobile cations; cations with a high hydration number will produce greater deformation than cations with a low hydration number [7]. For this reason, in practical applications, the hydrogen ion of the Nafion$^{\circledR}$ molecule is substituted, via an ion exchange process, with $Na^+$, $Li^+$, etc.

The acting and sensing mechanisms, involved in the electromechanical conversion exploited by the IPMCs, are shown in Fig. 2.



**Fig. 2.** The mechanisms involved in the electromechanical conversion exploited by the IPMCs. (**a**) for the actuator case: applied voltage into caused deformation; (**b**) for the sensor case: forced bending into voltage

For the actuator case two phenomena occur: the first is known as water pumping and depends on IPMC water content, the second is due to the interaction between charges relative to the particular metal electrodes and the polymer substrate. Assuming to impose a voltage across an IPMC strip an in deep description of involved phenomena is given.

Sodium ions (or the mobile cation whatever it is), following the electric field, move towards the cathode carrying with them water molecules, in a parasitical way. This *water pumping* causes the IPMC cathode side to spread, while the anode region decreases its volume and hence contracts. The result is the bending of the strip towards the anode. The higher are the ion concentration and the number of water molecules that the mobile ion can carry with it, the bigger will be the resulting deformation [8]. Figure 2(a) shows the cross section of the membrane, the two external lines correspond to the metal depositions that act as the electrodes where the voltage signal is applied.

If the applied voltage is constant, the cation electric drift produces a concentration gradient that generates a consequent back diffusion of water molecules (cation back diffusion current is balanced by the drift current due to the electric field). The macroscopic effect is the partial relaxation of the IPMC strip that goes on until water distribution inside it becomes uniform [8].

If the applied voltage is time variant, for instance a sinusoidal voltage, the strip starts oscillating and the maximum displacement depends on both the amplitude and the frequency of the exiting input signal [9]. The IPMC actuator has a good low frequency response; at high frequency a number of phenomena do not allow it to follow the rapid signal variations. The frequency range where the deformation of the IPMC is appreciable was experimentally found to be approximately from 0 to 30 Hz. For higher frequency reported studies show that deformation reduces remarkably [10].

The second contribution to the IPMC deformation, when a voltage is applied, is due to the coulombian interactions that arise between charges in metal electrodes and the fixed charged groups belonging to the polymer matrix. The porous nature of Nafion$^{\circledR}$ makes the deposed metal to penetrate inside the membrane, forming dendritic electrodes [11]. The external voltage charges these dendrites that tends to go away one from each other, because of electric repulsion, causing the electrode to expand. The charge of the anionic sulfonic groups inside the polymer chain amplify this repulsive phenomenon near to the cathode region, while it reduces it in the anion region. The overall effect is a further expansion of the cathode side.

For the sensor case, by mechanically bending the material it is possible to change the distribution of the charges with respect to the membrane neutral axis (see Fig. 2(b)): the applied stress will contract one side of the membrane while will spread the other, the mobile ions will move consequently toward the region characterized by a lower charge density parasitically carrying the solvent molecules (i.e. deionised water). A deficit of negative charges and an excess of the positive ones will therefore result in the expanded side. In the contracted side the opposite will occur. This phenomenon produces a voltage gradient collected at the metal electrodes. It is intuitive as this property results in a sensing capability [12].

## 3 IPMCs Modelling

To exploit the potentiality of IPMCs the development of complete model able to fully describe their behaviour, both as motion actuators and sensors, is mandatory. This task has been accomplished by using a grey-box approach modelling: models are obtained on the basis of few well understood phenomena ruled by physical parameters that can be estimated on the basis of experimental observations on the devices. The models are scaled by using parameters that are under the control of the designer.

Several IPMC samples were built following the standard procedure made by Oguro [13]. The samples used are made of Nafion$^{\circledR}$ 117 and 115 as ionic polymer (produced by Dupont$^{\circledR}$ and distributed by Sigma-Aldrich$^{\circledR}$) with Lithium and Sodium as counterion, while Platinum has been chemically deposed to form the electrodes.

**Fig. 3.** The cutter used to obtain strips of IPMC of predefined size. The cutter was designed ad-hoc

Obtained sheets were cut into strips of different size, using an *ad-hoc* designed and built precision cutter, in order to obtain the strips with repeatable dimensions. A photo of the cutter is shown in Fig. 3.

As far as the actuator is concerned the model input is the imposed voltage, the output is the available force (blocked force) or the displacement of a given point (free deflection).

As far as the sensor is concerned the model input is represented by a deformation imposed and the output is the electrical produced signal.

The models have been developed referring to the largely known model of the Euler-Bernoulli beam, pinned by one end. This is the natural choice since IPMCs electromechanical conversions occur in the transverse direction. A scheme of the used cantilever configuration is shown in Fig. 4.



**Fig. 4.** A scheme of the cantilever configuration used to model the IPMC

**F**   is the force;

$\delta$   is the transversal deformation;

**Lc**  is the length of the clamped part of the beam;

**Lt**  is the total free length of the beam (without considering the length of the pinned part);

**Ls**  is the point where the developed mechanical reaction is desired (actuator), or the mechanical (force or displacement) stimulus is applied;

**w, t** ware the dimensions of the beam cross section.

## 3.1 As Actuators

The model, for the actuator, is organized in accordance with the following Figs. 5a and 6a. The corresponding scheme of the measuring facilities used to acquire relevant data are reported in Figs. 5b and 6b, respectively.



*(a)*                                                        *(b)*

**Fig. 5.** A block scheme for the actuator model in the case of the blocked force (**a**) and the scheme of the measuring facility used to acquire relevant data (**b**)

The first block is a nonlinear one that is introduced to convert the imposed voltage into the absorbed current. The introduction of the nonlinearity is necessary because of the non linear effects that are widely reported in the literature [14] and that have been observed in experimental data acquired on the used membranes. Figure 7 shows a typical voltammogram relative to one sample of IPMC Nafion$^{\circledR}$ 117, with Li$^+$ as counterion, with length 25 mm and width 6 mm when a sinusoidal voltage with 6 Vpp and frequency 200 mHz is applied.

The modelling of the electromechanical conversion has been obtained by the following assumptions:

1. the cause of the membrane bending is represented by charge accumulation close to the electrodes;
2. the electromechanical coupling phenomenon can be described by adapting a model already accepted for piezoelectric materials;

*(a)*                                                              *(b)*

**Fig. 6.** A block scheme for the actuator model in the case of the free deflection (**a**) and the scheme of the measuring facility used to acquire relevant data (**b**)



**Fig. 7.** A typical voltammogram relative at Nafion$^{®}$ 117 Li$^{+}$

3. equations ruling the mechanical deformation of the membrane are those of a pinned beam subjected to a distributed moment.

The first conversion, i.e. from applied voltage to absorbed current has been solved by using a lumped parameter circuit.

Equivalent-circuit models with lumped parameters are advantageous because they can provide an intuitive, graphical representation of the governing equations of a system. Defined properly, the individual circuit elements can have clear physical interpretations, and the user can examine the relationship between the various elements without being forced to study the underlying equations.

The proposed circuit is a generalization of similar models reported in the literature [15–18]: a nonlinear branch has been introduced that describes the

nonlinearity observed in the dependence between the voltage applied to an IPMC and the absorbed current [19].

The absorbed current is the cause of the mechanical reaction (as widely reported in the literature) because of charges/water redistribution [20–22]. The membrane is used in a pinned configuration, hence it produces a force in the direction orthogonal to its length and/or the deformed membrane shape corresponds to a circumference arc [23, 24]. Deformation results in a change of the sample curvature $(\text{mm}^{-1})$.

## 3.2 As Sensors

The state of the art about IPMC models in the sensing working mode is even less developed than the state of the art of IPMCs as actuators. The model of the IPMCs working as sensors considered in this paper has been determined by using the same approach developed for the actuators. In particular to determine the parameters that characterise the IPMC models experiments have been performed by using ad-hoc measuring facilities, designed and realised at the DIEES laboratory. The instruments used at this stage are designed to measure both the imposed force and/or deformation and the sensing signal, produced by the IPMC. The model input is the imposed displacement at the point $L_s$, the output can be either the charge (with no drawn current) or the current (in short circuit conditions).

The sensor model is scaled as a function of relevant sensor parameters (counterion, length, thickness, width).

The models, for the sensor, are organized in accordance with the following Figs. 8a, 9a. the block scheme of the measuring tools used to acquire data on relevant electromechanical parameters are reported in Figs. 8b and 9b.



*(a)*      *(b)*

**Fig. 8.** The block scheme of an IPMC based sensors when the produced charge is measured (**a**) and the corresponding block scheme of the measuring apparatus (**b**)

Note that in this case the non linear effect has been not included because of the experimental evidence, obtained by considering the sensing signal produced by a sinusoidal mechanical input signal, no significant voltage distortion was observed. An example of typical input-output signals are reported in Fig. 10.

**Fig. 9.** The block scheme of an IPMC based sensors when the produced current is measured (**a**) and the corresponding block scheme of the measuring apparatus (**b**)



**Fig. 10.** The sensing output obtained measuring the short sensing current (green line) when a deformation (blue line) is imposed at the membrane

The Fig. 10 shows the output obtained measuring the short sensing current (green line) when a deformation (blue line) is imposed to the membrane. In the considered case the frequency of the oscillations imposed to the sample is about 24 Hz.

This gives evidence that IPMCs working as sensors can be modelled by using a linear approximation.

Also the hydration working conditions are significantly different from those used for the actuator. In fact, the sensing behaviour of IPMCs is better when the membrane is dry. A number of experimental observations revealed that the sensor works better when the water excess is removed. Figure 11 shows this behaviour.

Moreover, the sensing signal maintains the same characteristics after hours of working.

**Fig. 11.** (**a**) The produced signal for an highly hydrated sample (**b**) the same signal after 1 h of working with the same applied deformation

## 3.3 The Instruments for Testing the Actuator

To identify the actuator model, measurements of a set of quantities are required. In particular the identification of the electric model requires to measure both the voltage applied across the thickness of the membrane and the corresponding absorbed current. For this reason an *ad-doc* circuit, able to sense both the absorbed current and the voltage applied has been developed. Figure 12 shows the realized prototype. The circuit is composed by an instrumentation amplifier that detects the voltage across a shunt resistance of value 0.1 $\Omega$ and tolerance 1%, in series to the membrane.



**Fig. 12.** The prototype of the circuit able to measure the voltage across the thickness of the membrane and the absorbed current

As previously mentioned, the electromechanical model transforms the applied voltage into the blocked force. An instrument to capture the developed force and the stimulus signal has been designed. The IPMC strip is clumped edgeways in a cantilever configuration and is forced by a voltage stimulus. Relevant signals are conditioned by using the circuit shown in Fig. 13 and

(a)



(b)

**Fig. 13. (a)** A top view of the instrument used to measure the developed force of the membrane and the stimulus applied. **(b)** A typical plot of the acquired signals: the input stimulus (blue line) and the corresponding blocking force (green line)

are successively sent to an acquisition card mounted in a PC. Contemporary the force produced by the sample working as an actuator, is measured by a load cell model GSO-10 by Transducer Techniques for micro-force detecting, equipped with its conditioning circuitry. The load cell has to be calibrated by setting the gain of the conditioning circuitry in order to have an output voltage equal to the measured force. This result is obtained by using a known weight, following a procedure suggested by the load cell producer.

Figure 13a is a top view of the instrument used to acquire both blocking force and stimulus signals. The result of a typical acquisition is reported in Fig. 13b.

The deformation, produced by the membrane when a signal stimulus is applied, is another important quantity, used to identify the electromechanical model of the actuator. The absorbed current and the applied voltage are acquired by using the conditioning circuit introduced before. The deformation of the membrane is measured by using a system based on infrared (IR) sensors. In particular, the IR system is composed of a transmitter and a receiver. It is able to reveal the distance of a target by measuring the intensity of the reflected ray that is obviously modulated by the target movement. The IR system together with the required conditioning circuitry (amplifiers, demodulating stage, and filters) is enclosed in a screened box to avoid any interference. Figure 14a is a view of the whole system realized to perform the required measuring surveys, reported in Fig. 14b [25].

### 3.4 The Instruments for Testing the Sensor

As far as the sensor model is concerned, an *ad-hoc* instrument has been realized in order to measure the sensing signal produced by the membrane when a force produces a deformation [26]. The produced signal is detected by using the condition circuit previously described. In particular the adopted circuit allows the measure of the short-circuit current. By adequately design the circuit it is possible to extend in a wide range of frequency the measure of this quantity. On the contrary, the measure of the produced charge when the membrane is deformed requires the design of the circuit in order to chose the range of frequencies in which the acquisition will be done. This induces to prefer a measure of the short-circuit current because the circuit is more accurate than the charge amplifier. A top view of the card used to measure the short-circuit current produced by the membrane, working as a sensor in shown in Fig. 15.

Also for the case of the sensor model identification, a set of electromechanical quantities needs to be measured. In particular the quantities of interest are: the applied deformation, the short-circuit current and the blocked force.

A view of the real system is given in Fig. 16a. The force is detected by using the load cell before mentioned, the deformation is acquired by using a system based on IR sensors. Typical acquired signals are shown in Fig. 17.

## Conclusions

In this contribution the actual trends on IPMC technology is outlined and some new results are included. The approach adopted to model the IPMCs both as motion actuator and sensor is presented with several instruments ad hoc built to address the characterization step.

(a)



(b)

**Fig. 14.** (**a**) A top view of the instrument used to measure the deflection of the membrane and the stimulus applied; (**b**) Typical signals acquirable by the instrument, from the *top*: the applied voltage, the current absorbed by the IPMC and its consequent deformation

**Fig. 15.** The prototype of the circuit for the measurement of the short-circuit current produced by the membrane working as sensor



**Fig. 16.** (**a**) A view of the system able to perform the measurements required for the identification of the sensor model. (**b**) A particular of the system that shows the membrane positioning

## Acknowledgements

## References

1. M. Shainpoor, Y. Bar-Cohen, J.O. Simpson, and J. Smith, "Ionic Polymer-Metal Composites (IPMC) as Biomimetic Sensors, Actuators, and Artificial Muscle - A Review", Int. J. Smart Materials and Structures, Vol. 7, pp. R15-R30, (1998)
2. K. Oguro, Y. Kawami, and H. Takenaka, "Bending of an Ion-Conducting Polymer Film-Electrode Composite by an Electric Stimulus at Low Voltage", Trans. Journal of Micromachine Society, Vol. 5, pp. 27–30, (1992)

**Fig. 17.** The signals involved in the sensor model identification, acquired by the instruments shown in Fig. 16: (**a**) the mechanical stimulus applied to the membrane, is a pulse train followed by a frequency sweep and a noise signal; (**b**) the corresponding blocking force; (**c**) the relative sensing current

3. K. Sadeghipour, R. Salomon, and S. Neogi, "Development of a Novel Electrochemically Active Membrane and 'Smart' Material Based Vibration Sensor/Damper", Smart Materials and Structures, Vol. 1, pp. 172–179, (1992)
4. Y. Bar-Cohen, "Electroactive Polymers as Artificial Muscles – Capabilities, Potentials and challenges", HANDBOOK ON BIOMIMETICS, Section 11, Chapter 8, "Motion" paper #134, NTS Inc. Aug. 2000
5. Nafion$^{\circledR}$ Resins, Aldrich Technical Information Bulletin AL-163
6. Y. Bar-Cohen, S. Leary, J.O. Harrison, J. Smith, "Electro-Active Polymer (EAP) actuators for planetary applications", Proc. of SPIE's $6^{th}$ Annual International Symposium on Smart Structures and Materials, Paper N° 3669-05, Newport Beach, CA, Mar. 1999
7. M. Shahinpoor, K.J. Kim, "The Effect of Surface-Electrode Resistance on the Performance of Ionic Polymer-Metal Composite (IPMC) Artificial Muscle", Smart Material and Structures Journal – Vol. 9, N° 4, pp. 543–551, (2000)
8. Sia Nemat-Nasser, "Micromechanics of Actuation of Ionic Polymer-metal Composites", Journal of Applied Physics, Vol. 92, num. 5, Sep. (2002)
9. M. Shanipoor, "Electro-Mechanics of Iono-Elastic Beams as Electrically-Controllable Artificial Muscles", Proceedings of SPIE's 6th Annual International Symposium on Smart Structrures and Materials (1999)

10. M. Shahinpoor and K. J Kim, "Ionic polymer–metal composites: I. Fundamentals", Smart Mater. Struct. 10 819–833 (2001)

11. S. Nemat-Nasser and Y. Xian Wu, "Comparative experimental study of ionic polymer–metal composites with different backbone ionomers and in various cation forms", Journal of Applied Physics Volume 93, Num. 9 (2003)

12. M. Shahinpoor, "A New Effect in Ionic Polymeric Gels: The Ionic "Flexogelectric Effect," Proc. SPIE 1995 North American Conference on Smart Structures and Materials, February 28-March 2, 1995, San Diego, California, vol. 2441, paper no. 05, (1995)

13. Keisuke Oguro, Preparation Procedure Ion-Exchange Polymer Metal Composites (IPMC) Membranes, on web at http://ndeaa.jpl. nasa.gov/nasande/lommas/eap/IPMC_PrepProcedure.htm

14. S. Leary and Y. Bar-Cohen, "Electrical Impedance of Ionic Polymeric Metal Composites", Proceedings of SPIE's 6th Annual International Symposium on Smart Structures and Materials, 1–5 March, 1999, Newport Beach, CA. Paper No. 3669-09 (1999)

15. K.M. Newbury, "Characterization, modeling, and control of ionic polymer transducers", PhD Dissertation Virginia Tech etd-09182002-081047

16. X. Bao, Y. Bar-Cohen, S.-S. Lih, "Measurements and Macro of Ionomeric Polymer-Metal Composites (IPMC)", Proceedings of the SPIE-EAPAD Conference (2002), paper 4695-27.

17. J. Paquette, K.J. Kim, J.-D Nam, Y.S. Tak, "An equivalent Circuit Model for Ionic Polymer.Metal Composites and Their Performance Improvemet by Clay-Based Polymer Nano-Composite Technique", J. of Int. Mat. Sys. and Struc., 14, 2003, pp. 633–642.

18. R. Kanno, S. Tadokoro, T. Takamori, M. Hattori, "Linear Approxiamte Dynamic Model of ICPF (Ionic Conducting Polymer Gel Film) actuator", Proc. IEEE. Int. Conf. on Robotics and Automation, Minneapolis, MN, Apr. 1996, 219–225

19. C. Bonomo, L. Fortuna, P. Giannone, S. Graziani, "A Circuit to Model the Electrical Bahaviour of an Ionic Polymer-Metal Composite", accepted for the publication on IEEE CAS I.

20. M. Shahinpoor, Y. Bar-Cohen, T. Xue, J.O. Simpson, J. Smith, "Ionic Polymer-metal Composites (IPMC) As Biomimetic Sensors and Actuators", Proc. SPIE-EAPAD 1998, Conference (2002), paper 3324–27.

21. S. Nemat-Nasser, "Micromechanics of actuation of ionic polymer-metal composites", Journal of Applied Physics, vol. 92 number 5 (2002) 2899–2915.

22. S. Nemat-Nasser, J. Yu Li, "Electromechanical response of ionic polymer-metal composites", J. of Applied Phys. 87, 7, 3321–3331, 2000.

23. Z. Chen, X. Tan, M. Shahinpoor, "Quasi-static Positioning of Ionic Polymer-Metal Composite (IPMC) Actuators", http://www.egr.msu.edu/~xbtan/Papers/aim05.pdf.

24. Y. Bar-Cohen, X. Bao, S. Sherrit, S.-S. Lih, "Characterization of the Electromechanical Properties of Ionomeric Polymer-Metal Composite (IPMC)", Paper 4695-33, Proc.of the SPIE Smart Structures and Materials Symposium, EAPAD Conference, San Diego, CA, March 18–21, 2002.

25. C. Bonomo, L. Fortuna, P. Giannone, S. Graziani, "Frequency response analysis of IPMC actuators by an IR system", Proceedings of SPIE – Volume 5759 Smart Structures and Materials 2005: Electroactive Polymer Actuators and Devices (EAPAD), May 2005, pp. 41–48.

26. C. Bonomo, L. Fortuna, P. Giannone, S. Graziani, S. Strazzeri, "A Method to Characterize IPMC Membrane Sensor" Proceedings of IMTC 2005 – Instrumentation and Measurement Technology Conference Ottawa, Canada, 17–19 May 2005

# Contributed Papers

# Pattern Formation Stability and Collapse in 2D Driven Particle Systems

M.R. D'Orsogna[1], Y.-li Chuang[2], A.L. Bertozzi[1], and L.S. Chayes[1]

[1] Department of Mathematics, UCLA, Los Angeles, CA 90095
    dorsogna@math.ucla.edu
[2] Department of Physics, Duke University, Durham, NC 27708

Interacting, multi-robot systems show increasing promise for advances in exploration and defense applications. Here, we model a non-linear system of self-propelled individuals interacting via a pairwise attractive and repulsive potential. Depending on the interaction parameters, the agents may disperse, accumulate into self-organizing structures such as flocks and vortices, or collapse onto themselves. Borrowing tools from Statistical Mechanics, we discuss the connections between the H-stable nature of the interaction potential and resulting aggregating patterns and asymptotic behaviors.

## 1 Introduction

Designing and controlling robot assemblies to achieve specific collective goals has drawn considerable interest in recent years [1–4].An individual agent may be programmed to be fully autonomous and independent, but because of physical and resource constraints, its abilities may be limited. On the other hand, groups of individuals exchanging information and optimally self-organizing may have a much broader range of capabilities. Natural examples of interacting "swarms" abound: fish, birds, bacteria and insects communicate to create complex patters with new and useful group properties [5, 6]. These structures often form without the aid of a leader or of mediating chemical fields and only in response to local interactions. The underlying idea of robot swarming is to coordinate agent motion in a similiar, intelligent manner, sometimes looking at nature for inspiration [7]. Teams of interacting artificial agents may someday be routinely used in underwater or space exploration missions, or for the completion of military and other dangerous tasks, such as land-mine detection and removal or earthquake recovery [1].

A major issue is the control of swarm size and stability with respect to constituent number. Given $N$ agents that interact through simple rules, how does the size of the swarm depend on $N$? Can we design interactions so that swarms are stable as the number of constituents increases and inter-agent

**Fig. 1.** H-stable and catastrophic behaviors. The pairwise panels I, II show the two-body interaction between agents. In both cases a minimum separation distance exists. If agents are placed on an infinite triangular lattice of spacing $d$, the overall energy is minimized in dramatically different ways. For the pairwise I case, in the large number limit, agents are separated by a finite distance. For the pairwise II case, the infinite system energy is minimized with lattice constant $d \to 0$. This is the catastrophic case where agents collapse upon themselves

spacings are fixed without collapsing or dispersing as $N$ grows? Obviously the answer depends on the particular agent-agent interaction. Here, we discuss stability and pattern formation of a prototypical system of $N$ self-propelled individuals interacting through pairwise attractive and repulsive potentials. We have addressed most of these issues in [12].

## 2 Stability and Collapse

Many-body, interacting individuals are often encountered in physical systems at the microscopic scale. Statistical mechanics aims at describing such aggregates in a more macroscopic, "thermodynamic" way [8]. To make the passage to the macroscopic world meaningful so that thermodynamics can be fully recovered, the pairwise interactions must obey specific requirements. In particular, for any arbitrarily large number $N$ of agents, if a constant $B \geq 0$ exists such that $\sum_i U(\mathbf{x}_i) \geq -NB$ the microscopic agents will not collapse onto themselves and a typical distance between individuals will be well defined.

**Fig. 2.** H-stability phase diagram of the Morse potential (taken from [12]). Catastrophic and stable behavior are predicted as a function of the parameter ratios $\ell = \ell_r/\ell_a$ and $C = C_r/C_a$. Extrema of the potential $d_{min}$ exist only for $\ell > \max\{1, C\}$ and for $\ell < \min\{1, C\}$. In these cases $d_{min} = \ell_r \log\left(\ell/C\right)/\left(\ell - 1\right)$. For each region a qualitative pattern outcome is shown. In region V agents disperse to infinity

This is the fundamental property of *H-stability*. Systems that do not obey this constraint are called catastrophic. Upon increasing the number of agents, the latter will tend to accumulate at the same region in space. For catastrophic systems the thermodynamic limit cannot be defined. We will apply these concepts to our multi-vehicular ensembles.

Mathematically, H-stability translates to many conditions on the pairwise potential [9, 10]. For example if its spatial integral is negative, the system is proven to be catastrophic and collapse will occur. An illustration is given in Fig. 1, where two cases of pairwise potentials and their corresponding many-body energy on an infinite lattice are shown. The two pairwise curves are qualitatively similar: both are soft-core, have a minimum and decay to zero exponentially. The pairwise potential on the right, however, subtends a negative area, signaling catastrophic behavior. This is seen, for example, as agents are placed on an infinite triangular lattice of variable lattice constant $d$. For the H-stable potential, depicted on the left, a typical spacing emerges that minimizes the energy: agents assemble at this finite distance even as $N \to \infty$. In particular, the system behaves extensively and agents do not collapse onto each other. The right panels correspond to the catastrophic regime. Here, a minimum exists for the two body potential, as seen in the upper figure. However, as $N \to \infty$ global energy minimization occurs for inter-agent spacing $d = 0$. The system will now collapse onto itself in the large agent limit. Similar trends persist when agents assemble on a square lattice. Other rules for H-stability are given in D. Ruelle's book [9]. Lennard-Jones and hard-core potentials are always stable: for systems interacting according to either, catastrophic behavior never occurs.

One other requisite for thermodynamic, macroscopic behavior, is that of temperedness, namely, that the long range attractive part of the potential should not be too strong. In particular, it can be shown that in $d$ dimensions the attractive part of the potential should go to zero faster than $r^{-d}$. Three dimensional gravitational forces scale as $r^{-1}$: the universe does not obey the laws of thermodynamics! The same can be said about Coulomb forces between same charge particles [9, 10]. The generalized Morse potential of (3) decays exponentially, and satisfies the temperedness condition regardless of the potential parameters. Macroscopic behavior therefore can be violated only in the catastrophic regime.

## 3 The Model

We let the $1 \leq i \leq N$ discrete swarming agents be governed by the following equations of motion [11, 12]:

$$\frac{d\mathbf{x}_i}{dt} = \mathbf{v}_i \,, \tag{1}$$

$$m_i \frac{d\mathbf{v}_i}{dt} = (\alpha - \beta|\mathbf{v}_i|^2)\mathbf{v}_i - \boldsymbol{\nabla}_i U(\mathbf{x}_i) \,. \tag{2}$$

The self-accelerating force $\alpha\mathbf{v}_i$ tends to balance the self-decelerating term $\beta|\mathbf{v}_i|^2\mathbf{v}_i$, giving individuals the tendency to travel close to the self-propelling speed $|\mathbf{v}_i| = \sqrt{\alpha/\beta}$, a mechanism first introduced by Rayleigh [14, 15]. Note that the self-propelling mechanism tends to fix the magnitude of the velocity but not the direction of motion. Two major dynamical outcomes, mostly dependent on the initial conditions, are then possible: circular motion or a coherent agent drift. Interactions follow the generalized Morse potential:

$$U(\mathbf{x}_i) = \sum_{j \neq i}[C_r e^{-|\mathbf{x}_i - \mathbf{x}_j|/\ell_r} - C_a e^{-|\mathbf{x}_i - \mathbf{x}_j|/\ell_a}] \,, \tag{3}$$

where $\ell_a, \ell_r$ represent the range of the attractive and repulsive part of the potential and $C_a, C_r$ are their respective amplitudes. As we shall see, the existence of a preferred speed coupled with pairwise interactions among agents leads to interesting aggregation patterns. The potential of (3) is used as an example of soft-core potential; combinations of other attractive and repulsive terms lead to collective trends that can be easily understood through this specific example. For simplicity, we shall only consider a 2D description and focus here on identical mass agents: $m_i = m$.

Actual realizations of self-propelled vehicles interacting according to virtual Morse potentials have been already introduced in the robotics literature [1, 13]. For example, in [13], "Kelly" vehicles were built to be self-propelled by a set of two fans separated by a variable distance and driven by virtual

**Fig. 3.** A "Kelly" vehicle guided by two virtual potentials of the Morse type: one to guide the vehicle towards a target site, the other to avoid an obstacle. From [13]

attractive and repulsive forces. The equations of motion for speed and angular velocity of the Kellys were directly mapped onto equations similar to those presented in (1)–(3). In particular, changing the self-propulsion and the Morse parameters enabled the vehicle to effectively orient itself towards an attractive target, or to avoid multiple stationary obstacles. Further cooperative strategies could also be implemented so that a many-vehicle system could be used to effectively search more than one location. An example of the Kelly motion is given in Fig. 3.

In the remainder of this paper we discuss the consequences of decentralized control of the form presented in (1–3) when applied to a large system of vehicles interacting with each other.

## 4 Pattern Formation

In this section we apply the criteria for H-stability to the Morse potential of (3). The choice of the parameters $C \equiv C_r/C_a$ and $\ell \equiv \ell_r/\ell_a$ determines stable or catastrophic behavior [12], as seen in the phase diagram of Fig. 2. A similar analysis can be extended to most potentials introduced in the robotics literature, for example that of [16, 17], where catastrophic behavior is seen.

By numerically integrating (1–3) [18] with free boundaries and random initial conditions, we can distinguish different aggregation regimes in the $\{C, \ell\}$ phase space. All are consistent with the stable or non-stable predictions of Fig. 2. Regions I through IV of the phase diagram define catastrophic potentials and structures decrease in size as $N$ increases. In Fig. 4 we show different aggregating patterns. Clumps form whenever the pairwise interaction admits a minimum (region I) and rings occur when that minimum is zero (region II). In regions III and IV, where the pairwise potential does not allow for a minimum, clumped rings form (not shown) as a way to minimize the total

**Fig. 4.** Aggregating geometries for different Morse potential parameters and for $N = 100$ agents. From left to right: clumps (region I), a ring (region I), a vortex and a flock (region VII)

energy while keeping a constant speed. Region V corresponds to dispersive behavior, where the agents occupy the entire volume.

Regions VI and VII are the most interesting of the phase diagram. In the stable region VI, coherent structures can form only at relatively low values of $\alpha/\beta$, when the kinetic energy of the agents is comparable to the confining interaction potential. Swarming individuals assemble in a flock or in a disk, depending on the initial conditions; generally spacings are well defined, the motion is rigid-body like and the structures are extensive. For the case of rigid-body motion, the agents do not define a stationary center of mass, but rather the latter executes a non trivial trajectory. For larger values of $\alpha/\beta$, individuals disperse.

In the catastrophic region VII, in addition to the two extreme behaviors seen in region VI – rigid-like motion and dispersion at low and high values of $\alpha/\beta$, respectively – rotating vortices may be generated for intermediate values of $\alpha/\beta$. Here, the swarming individuals travel close to the characteristic speed $|\mathbf{v}_i|^2 \sim \alpha/\beta$ and for random initial directions of motion, the center of mass is fixed. As shown pictorially in the left hand panel of Fig. 5 vortex size *decreases* dramatically with agent number. In the right hand side of Fig. 5 we also plot vortex area scaling with number of constituents $N$, for various $\alpha$ at fixed $\beta$ in the catastrophic regime. Note that as $N$ increases the area dramatically decreases.

In the case of vortex motion, the dynamics of the center of mass, moving with velocity $V$, can be obtained by summing (2) over all particles $i$ with $|v_i|^2 = \alpha/\beta$. We obtain:

$$m \sum_i \dot{\mathbf{v}}_{\mathbf{i}} = Nm\dot{\mathbf{V}} = \sum_i \nabla_i U(x_i) = 0 \tag{4}$$

where the latter equality arises from the distance dependence of the potential in (3) and the double sum in all pairs $\{i, j\}$. Vortices thus are localized in space, to the contrary of rigid-body structures for which the equality $|v_i|^2 = \alpha/\beta$ does not hold.

Structures such as clumps, vortices and rings generally may rotate counter-clockwise and clockwise, depending on the initial conditions. In the catastrophic

**Fig. 5.** *Left*: Swarms in the catastrophic region VII of Fig. 1. From *bottom* to *top* $N = 100, 200, 300$. Due to the catastrophic nature of the potential, the vortex area decreases dramatically with $N$, and the density increases. The specific potential parameters are chosen as: $C_a = 0.5, C_r = 1, \ell_a = 2, \ell_r = 0.5$. The self-propelling values $\alpha$ and $\beta$ are set as: $\alpha = 1.6$ and $\beta = 0.5$. *Right*: Vortex area as a function of $N$ for various $\alpha$ in log-log scale for the same values of potential parameters. Note the collapsing trend as $N \to \infty$

regime however, under particular choices of the initial conditions, the two rotational directions may coexist, with a portion of the agents going clockwise and the rest counterclockwise. This left and right rotational coexistence is not present in the H-stable regime, where agents keep a fixed distance from each other; in the catastrophic regime, on the other hand, it persists over very large simulation times and can be considered one of its hallmarks.

It is also important to note that the equations of motion giving rise to patterns of the type shown in Fig. 4 treat all individuals in the same way: there is no central commander and the motion of each agent depends only on its position relative to other members of the swarm. Our modeling is thus consistent with the fact that natural swarming occurs in a "democratic" fashion with no leader emerging from the aggregate.

## 5 Variations

In this section we present more complex realizations of the model in (1–2). One interesting scenario is the usage of vehicle swarms of different masses $m_i$. Segregation is likely to occur. For the vortex scenario of region VII, for instance, all agents $i$ rotate at a constant speed $|v_i|^2 = \alpha/\beta$ and experience a mass dependent, centripetal force. This force must be balanced by the mass independent centrifugal term arising from the Morse potential. We obtain:

$$\frac{m_i \alpha}{\beta r_i} = |\nabla_i U(\mathbf{x}_i)| . \tag{5}$$

**Fig. 6.** Catastrophic vortex in the case of variable masses $m_i$, ranging continuously from 1.0 to 3.5. Only four colors are shown for these mass subgroups: from 1.0 to 1.7 blue; from 1.7 to 2.3 green; from 2.3 to 2.9 yellow; from 2.9 to 3.5 red. The red agents concentrated on the outer periphery, with higher masses, circulate with a larger radius. Those in blue are the lightest and describe smaller circles. Similar patterns arise in the case of uniform masses but variable self accelerations $\alpha_i$ or friction $\beta_i$. The interaction parameters are the same as in Fig. 5

In this expression, the first term represents the centripetal force $m_i|vi|^2/r_i$. From (5), it is evident that the quantity $m_i/r_i$ must be constant: agents with larger masses segregate and tend to describe larger radii of motion, as seen in Fig. 6. Similar arguments can be applied to the case of same-mass agents when the $\alpha$ parameter is not chosen uniformly for all vehicles, but from a distribution $\alpha_i$ (or equivalently for $\beta$). In this case $\alpha_i/r_i$ is constant and similarly, agents with larger self-propulsion and faster speeds tend to describe larger radii of motion. For systems where the self-propulsion is programmed to be inversely proportional to the mass of the agent, vehicle segregation is no longer observed, in agreement with (5).

# 6 Continuum Limit

Numerical simulations of the set of coupled, discrete equations of motion 1–3 are time consuming and generally for the systems shown here, one cannot go beyond a few thousand of interacting agents. In this section, we briefly outline the methodology for passing to the continuum limit, where discrete quantities are replaced by continuously varying fields that can offer insight for the collective behavior of large vehicle numbers. The continuum limit might also be useful to study systems for which the exact details of individual constituent motion are not necessary and a macroscopic, statistical picture is of more pertinent interest. We coarse grain the equations of motion by following the

work of Irving and Kirkwood [19] who first derived the macroscopic equations of hydrodynamics for a set of particles in a fluid. The density of the system is defined as:

$$\rho(\mathbf{x}) = \sum_i m_i \langle \delta(\mathbf{x} - \mathbf{x}_i) | f \rangle \tag{6}$$

where brackets signify averaging over many configurations in the phase space defined by the position and momentum of the $i$ particles. The associated phase space probability distribution is denoted by $f$. Similarly, the velocity of the swarm $\mathbf{u}$ can be defined as:

$$\rho(\mathbf{x})\mathbf{u}(\mathbf{x}) = \sum_i m_i \langle \mathbf{v}_i \delta(\mathbf{x} - \mathbf{x}_i) | f \rangle \tag{7}$$

Using these definitions, the dynamics of individual agents, and ensemble averaging, one obtains the following conservation and transport laws [11, 19]:

$$\frac{\partial \rho}{\partial t} + \nabla \cdot (\mathbf{u}\,\rho) = 0 \tag{8}$$

$$\frac{\partial \rho}{\partial t} + (\mathbf{u} \cdot \nabla)\mathbf{u} = \alpha\mathbf{u} - \beta|\mathbf{u}^2|\mathbf{u} - \int \rho(x')\nabla U(\mathbf{x} - \mathbf{x}')dx' \tag{9}$$

where $U(\mathbf{x} - \mathbf{x}')$ is the two body Morse interaction. The derivation of these equations relies on the assumption that agents are not intrinsically correlated, so that the pair density $\rho^{(2)}(\mathbf{x}, \mathbf{x} + \mathbf{x}', t)$ can be written as the product of mass densities $\rho(\mathbf{x}, t)\rho(\mathbf{x} + \mathbf{x}', t)$. Furthermore it is assumed that velocity fluctuations and the energy flux into the system are negligible. In this case, local deviations of each agent velocity from the macroscopic group velocity are considered small. These assumptions are peculiar to the swarming system under consideration, where, for example interaction terms are much stronger than statistical fluctuations. On a more fundamental note, the system described through (1–3) is of a non-conserved type, whereas the standard Irving-Kirkwood derivation of the hydrodynamics equations assumes strictly Hamiltonian type forces. Coarse graining these systems depends on the classical Liouville equation for the time evolution of the position and momentum phase space density function $f$. In the present case, where there are dissipative terms, a strict usage of the Liouville equation is not possible, and a variant is needed.

The continuum model reproduces many features observed in the discrete case, and the phase diagram behavior anticipated in Fig. 2 is recovered. Counterclockwise and clockwise rotating structures are not observed anymore however, due to coarse graining. In particular, steady state solutions for flocking and vortex behavior may be found. In the first case, flocking will arise for $\mathbf{u} = \alpha/\beta\hat{z}$ where $\hat{z}$ is the flock direction and $\rho(r) = \text{constant}$; vortex solutions correspond to $\mathbf{u} = \alpha/\beta(-\sin\theta, \cos\theta)$. In this case, the momentum equation gives an intrinsic formula for the density $\rho(r)$ as follows:

$$\int_0^\infty \rho(\mathbf{r}')\nabla U(\mathbf{r} - \mathbf{r}')d\mathbf{r}' = -\frac{\alpha}{\beta}\frac{\mathbf{r}}{|\mathbf{r}|^2} \ . \tag{10}$$

This equation can be inverted by going to Fourier space and utilizing the convolution theorem. In the flocking case, a linear stability analysis around the uniform solution can be easily performed [20]. In the unstable regime, which corresponds to the catastrophic regime of potential parameters, perturbations will lead to the emergence of multiple, smaller rotating vortices, regularly spaced. We thus expect multiple spirals to arise in the discrete case as well, in the limit of large agent numbers – larger than the numbers used in the discrete systems analyzed here. A linear stability analysis can be performed in the case of vortex solutions as well, although the analysis is more complex. More details will be presented in [20].

# 7 Further Work and Conclusions

Many issues still need to be addressed, such as determining typical swarm scaling as a function of $N$, both in the stable and catastrophic regimes, as well as the time required for pattern formation and whether these times can be cleverly expedited. In addition, when the interaction parameters are chosen to lie between the H-stable to catastrophic transition the time for pattern formation increases and breathing phenomena occur within the emerging structures. It would be interesting to characterize these breathing modes. Furthermore, this study is of a deterministic type, but actual implementations of this model would involve the presence of stochastic terms, due to inevitable random sources of noise. It would be interesting to understand the level of robustness of this multi-vehicular ensemble, both in the stable and catastrophic regimes, as well as its response to individual failure or in the presence of unwanted fields. Delays in communication between agents should be considered as well.

The goal of this work was to present various phases of aggregation for agents interacting through a tunable soft-core potential. Under certain conditions catastrophic behavior arises and in the limit of infinite constituents, the system collapses upon itself. It is easy to avoid such behavior: the insertion of a short-range infinite hard-core repulsion will make any potential stable. Collapsing behavior however, in spite of its name, might be quite useful in designing robot interactions: programming a stable to unstable crossover might lead the robots to change from a dispersive, searching mission to convergence at a specific site.

## Acknowledgments

## References

1. N E Leonard, E Fiorelli (2001) Proceedings of the 40th IEEE International Conference on Decision and Control 2968–2973 IEEE, Orlando
2. J Desai, V Kumar, J Ostrowski (1998) Proceedings of the IEEE International Conference on Robots and Automation 2864–2869
3. D Fox, W Burgard, H Kruppa and S Thrun (2000) Autonomous Robots 8:3
4. E W Justh, P S Krishnaprasad (2003) Proceedings of the 42nd IEEE International Conference on Desicion and Control 3609–3614
5. S Camazine et al (2003) Self organization in biological systems Princeton Univ. Press, Princeton
6. J K Parrish, L Edelstein-Keshet (1999) Science 284:99–100
7. E Bonabeau, M Dorigo, G Theraulaz (1999) Swarm intelligence: from natural to artificial systems Oxford Univ Press, Oxford
8. K Huang (1987) Statistical Mechanics Wiley, New York
9. D Ruelle (1969) Statistical Mechanics, Rigorous results, W A Benjamin Inc, New York
10. A Procacci, Cluster expansion methods in rigorous statistical mechanics (www. mat.ufmg.br/aldo/papers/book.pdf)
11. H Levine, W J Rappel, I Cohen (2000) Phys Rev E 63:017101
12. M R D'Orsogna, Y L Chuang, A L Bertozzi, L S Chayes (2005) http://arxiv.org/abs/cond-mat/0509502
13. B Q Nguyen et al (2005) Proc American Control Conference 1084–1089
14. R Hilborn (2001) Chaos and Nonlinear Dynamics Oxford Univ Press, Oxford
15. F Schweitzer, W Ebeling, B Tilch (2001) Phys Rev E 64:021110
16. V Gazi, K Passino (2002) Proceedings of the 41st IEEE International Conference on Decision and Control 2842–2847 IEEE, Las Vegas
17. V Gazi, K Passino (2003) IEEE Transactions on Automatic Control 48:692–697
18. G H Golub, J M Ortega (1992) Scientific Computing and Differential Equations: An Introduction to Numerical Methods Academic Press, New York
19. J H Irving, J G Kirkwood (1950) J Chem Phys 6:817–829
20. Y L Chuang in preparatino

# Uncertainty Sources in RTD-Fluxgate

B. Andò[1], S. Baglio[1], V. Sacco[1], A. Bulsara[2], and V. In[2]

[1]  DIEES Università degli studi di Catania Catania, Italy
    vincenzo.sacco@diees.unict.it
[2]  Space and Naval Warfare Systems Center San Diego, CA
    bulsara@spawar.navy.mil

**Abstract.** Models and theoretical findings of Residence Times Difference (RTD) Fluxgate have been already presented in previous papers. A very simple sensor structure, negligible onboard power requirements and the intrinsic digital form of the readout signal are the main features of the proposed strategy. In this paper we aim to investigate main sources of uncertainty, including noise, and possible strategies to limit their effects on the devices; finally results on noise characterization are presented.

## 1 Introduction

Fluxgate magnetometers represents the best compromise between SQUIDs and Magnetoresistive devices to sense weak magnetic field with a resolution of 100pT at room temperature. These devices find applicability in many fields; very good examples of past and emerging applications of fluxgate magnetometers can be found in references [1, 2]. These magnetometers are, most often, operated via a readout based on the spectral amplitude of the second harmonic of an applied time-dependent reference signal; the presence of a target dc magnetic flux signal leads to the appearance of peaks at even harmonics of the reference signal frequency, with the peak amplitudes a function of the target signal [3].

Recently, the possibilities offered by new technologies and materials in realizing miniaturized devices with improved performance, have lead to a renewed interest in seeking solutions for how to reduce cost and size of Fluxgate sensors. Recent examples of sensors in PCB [4,5] and CMOS [6,7] can be found in the scientific literature.

In [8, 9] we presented the Residence Times Difference (RTD) Fluxgate magnetometer, a time domain operated Fluxgate that resembles the pulse modulation method [10]. A very simple sensor structure, negligible onboard power requirements and the intrinsic digital form of the readout signal are the main advantages of this strategy, over the conventional readout. However,

lacks of technology to develop valuable ferromagnetic materials and suitable readout electronics limited in the past the device features in terms of noisefloor and sensitivity.

In this paper, we focus on optimizing Fluxgate performances, investigating the main sources of uncertainty including noise, and identifying possible strategies to limit their effects on the device.

## 2 An Overview of the RTD-Fluxgate

### 2.1 Working Principle

An RTD Fluxgate is based on a two-coils structure (a primary coil and a secondary coil) wound around a suitable ferromagnetic core having a hysteretic input-output characteristic. A periodic driving current, $I_e$, is forced in the primary coil and generates a periodic magnetic field, $H_e$. A target field $H_x$ is applied in the same direction of $H_e$; the secondary coil is used as pick-up coil and the output voltage $V_{out}$ is used to detect the target according to the following working principle.

We assume the magnetic core has two commutation thresholds and a two state output, whose behaviour can be described via bistable dynamics governing a double-well potential energy function $U(x)$ [11]. The magnetization $x$ of the magnetic core is governed by the excitation field, $H_e$, produced in the primary coil, and the bistable potential energy function $U(x)$, which underpins the crossing mechanism between the two steady magnetization states of the magnetic core. In order to reverse the core magnetization (from one steady state to the other one), the driving field ($H_e$) must cross the switching thresholds of the magnetic core.

In the case of a time-periodic excitation having amplitude large enough to cause switching between the steady states and in the absence of any target field, the hysteresis loop (or the underlying potential energy function $U(x)$) is symmetric and two identical Residence Times are obtained (see Fig. 1 for representation of potential energy function $U(x)$, its stable states and the relation with hysteresis cycle).

The presence of a target dc signal ($H_x$) leads to a skewing of the loop with a direct effect on the Residence Times, which are no longer the same. Finally, we assume a sharp hysteretic characteristic for the magnetic core; in turn, this allows us to infer that switching between the two stable states of the magnetization occurs instantaneously when the applied magnetic field exceeds the coercive field level $H_c$. Under these simplified conditions, the device operates almost like a static hysteretic nonlinearity, e.g. a Schmitt Trigger. Mathematically, this amounts to the assumption of a very small time-constant $\tau$; in fact, for calculation purposes, it may be assumed that the signal frequency is smaller than $\tau^{-1}$.

**Fig. 1.** Potential energy function $U(x)$, stable states and relation with hysteresis cycle

## 2.2 Technology Issue on RTD-Fluxgate

The technological approach we present is based on the use of metallized FR4 layer 1.6 mm thick. Cost and easy assembling are the most important features of the proposed architecture that have let easy and cost effectiveness development of different testing prototypes.

Devices constructed using a three-layer architecture are presented in Fig. 2. The two outer layers consist of the metallized FR4 layer (very similar to the layer adopted in standard PCB, but thinner) with copper wirings printed on them; the middle layer is a ferromagnetic material. An Amorphous Metal (also known as metallic glass alloy) has been chosen as ferromagnetic core for its suitable hysteretic characteristic. This materials have a non-crystalline structure and possesses, other than useful magnetic properties, interesting physical characteristics that combine strength and hardness with flexibility and toughness and that, even if non relevant magnetic parameters, characterize this prototype realization and allow for considering many interesting applications for this sensor. In particular the Magnetic Alloy (Cobalt-based) 2714A and 2705, by Metglas$^{®}$ [12], has been adopted. Its magnetic characteristic is very sharp and can be suitably approximated via a two state bistable hysteretic description.

There are two sets of printed wirings on each of the PCB's, corresponding to the driving (excitation) coil and the sensing coil. At the end of each wiring set are two small holes that are used for soldering to complete the winding circuitry when the components are put together. The middle layer is a sheet of the ferromagnetic material cut into a specific shape, with outer dimensions varying according to prototypes. The wound core is, then, realized by putting each ferromagnetic core in between the two PCB boards and aligning the holes at the ends of the printed copper wirings of the two boards. A small copper

**Fig. 2.** Two RTD-Fluxgate prototypes and a sample of the magnetic material

wire is passed through the holes from one board to the other and solder is used to fuse the wire in place. The two PCB boards are fused together in this manner to complete the windings for the sensing and excitation coils.

## 2.3 Performance

The set of experiments performed in our laboratories [8,9,13] has showed very good agreement between expected behaviour of RTD-FG and actual behaviour; A number of prototypes have been developed with different geometry.

Table 1 reports the performance obtained for an optimized rectangular planar, PCB Fluxgate using the 2705 Metglas core. Results confirm the suitability of the developed prototype for security applications (metal objects detection) and biomedical ones (Immunoassays techniques) showing features comparable to conventional Fluxgate, with the intrinsic advantages of the strategy proposed (digital output, low power, low cost).

The Operating Range is evaluated as the range of the applied field within the linear region and the Sensitivity value varies from prototypes to prototypes according to the geometry (core form factor).

## 2.4 Readout Circuits

Time-coded nature of the sensor output voltage make it intuitive to operate the RTD-Fluxgate as a digital magnetometer. In this section we aim in fact to

**Table 1.**

| Sensor performance | |
|---|---|
| Operating Range | $20\,\mu T$ |
| Sensitivity | 20000 (s/T) |
| Power Consumption | $1\,mW$ |
| Sensor bandwidth | $600\,Hz$ |
| Resolution | $3\,nT$ |

present a digital strategy to precisely quantify the RTD from sensor output, pointing on low-cost, integrated solutions.

We start the discussion with the general consideration that the pickup coil ($V_{out}$), has a peak to peak value of few mV (depending on the pickup coil geometry, cross-sectional area and number of turn) and can be affected by a considerable amount of noise. In order to increase the signal to noise ratio and adapt the voltage levels to the next stage of the circuit, the first part of the conceived circuit consists of an instrumentation amplifier. Such stadium serves also to adapt the $V_{out}$ offset the next electronic block. The amplifier is in fact followed by a a Schmitt trigger in order to modifies the spike train into a square wave (i.e. a digital signal containing the RTD information). This analog to digital conversion accomplished by the trigger may be also interpreted as an easy way to convert back to the magnetization form the voltage at the pickup coil.

The RTD $= T^+ - T^-$, now contained in the square wave signal indicated as "FGtrigg", is detected using a circuit, hereafter the "RTD Counter", which is based on a digital Up/Down n bit counter (the schematic of such circuit is presented in the text later on); the circuit counts up during $T^+$ and down during $T^-$, hence the decimal representation of the counter output in two's complement:

$$b_1(-2^{n-1}) + b_2(2^{n-2}) + \ldots + b_n(2^0) \tag{1}$$

($n$ number of bit $b_i$ value of the bit at position i) represents the ratio between the RTD and $T_c$.

The two's complement representation has been adopted to allow a positive and negative RTD readout; however we must point out that to maintain a correct sign information we must guarantee that the applied field does not produce an RTD value exceeding $(2^{n-1} - 1)/f_c$, (this in fact would generate a logic 1 for the most significant bit also for positive field) which results in fix the operating range of the sensors to $(2^{n-1} - 1)/(f_c \times S)$ where we have indicated with $S$ the sensitivity of the Fluxgate. Finally the $n$ bit "RTD Counter" output is multiplexed to reduce the number of pin in the design.

The following design rules can be then defined: given the required resolution $r$ and the device sensitivity $S$ the minimum frequency clock is fixed:

$$f_c^{\min} = \frac{1}{S \times r} \tag{2}$$

Concerning the relationship between the operating range O.R. and the number of bit the following expression can be used:

$$\text{O.R.} = \frac{2^{n-1} - 1}{f_c \times S} = (2^{n-1} - 1) \times r \tag{3}$$

this expression allows either to estimate the bit number when the O.R. has been fixed or to estimate the O.R. for a fixed bit number.

As an example, in the case of the prototype with S = 20000 s/T we obtain a minimum clock frequency of 16 kHz.

Considering that increasing the clock frequency shrinks the O.R. (or increase the bit number required) a good compromise between resolution, costs (bit number) and the O.R. must be chosen. As an example, if a 20 kHz clock is used, to cover the operating range of the sensor (20 μT) the number of bit should be set to $n = 14$ bit.

Despite the constraints obtained for the readout electronic design, in the following layout a 8 bit counter operated with a 20 kHz clock will be adopted for the sake of convenience (low cost). The latter choice assures the resolution required and fix the O.R. to 317.5 nT, exploiting advantages deriving from a 8 bit topology at the expense of a tight operating range of the device. It should be stressed that for the considered application resolution and costs are much more valuable than a wide operating range.

Figure 3 shows the block diagram of the circuit implemented with a 5 V technology and $n = 8$ bit.



**Fig. 3.** Whole circuit block diagram implemented with 5 V technology and $n = 8$

The working principle of the "RTD counter" is illustrated in Figs. 4 and 5 via its block diagram and signal flowchart. The system is started via the digital inputs CE (count enable) and RS (reset active low); the first starts the counting operation, the second sets asynchronously the output to 0 if a reset operation is due. Being the CE signal asynchronous respect to the falling edge of the FGtrigg (where the counting operation should start) both the signals are passed through a the JK Flip-Flop: the JK output changes its state at any falling edge of the FGtrigg signal producing the signal CEeff; the rising and falling edge of CEeff, now synchronous with the FGtrigg, can be used respectively to start and stop the counting operation of the Up/Down 8 bit counter (made up by two Up/Down 4 bit counter cascade), we also point that at any counting period follows a standby one where the CEeff is low.

The measuring time is then related to the period of the FGtrigg signal (corresponding to that of the driving field T) according to the relation $T_m = 2 \times T$. Finally the counter output is latched using 8 D Type Flip Flop rising edge sensitive (the CEeff is inverted to update the output at its falling edge).

**Fig. 4.** Schematic of "8-bit RTD counter"



**Fig. 5.** Signal flow of the "8-bit RTD counter

Design and simulation of the circuit have been implemented using the Hit-Kit Tool from Cadence; in order to include the dynamical behaviour of the Fluxgate sensor into the simulation, a model of the device in Verilog-a has been realized.

The tool Analog Artist is used for the transistor level simulation of the instrumentation amplifier and the Schmitt trigger, and for the description level simulation of the whole system. The simulation phase has demonstrated the validity of the working principle; moreover, the maximum clock frequency (for intrinsic resolution estimation) and the instrumentation amplifier static parameters (for compatibility with the sensor characteristic) have been simulated.

The Cells used in the design of the circuit belong to the Austria-microsystems [14] Standard Cell families available for $0.8\,\mu$m double metal CMOS process technologies. The circuit has been realized in the context of a multi-project wafer, that is why in the same die coexists two designs. Figures 6 and 7 shows respectively the layout and a picture of the final device; the circuit in the picture is underline by the white square.



**Fig. 6.** Layout; the circuit is framed in the white square; Die picture; the circuit is framed in the white square



**Fig. 7.** Uncertainty $\Delta$T in the crossing time of a threshold produced by a sinusoidal electric noise affecting the bias fiel

The results of the testing phase are reported in Table 2; the Instrumentation Amplifier Static Gain and Bandwidth have been evaluated using a scope and a function generator.

**Table 2.** Measured circuit performance

| | |
|---|---:|
| Circuit Dimension | $600\,\mu m \times 600\,\mu m$ |
| IA Static Gain | 1000 |
| IA Bandwidth | $3.2\,kHz$ |
| Trigger thresholds | Adjustable ($V_{ee}$) |
| Power Consumption | $70\,mW$ |
| Min Power Supply | 4.7 |
| Max Power Supply | 6 |
| Max operating $f_c$ | $20\,MHz$ |

The total power consumption of the circuit is predominant respect to the sensor itself so for low power application, an optimization of the circuit need to be addressed. The circuit is flexible respect to the power supply fluctuation, and as expected it is possible to adjust the trigger thresholds chancing the pin voltage $V_{ee}$. The value $20\,MHz$ is fully compatible with the noise level of the systems, and according to the above design consideration a frequency clock of $20\,kHz$ is set.

The results shows the possibility to realise hybrid devices exploiting the PCB and the integrated technology for applications requiring high resolution in a tight operating range such as security and traffic monitoring being the performance suitable to detect the presence, or the transit of metal object via their interaction with the geomagnetic field.

# 3 Uncertainty Sources in RTD-Fluxgate

## 3.1 Noisy Sources of Uncertainty

In this section of the work we aim to focus on the uncertainty sources affecting the $H_x$ estimation including noise, identifying possible strategies to limit their effects.

### Electric Noise Affecting the Bias Signal

Assuming that an high frequency noise affects the bias signal, without delve into the details of the calculations, it is possible to shows that the effect of such noise is that to produce an uncertainty in the estimation of the three crossing times of the coercive field. Figure 5 shows graphically the uncertainty $\Delta T$ in the crossing time of a threshold produced by a sinusoidal electric noise

affecting the bias field. It is trivial to shows that $\Delta T$ may be reduced increasing the slope of the bias field in the point of threshold crossing; in principle this may be achieved both by increasing the amplitude or frequency of the driving field, yielding however to a less sensitive devices. (i.e. the device is less sensitive to the noise but also to the signal hence the overall Signal to Noise ratio is not increased.) Lastly, system shielding or bias field filtering is the only suitable approach to reduce the effects of the bias electric noise.

**Electric Noise Affecting the Output Signal**

As discussed earlier in the paper the readout electronics transforms the output signal (via the Schmitt trigger) into a dichotomous signal from where the PIC reads out the RTD; once again an electric noise on the output signal causes an uncertainty in the commutation instants of the Schmitt trigger that may be decreased by increasing the slope of the output signal in the point of thresholds crossing. reduction may be achieved with materials allowing well defined and sharp output spikes. In [13] we showed that in the same driving condition a device adopting the 2705 un-annealed material (as in the prototype investigated in this paper) presented a spikier output compared with a device using the 2714 as Cast or annealed materials.

**Magnetic Noise**

The magnetic noise is produced by small volumes of the ferromagnetic core, more difficult to magnetize than the rest, and so they are not necessarily saturated during each period of the bias field; the uncertainty of the magnetization state of these regions are often associated with structural or surface imperfections of the core; it has been shown previously that improving the structural and surface quality by annealing or polishing reduces noise [1].

**3.2 Non-Noisy Uncertainty Sources**

Any of the noise sources above mentioned, represent uncertainty source for the estimation of the RTD; an other uncertainty source, not expressively referred to a noise source is the readout electronics uncertainty.

**Uncertainty in the RTD Estimation Due to the Readout Electronics**

The estimation of the RTD accomplished by the readout electronics described earlier is affected by the digit uncertainty (that may be decreased by increasing the frequency of the clock) and by the stability of the clock. In principle it is always possible to set the value of the clock to make the digit uncertainty negligible as respect to the other causes.

Lastly it is worth to point out that the estimation of the target magnetic field is made from the estimation of the RTD and the use of a model describing the link between this two physical quantities. The uncertainty introduced by the model must therefore be considered.

**Uncertainty Introduced by the Model**

In [8] we presented a model mapping the $H_x$ value to the RTD quantity in the case of a sinusoidal forcing signal.

$$RTD = \frac{2}{\omega} \left[ \arcsin\left(\frac{H_c + H_x}{H_e}\right) - \arcsin\left(\frac{H_c - H_x}{H_e}\right) \right] \qquad (4)$$

In the small target limit, $(0 - 0.1(\mathrm{A/m}))$, we may assume a linear model, around $H_x = 0$ and the following expression can be obtained for RTD as function of the magnetic field to be sensed [13]:

$$H_x = \frac{\pi f \sqrt{\hat{H}_e^2 - H_c^2}}{4} RTD \qquad (5)$$

However this model must be modified to formally include the demagnetizing effect into the transduction function. The demagnetising factor D multiplied by the relative permeability $\mu_r$ of the core material is approximately equal to the ratio between the target external field $H_x^{ext}$, and the smaller internal field $H_x^{int}$ inside the core material [15]:

$$H_x^{ext} = \frac{\pi f \sqrt{\hat{H}_e^2 - H_c^2}}{4} D\mu_r RTD \qquad (6)$$

where the quantity, $\frac{H_x^{ext}}{D\mu}$ represents the field inside the material. The uncertainty introduced by model (2) contributes to the overall uncertainty in the $H_x$ estimation.

**3.3  Noisefloor Measurements**

The Noisefloor is defined as the square root of the Power Density Spectrum (PDS) at 1 Hz, of the total system noise expressed in T or A/m; it is estimated from a time series (36 seconds) of the magnetometer output with the sensor placed in a three layer Metglas® magnetic shield, with no field applied; the actual device sensitivity around $H_x = 0$ is used to refer the fluctuation of the output (seconds) to the magnetic noise fluctuation (T or A/m). Figure 8 shows the time plot of the moving averaged magnetometer output over an observation window of 2 seconds (shifted by one sample), the p-p noise level was 380 nT and the power density was 530pT/sqrt(Hz)@1 Hz in the frequency range of 0.0278 Hz and 25 Hz. From Figure it is also possible to note the $1/f$ shape of the power density spectrum ($PDS$) below 1 Hz, typically present in all fluxgate sensor types.

**Fig. 8.** Time plot of the averaged magnetometer output over an observation window of 2s (*top*), PDS of the averaged Fluxgate output

## Acknowledgment

## References

1. P. Ripka: Magnetic Sensors and Magnetometers. Artech House, Boston, (2001).
2. F. Kaluza, A. Gruger, H. Gruger: New and future applications of fluxgate sensors. Sensors and Actuators A Vol. 106 (2003) pp. 48–51.
3. F. Primdahl: The fluxgate Mechanism, Part I: The Gating Curves of Parallel and Orthogonal Fluxgates. IEEE Trans. Magn. Vol. 6 (1970) pp. 376–383.
4. A. Tipek, P. Ripka, T. O'Donnell, J. Kubik: PCB technology used in fluxgate sensor construction. Sensors and Actuators A Vol. 115 (2004) pp. 286–292.
5. O. Dezuari, E. Belloy, S.E. Gilbert, M.A.M. Gijs: Printed circuit board integrated fluxgate sensor. Sensors and Actuators A Vol. 81 (2000) pp. 200–203.
6. P. Kejik, L. Chiesi, B. Janossy, R. S. Popovic: A new compact 2D planar fluxgate sensor with amorphous metal core. Sensors and Actuators A Vol. 81 (2000) 180–183.
7. L. Chiesi, P. Kejik, B. Janossy, R. S. Popovic: CMOS planar 2D micro-fluxgate sensor. Sensors & Actuators A Vol. 82 (2000) 174–180.

8. B. Andò, S. Baglio, A. Bulsara, V. Sacco: Theoretical and experimental investigations on residence times difference Fluxgate Magnetometers. Measurements Vol. 38/2, (2005) pp. 89–112.

9. B. Andò, S. Baglio, A. Bulsara, V. Sacco: "Residence Times Difference" Fluxgate Magnetometers. Sensors Journal, IEEE Vol. 5, Issue 5, (2005) pp. 895–904

10. P. Ripka: Review of fluxgate sensors. Sensors and Actuators A Vol. 33 (1992) pp. 129–141.

11. A. Bulsara, C. Seberino, L. Gammaitoni, M. Karlsson, B. Lundqvist, J.W.C. Robinson; Signal detection via residence-time asymmetry in noisy bistable devices. Phys. Rev. E67, 016120 (2003).

12. Honeywell METGLASS Solution, magnetic materials, METGLASS Magnetic Alloy 2114A datasheet, http://www.metglass.com

13. B. Andò, S. Baglio, A. Bulsara, V. Sacco: Effects of driving mode and optimal material selection on RTD-Fluxgate. IEEE Transaction on Instrumentation and Measurement, Vol. 54, issue 4, pp. 1366–1373 (2005).

14. http://www.austriamicrosystems.com/

15. F. Primdahl, B. Hemandox, V. Nielseng, J. R. Petersens: Demagnetising factor and noise in the fluxgate ring-core sensor. J. Phys. E: Sci. Instrum. Vol. 22 (1989) pp. 1004–1008.

# Modeling and Design of Ferrofluidic Sensors

S. Baglio, P. Barrera, N. Savalli, and V. Sacco

DIEES, Università degli studi di Catania, Viale andrea Doria 6, 95125 Catania
`salvatore.baglio@diees.unict.it`

**Abstract.** Novel inertial sensor based on ferrofluids are presented in this paper. The proposed devices have a widely tunable operative range and high sensitivity. A ferrofluidic sample in aqueous suspension acts as inertial mass. The devices are constituted by one excitation coils and one differential sensing coil wound around a glass pipe where the ferrofluid is contained. The bias magnetic force, induced by the coil, attracts the ferrofluid in its centre thus acting like an equivalent spring. The acceleration to be measured reflects therefore in the inertial mass oscillation amplitudes that are sensed by using a differential transformer whose output voltage is a function of the ferrofluid position. Analytical models, simulations and experimental result are presented to demonstrate the suitability of the proposed approach.

## 1 Introduction

Accelerometers with high resolution and wide bandwidth are increasingly demanded in different applications to measure absolute motion, vibration and shock responses [1]. Specifications of wide operating range, high sensitivity and high resolution are strongly fixed by physical and mechanical parameters. Advances in ferrofluids properties investigation, offer opportunities for many novel applications and is giving rise to accelerometer applications, since they have been mainly conceived in the past, in this field, to control damper fluids viscosity [2]. Ferrofluids are synthetic compounds, in either aqueous or nonaqueous solutions, composed by colloidal suspensions of ultra-fine (5–10 nm) single domain magnetic particles [3]. If a magnetic field is applied, the fluidic state is maintained but particles align in the direction of the field and move in a more compliant position, thus causing viscosity variation. Applying suitably high magnetic field magnetic forces can be induced resulting in the entire ferrofluid motion. On the other hand the magnetic force applied to the ferrofluid, for example using a coil, is spatial dependent and function of ferrofluidic amount, applied current and coil configuration. It means that changing only one of the previous parameters magnetic force can be manipulated to act like a variable equivalent spring. Based on the magnetic properties

of ferrofluids and taking into account the magnetic field spatial distribution inside a coil it is possible to conceive accelerometers where the desired resonant frequency, damping coefficient and operative range can be controlled by modulating the applied magnetic force, acting on the ferrofluid. It is obtained a significant flexibility, as tunable operative range, that is "independent" by the structural configuration with respect to traditional accelerometers. The aim of this paper is to propose an approach to high-sensitivity, low-cost, devices for the detection of acceleration in the range of 0–80 μg. Analytical model, simulation and experimental results of the ferrofluidic accelerometers are presented in the following.

## 2 Governing Force and Equations

The analysis of the ferrofluid dynamics in a gradient magnetic field is used to design the accelerometer, evaluating its performance, and developing a method of high-resolution position control. The basic idea concern of detecting external acceleration by using a sensing coil wound on the ferrofluid glass support. Constant current is provided to the primary coil in order to place the particles aggregate in its initial position as depicted in Fig. 1. In order to model the described system behaviour and to evaluate the displacement of the ferrofluidic mass, induced by external acceleration, some assumptions must be made to simplify the computations. In particular, the ferrofluid has been assumed to be viscous and incompressible [4], and magnetic field does not affected by the ferrofluid displacement and/or magnetization [5]. The equilibrium between all forces acting on the particles aggregate must include magnetic and hydrodynamic forces in addition to the inertial force.



**Fig. 1.** Ferrofluidic accelerometer prototype

Taking into account all the acting force the motion of the ferrofluidic sphere can be described by Newton's second law as follow:

$$F = V(\rho_f - \rho_l)g - \mu_0 VM\nabla H - 6\pi\eta Rv \tag{1}$$

where, $V$ is the volume of the ferrofluid sphere, $\rho_f$ and $\rho_l$ are the densities of the particles and the liquid respectively, $g$ is equal to $9.81\,\mathrm{m/s^2}$, $H$ and $M$ are respectively magnetic field and magnetization, $\eta$ is the viscosity of the fluid, $R$ is the radius of the ferrofluidic sphere and $v$ is its velocity. If all the applied forces are considered, while the ferrofluidic accelerometer is assumed to have constant mass, it can be described by a second order system in which the only non linear component derives from the magnetic force (Fig. 2a). If the operating range is restricted to the range in which the magnetic force has a linear spatial dependence the accelerometer model reduces to the classical one (Fig. 2a).

## 3 Simulation and System Setup

Assuming that the ferrofluidic mass is concentrated in the central region of the coil, the steady-state system response to acceleration in the range $0\,\mu\mathrm{g}$–$80\,\mu\mathrm{g}$ (where g is equal to $9.81\,\mathrm{m/s^2}$), with a ferrofluid (suspended in aqueous solution) density equal to $1200\,\mathrm{kg/m^3}$, has been evaluated. As an example in the range $0.015$–$0.03\,\mathrm{m}$ the magnetic force can be assumed linear with proportionally coefficient of $6.5*10^{-7}$. In Fig. 2b the simulated inclinometer steady-state responses (ferrofluid steady-state displacement) to accelerations in the range $0\,\mu\mathrm{g}$–$80\,\mu\mathrm{g}$ is shown. It can be observed that as the acceleration is equal to $50\,\mu\mathrm{g}$ the ferrofluidic sphere reaches the coil limit and goes out the sensor active area (see black curve in Fig. 2b). In order to enlarge the operative range the excitation current is incremented changing the magnetic force slope and so the spring constant value (see red curve in Fig. 2b). The higher magnetic force as a result imposes the ferrofluidic sphere to undergo toward the coil centre and a new equilibrium position is achieved under the action of the external acceleration. A linear relationship between magnetic force, proportional to the square value of the current applied to the primary winding, and the resultant operating range can be observed by simulation results.

## 4 Experimental Results

Two sensing coils are placed, in a differential configuration, inside the primary coil (Fig. 1) as a part of the readout circuit; following the device realization a null peak-to-peak voltage will correspond to the equilibrium position of the ferrofluidic sphere. The primary coil is excited with a signal obtained by

**Fig. 2.** (**a**) Magnetic force along the coil axis, (**b**) Steady-state ferrofluid response (displacement) to accelerations in the range 0–80 µg with different "magnetic" spring constant. Incrementing current and so the spring constant the operative range is enlarged

summing a constant and a pulsating bias current as shown in Fig. 3. The bias signal has the scope of attracting and displacing the ferrofluidic aggregate in the coil centre therefore determining the value of the equivalent spring, the pulsating, having a frequency in the range of kHz and amplitude ten times smaller than the bias current is used to indirectly excited the secondary coils. The choice of the alternate signal frequency (5–10 kHz) has been made in order to have a suitable magnetic coupling between the primary and the secondary windings but such to be far from the single magnetic particles bandwidth therefore not affecting the ferrofluidic mass displacements.

As previously stated, the primary coil generates therefore a magnetic flux that induces equal, but opposite, voltages at the secondary coils, due to their opposite winding sense, therefore the differential output voltage is zero when the ferrofluidic sphere is in the centre between the two secondary windings. The unbalance in the ferrofluid mass position will cause a redistribution of the magnetic flux lines, which will result in a nonzero differential output voltage. The output voltages are then acquired. This approach allows a high sensitivity; because the voltage at the secondary coil, is proportional to the product of the number of turns of the primary and the secondary coils. The read-out circuit is shown in Fig. 3.

Assuming that the ferrofluidic mass (the seismic mass) is concentrated in the central region of the detection coil, the system response to acceleration in the range 0–0.35 g with a ferrofluid (suspended in aqueous solution) having density of about $1200 \, \text{kg/m}^3$ and volume equal to $10^{-7} \, \text{m}^3$ has been measured. As the acceleration is equal to 0.35 g the ferrofluid is lost (the equilibrium position falls outside the coil), so it is necessary to enlarge the sensor operating range. This can be done by increasing the constant current supplied to the container coil, in this way the equivalent spring constant is changed. The experimental result are shown in Fig. 4 a where voltage peak-to-peak variation are reported versus the acceleration in the range 0–0.35 g. The voltage

**Fig. 3.** Schematic of the readout circuit

peak-to-peak variation is a measure of the ferrofluid equilibrium position. It can be observed a linear voltage peak-to-peak increment. As the acceleration is equal to 0.35 g the ferrofluid is loss in each cases so it is necessary to enlarge the sensor operating range by increasing the constant current supplied to the container coil (Fig. 4b), so varying the equivalent spring constant and resonance frequency. To simple estimate the spring constant it is considered the force balance in correspondence of an acceleration equal to 0.35 g. If the magnetic force is considered as an equivalent spring the equilibrium condition is established by the follow equation in which the gravitational force are taken into account.

$$k = \frac{V_g(\rho_0 - \rho_a)a}{x} \tag{2}$$

where $x$ define the limit of the coil. Simulation result reported in previous work [6], shows that this value establish the limit for which ferrofluidic sphere is contained inside the coil.



**Fig. 4.** Variation of the output voltage measured across the sensing coil versus acceleration with an applied voltage (**a**) equal to 1 V, (**b**) equal to 1.6 V

**Fig. 5.** Experimental result: (**a**) 0.05 g Step Response, (**b**) 0.1 g Step Response, (**c**) 0.15 g Step Response, (**d**) 0.2 g Step Response



**Fig. 6.** Resonant accelerometer

## 5 Dynamic Accelerometer Performances

In order to evaluate the accelerometer dynamic performances the step response due to a constant acceleration in the range 0.05–0.2 g are measured. The experimental results are shown in Fig. 5a, b, c, d. Because of the accelerometer can be expressed as a second order system its transfer function can be expressed as follow:

$$H(s) = \frac{w_n^2}{s^2 + 2\xi w_n s + w_n^2} \tag{3}$$

in which $w_n$ and $\zeta$ values, reported in Table 1, are in good agreement with simulation result

**Fig. 7.** Frequency response of the system due to a pulsating sinusoidal magnetic force applied through the driver coil

**Table 1.** Step response experimental results

|  | Values | | | | Units |
|---|---|---|---|---|---|
| Acceleration | **0.2** | **0.15** | **0.1** | **0.05** | m/s$^2$ |
| $T_a$ (Settling Time) | 0.92 | 0.96 | 0.92 | 0.88 | s |
| $T_s$ (Rise Time) | 0.2 | 0.2 | 0.2 | 0.18 | s |
| $S\%$ | 13.3 | 6.6 | 5.97 | 3.17 | – |
| $\omega_n$ | 6.03 | 2.73 | 2.84 | 2.85 | rad/s |

# 6 Resonant Accelerometer

Two separately excited coils, wound on the glass pipe, are used to induce ferrofluid motions (Fig. 6). Like the ferrofluidic inclinometer constant current is provided to the primary coil in order to attract and displace the particles in the coil centre as a sphere and the two sensing coils are placed in a differential configuration across the equilibrium position of the ferrofluid. A sinusoidal current with a frequency equal to 300 mHz and a voltage peak-to-peak equal to 2.5 V is provided to the driver coils in phase opposition to promote ferrofluid oscillation (see Fig. 6) with this choice the ferrofluid is in the resonance range so sensitivity is enhanced (see Fig. 7).

Assuming that the ferrofluidic mass is concentrated in the central region of the sensing coils, the system response to acceleration in the range [0 g– 0.35 g] with a ferrofluid (suspended in aqueous solution) having density equal to 1200 kg/m$^3$, has been measured. The results are shown in Fig. 8 where peak-to-peak voltage (black line) and offset (red line) are reported versus the applied acceleration. The offset variation is a measure of the ferrofluid equilibrium position whereas the amplitude of its oscillation is revealed by the volage peak to peak variation.

**Fig. 8.** Preliminary experimental result of the symmetric resonant accelerometer

As it can be visualized a more linear and symmetric response is obtained and a sensitivity three time higher than inclinometer. The ferrofluidic accelerometers here can be considered a preliminary step toward the realization of such ferrofluidic accelerometers sensors, the principle of operation of the ferrofluidic accelerometers macro-prototype is demonstrated.

## 7 Conclusion

Different ferrofluidic accelerometers configuration, based on the use of magnetic ferrofluid properties and magnetic field spatial dependence, for the detection of 0–0.5 g acceleration, has been described here. Analytical models have been derived for the sensor with respect to acceleration and magnetic force changes. Ferrofluid displacement has been simulated for the detection of $0\,\mu g$–$80\,\mu g$ acceleration and tested in the range of 0 g–0.5 g acceleration. A suitable signal conditioning circuit has been also realized to reveal the detected acceleration as electric output. The proposed approach has been firstly validated by means of macroscopic sensor prototypes realized with suitable coil wound around a glass pipe. Inclinometer and resonant ferrofluidic accelerometer have been tested to explore ferrofluidic accelerometers performance. The performance of the devices presented here encourages further efforts for the development of low-cost, high-sensitivity, simple designed, novels accelerometers.

## References

1. N. Yazdi, F. Ayazi, and K. Najafi, "Micromachined Inertial Sensors", Proceedings of the IEEE, vol. 86, pp. 1640–1659, August 1998.

2. G.Q. Hu, W.H. Liao, "A feasibility study of a microaccelerometer with magnetorheological fluids", Proceedings of the 2003 IEEE International Conference on Robotics, Intelligent Systems and Signal Processing, Changsha, China, pp. 825–830, October 2003.
3. J. Popplewell and R.E. Rosensweig, "Magnetorheological fluid composites", J. Phys. D, Appl. Phys., vol. 29, pp. 2297–2303, January 1996.
4. A. Hatch, A.E. Kamholz, G. Holman, P. Yager, and K.F. Böhringer, "A ferrofluidic magnetic micropump", Journal of Microelectromechanical System, vol. 10, pp. 215–221, June 2001.
5. M.A.M. Gijs, "Magnetic bead handling on chip: new opportunities for analytical application", Microfluid Nanofluid, vol. 1, pp. 22–40, October 2004.
6. S. Baglio, P. Barrera, P. Liseo, N. Savalli, "Ferrofluidic accelerometers", proceeding to Eurosensor XIX, 11–14 September 2005, Spain.

# Thermocromic Materials for Temperature Sensors in New Applications

A. Boscolo, E. Menosso, B. Piuzzi, and M. Toppano

Artificial Perception Laboratories, DEEI Università degli Studi di Trieste, Via A. Valerio, 10 34127 Trieste, Italy

**Abstract.** The development of a temperature sensor with a non-linear behaviour is described in this work. The device acts like a thermal switch with hysteresis and is able to return to the initial conditions. Its on-off settings are the passage on fixed temperature that corresponds to the critical temperature (Tc) of thermocromic materials, in this case vanadium dioxide. The final aim is to modify the device switching temperature by changing the thermocromic material critical temperature, following the needs of the specific application.

A critical temperature variation can be induced through a doping process with elements having a great vanadium affinity. For this application molybdenum, tungsten, niobium have been employed, and experimental samples have been realized by thin films deposition. A model of Tc variation and doping concentration has been obtained for $VO_2$-Mo film, which good results have been obtained for, to support the programmable device development.

## 1 Introduction

This paper presents the development of a switching sensor based on vanadium dioxide. $VO_2$ is a thermocromic material, which is a material having a behaviour transition depending on temperature. When temperature reaches a value characteristic of thermocromic material, called "critical temperature" (Tc), there are changes in optical and electrical characteristics. For $VO_2$, the critical temperature is $68°C$. In order to change Tc, a study of chemical and structural properties showed that the most interesting doping materials are molybdenum, tungsten and niobium [1]. In particular, obtained results on the $VO_2$-Mo films give an interesting relationship between the used dopant quantity and the critical temperature variation. A Tc setting could be interesting for many applications, where the $VO_2$ critical temperature is too high. For example, the choice of on-off temperature could be useful in safety systems productions for countries with different legal requirements on fire alarm value. An advantage of the proposed approach is that this could be done keeping on

the same productive processes of sensor, only setting the doping process of thermocromic materials in order to adjust critical temperature.

## 2 Materials

Thermocromic properties are due to a semiconductor-metal transition of the material behaviour. Some of these materials are vanadium oxide ($VO_2$, $V_2O_3$) and iron oxide ($Fe_3O_4$). In particular vanadium dioxide has a monocline structure (with a centred face) and a semiconductor electrical behaviour, for temperatures below Tc. For higher temperatures, it has a tetragonal structure and a metal behaviour.

The electrical resistance trend is characterized by a hysteresis that presents, on a temperature range of some tenth of degrees, a resistance variation of about two or three magnitude orders (see Fig. 1).



**Fig. 1.** Sheet resistance versus temperature of RF sputtering produced $VO_2$

In the $VO_2$ doping process, dopant concentrations higher than these typical of usual electronic processes were employed.

The described temperature sensor, and consequently the developed material, needs a sheet resistance variation as higher as possible. In this way the advantage to minimize the use of complicated conditioning systems to detect temperature variation could be maintained. The hysteresis loop width is another important aspect to look at, for its direct action on sensitivity of the device. When a steep sheet resistance occurs, the range of temperature around the critical temperature is really thin. Devices will be characterized in this case by a good sensitivity.

# 3 Processes

There are many usually employed methods in vanadium oxide production. Sol-gel method requires alkoxides employment [2]. Reactive sputtering [3, 4] that could be DC-, RF- or magnetron-sputtering [5], is used in metallic vanadium deposition and is performed in an argon plus oxygen atmosphere. PVD (Physical Vapour Deposition) techniques [6] are also used for metallic vanadium deposition, and are followed by a thermal oxidation. For this application, vanadium dioxide and doped vanadium dioxide have been obtained by RF-sputtering of metallic vanadium, co-sputtering of vanadium plus doping materials and a following oxidation and thermal annealing. The choice of thin film technology permits to control and minimize structural damages due to thermal stress at critical temperature and mechanical stress due to semiconductor–metal transition. Furthermore separation between thin film deposition and thermal oxidation processes brings to a better productive process control and flexibility.

The semiconductor-metal transition occurring at the critical temperature is strictly linked with a $VO_2$ film structural modification. So it has to be paid attention, in the process tuning, to the substrate material choice and to structural modifications caused by the working process itself. The deposition and the thermal oxidation process in $VO_2$ production have to be worked out in an appropriate manner, to realize a really well-structured film [7]. For this reason, substrate must have non-amorphous structure [8] and the warming and cooling times of the thermal oxidation have to be quite long. The same problems exist in doping processes, and so the same working precautions have to be maintained also in them. Furthermore, in this last case, elements chemical affinity and behaviour in oxidation processes have to be considered.

A Leybold 550 VZK with a 3 kW RF generator has been used for the vanadium depositions. For this application, a $VO_2$ film has been deposited by means of RF sputtering with vanadium cathodes. For the doped vanadium film, co-sputtering process has been used. Two cathodes, the first a vanadium one and the second a dopant one, have been activated. In the vacuum chamber, the substrate was heated to clean it and to have a better film adhesion [9]. Furthermore, in order to obtain the desired structure, a silicon wafer was used. After the deposition, the oxidation and annealing processes were done in atmosphere, at a 500°C temperature, for 15 minutes.

The characterisation system has been realized with virtual instrumentation technology. The system acquires the sheet resistance and the temperature values at the same time. The temperature was continuously varied. First, in a heating phase, the temperature was increased slightly over the critical temperature, then in a cooling phase the temperature was lowered to the starting value. SEM with a microanalysis head with a $4.58*10^{-3}$ magnification was employed to verify the dopant level and the film homogeneity, obtaining a quantitative profile of elements in the film.

**Fig. 2.** Hysteresis loops comparison for thermocromic materials obtained by doping VO$_2$ processes

## 4 Results

Better results in the VO$_2$ film realization have been obtained using a metallic vanadium film with a sheet resistance of about 8 ohm/square. After the oxidation and the annealing processes, the sheet resistance presented the desired hysteresis loop, about tree orders of magnitude, in a range of ten degrees around a temperature value of 68°C.

Dopant materials nature and concentration have been chosen in order to modify the critical temperature without producing damages in the film structure. Therefore, chemical elements employed in the experiments have chemical affinity with vanadium. Their concentration is typical of alloys instead of a usual electronic doping process. Furthermore, these alloys have only a phase ($\beta$ phase) [10], by whatever concentration. Consequently, a similar behaviour could be waited for the three considered alloys, but thermal oxidation introduces another variable in the experimental setting. The material behaviour depends also on different oxidation evolution.

The better results have been obtained with VO$_2$-Mo alloy, that presents a modified Tc and an acceptable degradation of hysteresis behaviour compared with hysteresis of starting material (see Fig. 2). A model of sheet resistance versus film temperature and doping concentration has been obtained, and is reported in Fig. 3.

For the niobium doped vanadium there is a link between the critical temperature modification and the concentration, but hysteresis loop has been really worsened by the production process. The sheet resistance magnitude variation is less than one order and the temperature range is too width (see

**Vanadium Dioxide-Molybdenum**



**Fig. 3.** Relationship between Tc variation and doping concentration in VO$_2$-Mo film

Fig. 2). At last, tungsten [11] has had no noticeable results, for the differences in the oxidation mechanisms of vanadium and tungsten.

## 5 Conclusions

The research activity aims to evaluate the effectiveness of the doping in the VO$_2$ production to realize a temperature sensor with a settable working temperature. The bests results have been achieved with molybdenum; with this material a good relationship between Tc and concentration of doping material has been obtained. The VO$_2$-Mo shows a reduction of Tc value and its hysteresis presents a good change in the value of resistance versus a limited temperature range. In conclusion, it is a thermochromic material with settable Tc dependent on the dopant concentration, exploitable to develop a temperature sensor with working temperature centered on its Tc and useful to achieves the desired target in accuracy and sensitivity for a final temperature sensor.

## References

1. David R. Lide, editor. Handbook of chemistry and physics. CRC Press, 2001.
2. Songwei Lu, Lisong Hou, and Fuxi Gan. Surface analysis and phase transition of gel derived vo2 thin films. Thin solid films, (353):40–44, June 1999.
3. R.O. Dillon, K. Le, and N. Ianno. Thermochromic vo2 sputtering by control of a vanadium-oxygen emission ratio. Thin solid films, (398-399):10–16, 2001.

4. H. Miyazaki, F. Utsuno, Y. Shigesato, and I. Yasui. The structural character-istics of vox films prepared by he-introduced reactive rf unbalanced magnetron sputtering. Thin solid films, (281-282):436–440, 1996.

5. Jingzhong Cui, Daoan Da, and Wanshun Jiang. Structure characterization of vanadium oxide thin film prepared by magnetron sputtering methods. Applied Surface science, (133):225–229, 1998.

6. E. Lungscheider et al. Surface and coatings technology, (142-144):137–142, 2001.

7. Feliks Chudnovskiy, Serge Luryi, and Boris Spivak. Switching device based on first-order metal-insulator transition induced by external electric field. 2002.

8. Y. Muraoka and Z. Hiroi. Metal-insulator transition of vo2 thin films grown on tio2(001) and (110) substrates. Applied physics letters, 80(4):583–585, January 2002.

9. Aicha Elshabini-Riad and Fred D. Barlow III. Thin film technology handbook. Mc Graw Hill, 1998.

10. M.V. Sofin et al. Determination of phase equilibria in the ni-v-nb-ta-cr-mo-w system at 1375 k using the graph method. Journal of alloys and compounds, (321):102–131, 2001.

11. W. Burkhardt, T. Christmann, B.K. Meyer, W. Niesser, D. Schlach, and A. Scharmann. W-and f-doped vo2 films studied by photoelectron spectrometry. Thin solid films, (345):229–235, 1999.

# A SQUID Ring-Resonator Finate State Machine

P.B. Stiffell[1], M.J. Everitt[2,1], T.D. Clark[1], A.R. Bulsara[3], and J.F. Ralph[4]

[1] Centre for Physical Electronics and Quantum Technology, School of Science and Technology, University of Sussex, Falmer, Brighton, BN1 9QT, UK
`p.b.stiffell@physics.org`
[2] The British University in Egypt, Sheraton Heliopolis, Cairo, Egypt.
[3] Space and Naval Warfare Systems Center, US Naval Command, Code 2363, 53560 Hull Street, San Diego, California 92152-5001, USA
[4] Department of Electrical and Electronic Engineering, Liverpool University, Brownlow Hill, Liverpool, L69 3GJ, UK

## 1 Introduction

In this paper we consider a system comprising a highly hysteretic, single Josephson weak link SQUID ring (i.e. where $2\pi L_s I_c/\Phi_0 \gg 1$ for a SQUID ring inductance $L_s$ and weak link critical current $I_c$ with $\Phi_0 = h/2e$) coupled to an finite quality factor $LC$ parallel resonator (tank circuit), the schematic of which is shown in Fig. 1. The non-linear dynamics of this system has been the subject of a great deal of interest both for the quantum and classical operational regimes (see for example [1–3]). In this paper we confine ourselves to the semi-classical Resistively Shunted Junction plus Capacitance (RSJ+C) model of the SQUID ring [4]. For such a SQUID ring-tank circuit system, where the tank circuit is driven by a current $I_{\text{in}}$, the equations of motion are



**Fig. 1.** Schematic of a SQUID ring inductively coupled to a tank circuit

$$C_t \frac{\partial^2 \Phi_t}{\partial t^2} + \frac{1}{R_t} \frac{\partial \Phi_t}{\partial t} + \frac{\Phi_t}{L_t(1-k^2)} = I_{\text{in}}(t) + \frac{\mu \Phi_s}{L_s(1-k^2)} \tag{1}$$

$$C_s \frac{\partial^2 \Phi_s}{\partial t^2} + \frac{1}{R_s} \frac{\partial \Phi_s}{\partial t} + I_c \sin\left(\frac{2\pi \Phi_s}{\Phi_0}\right) + \frac{\Phi_s}{L_s(1-k^2)} = \frac{\mu \Phi_t}{L_s(1-k^2)} \tag{2}$$

where $C$ is capacitance, $L$ inductance, $R$ resistance and $\Phi$ is flux and the subscripts $t$ and $s$ refer, respectively, to the tank circuit and SQUID ring. Here the fraction of the flux that is coupled between the ring and the tank circuit is $\mu$ and $k = \mu\sqrt{L_t/Ls}$ quantifies the strength of the inductive coupling. We solve numerically this pair of coupled differential equations for the situation in which the tank circuit has a drive current of the form $I_{\text{in}} = A\sin(\omega_d t)$. We have found that for reasonable circuit parameters, and a suitable choice of the drive amplitude $A$, there exist several stable solutions to these equations that are separated evenly in the tank circuit voltage. We label this finite set of states $\alpha, \beta, \gamma, \ldots$ corresponding to increasing values of voltage.

We show that by adding any of a number of suitable voltage control pulses $V_p$ to the SQUID ring we can switch between any of the available levels in a deterministic way. Experimentally we have observed up to nineteen of these levels [1]. Hence, with this system not only do we have a finite alphabet $\Sigma = \{\alpha, \beta, \gamma, \ldots\}$ but we are also equipped with a set of state transition functions $\{V_p\}$ that enable us to move cyclically between the elements of $\Sigma$. These are more than sufficient to create a simple finite state machine. Moreover, here we can choose from a number of different pulses $V_p$ that will map the input alphabet onto itself and, as we can use inverted pulses to access the previous state, we also have an inverse.

Fabrication of thin film devices based on these ideas could very well lead to the development of finite-state machines and other multi-level logic circuits that might, in principal, be able to operate at very high drive frequencies and with concomitant speed.

## 2 Background

Experimental work done at Sussex [1] showed that a SQUID ring weakly coupled to a tank circuit resonator could exhibit a series of levels in the voltage-current dynamic of the resonator. After ensuring that this levels behaviour was not an experimental error the group began to investigate how these levels were caused. When some progress was made in identifying how these levels were caused we then began to consider a different question.

The question we wanted to answer was if we could find a mechanism to allow us to control the passage between these voltage levels. This idea actually centred around two important issues. Firstly, could we prevent the random motion between levels and secondly could we find a trigger that would allow us to make a stable transition from one level to another. The first of these two issues proved simple enough to deal with. Work done on similar systems [5,6]

**Fig. 2.** The voltage response of the tank circuit resonator against dimentionless time when an appropriate voltage pulse is applied to the SQUID ring. The circuit parameters used for this simulation were $C_s = 10^{-13}$ F, $L_s = 6 \times 10^{-10}$ H, $R_s = 10\,\Omega$, $I_c \approx 73\,\mu$A, $C_t \approx 7.6 \times 10^{-10}$ F, $L_t = 6.3 \times 10^{-8}$ H, $R_t \approx 4550\,\Omega$, a drive applitude of 222nA and a coupling strength $\mu \approx 0.0087$

showed no signs of this effect, so it seemed reasonable to assume that with sufficient experimental preperation the noise levels in our system could be reduced to the point where unwanted levels behaviour could be discounted.

Confident that we could prepare a system stable enough to eliminate random levels behaviour we then began to investigate if we could find any method by which we could trigger a voltage level transition. Further work [7] we did on this lead us to find that voltage pulses applied to the SQUID ring could indeed induce voltage level changes in the tank circuit. The work done in [7] on the stochastic origins of this behaviour also indicated that we were correct in our confidence that we could experimentally prepare a system free of stochastic levels behaviour.

## 3 Results

The most important result of this work was the discovery that an appropriately timed voltage pulse applied to a SQUID ring could be used to induce a change in the, otherwise stable, output voltage of a coupled resonator circuit. The effect of this kind of pulse is shown in the time evolution of the output voltage shown in Fig. 2. The fact that the output status of the resonator could be alter by the use of a simple pulse input meant that we had found a way to progress the output of our system through stable and definable changes with the use of a trigger function. The next important result to investigate was to ensure that this process was able to be continued throughout the voltage levels. To verify this we show, in Fig. 3 a series of voltage step changes from the

**Fig. 3.** The voltage response of the tank circuit resonator against dimentionless time when a series of voltage pulses is applied to the SQUID ring. The circuit parameters used for this simulation are the same as those used in Fig. 2

first stable voltage level progressing, sequencially, through the other 3 voltage levels in the system.

Having found we can step our system through a series of output states by the application of voltage pulses we saw that we could satisfy the basic criteria for a finate-state machine. However, we felt it was important to see if more complex sequences of output state could be achieved. To ensure this we felt we had to varify that our system could perform two more types of change. The first of these was the ability to move backwards through the output states. We found that this was, as we expected, possible by applying a negative voltage pulse to the system. The effect of this is shown in Fig. 4.

The last type of transition effect we wanted to verify was possible was to see if transitions could be made between non-nearest states with a single pulse. Altering the size of pulses used, within certain limits, did indeed allow us to activate direct transitions between non-neighbouring output states of our system. In fact, we show in Fig. 5 that we could access all the output states we found in this system from the initial stable state with a single pulse.

Thus far in our investigation we have used a system with four output states. This was an arbitrary fact that was due to the fact that the system parameters chosen for this model were used to give a fit to previous experimental work [1]. However, for a finate-state machine it is better to be able to have a variable number of output states. This then allows us to have a number of output states that equals the number of steps we want for our finite-state machine. To ensure this was possible, in theory, we altered the system parameters to try and find a configuration with more output states available. We found that a modification of the SQUID resistance allowed us to alter the number of output states we have available to us in our system. A series of output state transitions for a system with eleven output states is shown in Fig. 6.

**Fig. 4.** The voltage response of the tank circuit resonator against dimensionless time when a series of positive voltage pulses is followed by a negative pulse. The circuit parameters used for this simulation are the same as those used in Fig. 2



**Fig. 5.** Three seperate voltage responses of the tank circuit resonator against dimentionless time when the voltage pulses applied to the SQUID ring is of different applitudes. The circuit parameters used for this simulation are the same as those used in Fig. 2

## 4 Conclusion

We found that we could, at least in theory, define a system with arbitrary number of output states that can be traversed in any manner we choose, by the application of appropriately sized and timed voltage pulses. This means that this system can definately be looked upon as a viable engine for a finite-state machine. The system has an alphabet of output states and a deterministic control procedure for moving within that alphabet. Within reason, we can

**Fig. 6.** The voltage response of the tank circuit resonator against dimensionless time when an appropriate voltage pulse is applied to the SQUID ring. For this simulation the SQUID resistance $R_s$ is $45\,\Omega$, which causes the system to have a much greater number of potential voltage output levels

define a system with any number of output states such that it can be used to control any size of finate-state responce we could require of it.

At present the only issues with this principle is that the timing required on the pulses is very precise, and also dependent on the presence of a psuedo-random tansition between wells in the SQUID ring. This is due to the fact that the system is only responsive to the voltage pulses while the SQUID flux is exhibitting a transition between wells in the SQUID potential. The timing of these transitions is psuedo-random, but they can be made to occur with greater or lesser frequency by the use of appropriate drive currents in the resonator. For a more detailed analysis of these transitions is covered one of our recent papers about to be published [7]. It may also be possible to solve the problem of such delicate timing for the voltage pulses by instead using a voltage ramping that runs throughout the duration of one of these well transitions. However this will need to be addressed in further research on this topic, as we have not yet had sufficient time to verify this concepts thus far.

It is interesting to note that as well as being able to be used as the control status of a finite-state machine we also believe that these multi-level systems could be used as a way of fabricating multi-level logic components. If this can be done it would allow these systems to form the backbone of a new form of computing that could lead to significant speed increases in processing.

## Acknowledgements

# References

1. R.J. Prance, R. Whiteman, T.D. Clark, H. Prance, V. Schollmann, J.F. Ralph, S. Al-Khawaja, and M.J. Everitt. Nonlinear multilevel dynamics of a coupled squid ring-resonator system in the hysteretic regime. *Phys. Rev. Lett.*, 82:5401–5404, 1999.
2. A.R. Bulsara, K. Wiesenfeld, and M.E. Inchiosa. Nonlinear dynamics in a high-gain amplifier: the dc squid. *Ann Phys-Berlin*, 9:679–688, 2000.
3. J. Diggins, J.F. Ralph, T.P. Spiller, T.D. Clark, H. Prance, and R.J. Prance. Chaotic dynamics in the rf superconducting quantum-interference-device magnetometer – a coupled quantum-classical system. *Phys. Rev. E.*, 49:1854–1859, 1994.
4. K. Likharev. *Dynamics of Josephson Junctions and Circuits.* Taylor & Francis, 1986.
5. J.R. Friedman, V. Patel, W. Chen, S.K. Tolpygo, and J.E. Lukens. Quantum superposition of distinct macroscopic states. *Nature*, 406:43–46, 2000.
6. I. Chiorescu, P. Bertet, K. Semba, Y. Nakamura, C.J.P.M. Harmans, and J.E. Mooij. Coherent dynamics of a flux qubit coupled to a harmonic oscillator. *Nature*, 431:159–162, 2004.
7. P.B. Stiffell, M.J. Everitt, T.D. Clark, C.J. Harland, and J.F. Ralph. Control of multilevel voltage states in a hysteretic superconducting-quantum-interference-device ring-resonator system. *Phys. Rev. E.*, in publication, 2005.

# Invited Papers

# Suprathreshold Stochastic Resonance Mediated by Multiplicative Noise

N.G. Stocks, A. Nikitin, and R.P. Morse

School of Engineering, University of Warwick, Coventry CV4 7AL, UK

**Abstract.** We have investigated information transmission in an array of threshold units. Each unit receives a common input signal but independent multiplicative noise. We demonstrate a phenomenon similar to stochastic resonance and suprathreshold stochastic resonance and show that information transmission can be enhanced by a non-zero multiplicative noise level. Given that sensory neurons in the nervous system have multiplicative as well as additive noise sources, and they act approximately like threshold units, our results suggest that multiplicative noise might be beneficial in increasing information transmission in neural populations.

## 1 Introduction

In recent years there has been a significant increase in interest in the role of noise in nonlinear signal processing. This activity has largely been motivated by studies of stochastic resonance (SR) [1,2] but also by a desire to understand stochastic aspects of neural coding [3–8]. In particular, the study of signal coding in parallel arrays (populations) of nonlinear devices (neurons) has received considerable attention, in such arrays a new form of SR – termed suprathreshold stochastic resonance (SSR) – has been discovered [9,10]. In a similar vain to SR, SSR can lead to an improvement in information transmission when internal noise is added to the system. However, SSR, which can only occur in arrays of nonlinear devices, has a number of advantages over conventional SR. First, it occurs for all signal levels – it does not require that the signal be subthreshold. Consequently SSR can be used to improve information transmission for a broader class of signal than standard SR. Second, it provides an optimal method of enhancing information when the signal to be detected is comparable (or smaller) than the residual internal noise [11,12]. In contrast, for SR in a single device, greater information flow is usually obtained by simply increasing the input signal (or lowering the threshold if possible) rather then setting the signal to be subthreshold and utilising SR. For these reasons, the potential exploitation of SSR in technological applications is arguably greater

than for conventional SR. Possible applications are novel digital-to-analogue converters [9, 10, 13–15] and sonar arrays [10, 16].

Although SSR has now been studied in a wide variety of different contexts [6, 8–15, 17–21] all these studies have been undertaken assuming that the noise is additive to the signal. However, it is well established that in neural systems the noise may enter in a signal dependent or multiplicative fashion [24, 25]. Signal dependent noise, for example, is characteristic of the propagation of nerve impulses between nerve cells through the quantal release of neurotransmitter at synapses; Furukawa et al. [26] have shown experimentally that the variation of the neurotransmitter release into a synaptic cleft is proportional to the intensity of the stimulus.

Whereas additive noise, via the SSR effect, is potentially beneficial for signal coding in neural populations the role of signal dependent or multiplicative noise has not previously been studied. This was the motivation for the current study, in which we investigated the effect of multiplicative noise in an array of identical threshold units, each subject to a common input signal and each including an independent source of multiplicative noise.

## 2 Models

The model we study is loosely based on that proposed by Furukawa [26] for the modelling of synaptic transmission in hair cell transduction. However, the model, has been simplified to enable theoretical analysis to be undertaken and to make connection with other studies in SSR.

We have modelled an array of $N$ nerve fibres by an array of $N$ simple threshold units (level-crossing detectors). Each threshold unit was subject to the same input signal $x$, which for generality was transformed by the linear or non-linear function $F(x)$, but independent multiplicative noise. The multiplicative noise can be considered to be a simplified abstraction of the quantal release of neurotransmitter in the model of Furukawa. To be consistent with other studies on SSR the output, $y$, of each unit was given by the Heaviside function,

$$y_i = \begin{cases} 1 : v_i \geq U_i \ , \\ 0 : v_i < U_i \ , \end{cases} \tag{1}$$

where $U_i$ is the threshold of the unit and $i = 1 \ldots N$, and

$$v_i = DF(x)\eta_i \ . \tag{2}$$

Here, $\eta_i$ had a Gaussian noise distribution with zero average mean and unit dispersion and $D$ is the coefficient of the proportionality (the "noise intensity"). Through an appropriate choice of threshold settings and noise intensities, the array of threshold units can model a number of applications. For simplicity, however, we considered only the case where the thresholds were

identical $(U_1 = U_2 = \cdots = U)$ and all units have the same noise intensity. Moveover, the results were for a Gaussian signal, such that the probability density for the input signal was given by

$$P_x(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{x^2}{2\sigma^2}\right) , \tag{3}$$

where $\sigma$ is the intensity of the signal.

Because the noises $\eta_i$ were mutually independent, i.e. $\langle\eta_i\eta_j\rangle = 0$ if $i \neq j$, the variables $v_i$ were also mutually independent. The probability that exactly $j$ units were in the state 1 for a given value of $x$ was therefore given by the binomial distribution

$$Prob\left\{\sum_{i=1}^{N} y_i = j \middle| x\right\} = C_j^N q_x^j (1 - q_x)^{N-j}, \qquad C_j^N = \frac{N!}{j!(N-j)!} . \tag{4}$$

Given that the noise intensity for each unit was identical, $q_x$ was identical for all units too.

The overall output of the system was taken to be the sum of the outputs of the individual units. The system can therefore be taken to map the instantaneous amplitude of the signal, $x$, into $j$, the number of the units in the state 1 such that $j = \sum_{i=1}^{N} y_i$.

## 2.1 Simplified Model

We analyzed two versions of the model: one with a linear representation of the signal, i.e. $F(x) = x$ and one with a half-wave rectified signal. The rational for the rectified model is to mimic the half-wave rectification that is known to occur in hair cell transduction.

In the linear signal model, we took $F(x) = x$ and

$$v_i = Dx\eta_i . \tag{5}$$

With this condition, the conditional probability desnity of $v_i$ was Gaussian

$$P_{v|x}(v_i|x) = \frac{1}{\sqrt{2\pi D^2 x^2}} \exp\left(-\frac{v_i^2}{2D^2 x^2}\right) , \tag{6}$$

with an ensemble average, $\langle v_i \rangle$, of 0 and a standard deviation of $D^2 x^2$. The probability that $y_i = 1$ for a given signal level $x$ can be calculated as

$$q_x = Prob\{y_i = 1|x\} = Prob\{v_i > U|x\} = \int_U^\infty P_{v|x}(v_i|x)dv_i$$

$$= \frac{1}{2} - \frac{1}{2}erf\left(\frac{U}{\sqrt{2D^2 x^2}}\right) . \tag{7}$$

## 2.2 Model with Rectification

A more realistic model of hair cell transduction/synaptic transmission is one with rectification. It can be described by

$$v_i = DR(x)\eta_i \; , \tag{8}$$

where $R(x)$ is the rectification function,

$$R(x) = \begin{cases} x : x > 0 \; , \\ 0 : x \le 0 \; , \end{cases} \tag{9}$$

The variable $v_i$ is not Gaussian for this model over all values of $x$. When $x > 0$, the conditional probability density is

$$P_{v|x}(v_i|x) = \frac{1}{\sqrt{2\pi D^2 x^2}} \exp\left(-\frac{v_i^2}{2D^2 x^2}\right) \; , \tag{10}$$

but when $x \le 0$, it is

$$P_{v|x}(v_i|x) = \frac{1}{2}\delta(v_i) \; , \tag{11}$$

where $\delta(v_i)$ is the delta-function. The probability that $y_i = 1$ for a given positive value of $x$, is therefore

$$q_x = Prob\{y_i = 1|x > 0\} = Prob\{v_i \ge U|x > 0\} = \int_U^\infty P_{v|x}(v_i|x)dv_i$$

$$= \frac{1}{2} - \frac{1}{2}erf\left(\frac{U}{\sqrt{2Dx}}\right) \; , \tag{12}$$

and for negative $x$, $x \le 0$, it is

$$q_x = Prob\{y_i = 1|x \le 0\} = \int_U^\infty \frac{1}{2}\delta(v_i)dv_i = \begin{cases} 0 : U > 0 \; , \\ 1/2 : U \le 0 \; . \end{cases} \tag{13}$$

# 3 Mutual Information

Similar to other studies of SSR [9,10], mutual information was used as the performance measure. The mutual information between the instantaneous level of the input signal, $x$, and the number of units in state 1, $j$ is given by Shannon [27] to be

$$I = H(j) - H(j|x) \; , \tag{14}$$

where $H(j)$ denotes output entropy and $H(j|x)$ denotes the output entropy conditional on the input defined, respectively, by

$$H(j) = - \sum_{j=0}^{N} Prob\left\{ \sum_{i=1}^{N} y_i = j \right\} \log_2 Prob\left\{ \sum_{i=1}^{N} y_i = j \right\} , \qquad (15)$$

and

$$H(j|x) = - \int_{-\infty}^{\infty} dx P_x(x) \sum_{j=0}^{N} Prob\left\{ \sum_{i=1}^{N} y_i = j \middle| x \right\}$$

$$\times \log_2 Prob\left\{ \sum_{i=1}^{N} y_i = j \middle| x \right\} . \qquad (16)$$

where

$$Prob\left\{ \sum_{i=1}^{N} y_i = j \right\} = \int_{-\infty}^{\infty} dx P_x(x) Prob\left\{ \sum_{i=1}^{N} y_i = j \middle| x \right\}$$

$$= C_j^N \int_{-\infty}^{\infty} dx P_x(x) q_x^j (1 - q_x)^{N-j} , \qquad (17)$$

and

$$Prob\left\{ \sum_{i=1}^{N} y_i = j \right\} = - \int_{-\infty}^{\infty} dx P_x(x) \sum_{j=0}^{N} C_j^N q_x^j (1 - q_x)^{N-j}$$

$$\times \left[ \log_2 C_j^N + j \log_2 q_x + (N - j) \log_2 (1 - q_x) \right] . \qquad (18)$$

## 4 Results

The main results of this investigation are shown in Figs. 1 and 2. The results presented in Fig. 1 correspond to the simplified model (no rectification) and were obtained by computer simulation of (1) and (2) in which $F(x) = x$. This figure shows that without rectification there was a non-monotonic dependence of the mutual information, $I$, on the noise intensity, $D$, and that the mutual information could be optimized by a non-zero level of multiplicative noise. The phenomenon is therefore similar to stochastic resonance for additive noise [1,2]. As the number of threshold units, $N$, the mutual information also increased; thus the family of curves seen in Fig. 1 are similar to those observed in SSR [9,10]. However, there are two notable differences. First, in SSR information is transmitted even when the internal noise is zero. This occurs because the signal is suprathreshold and hence able to cross the threshold – these threshold crossings carry signal information. With multiplicative noise, zero noise level must always lead to zero infortmation regardless of the size of the signal. This is because the signal is multiplied by the noise before application to the threshold device – hence if the noise is zero then threshold

**Fig. 1.** The mutual information $I$ as a function of the noise intensity $D$ for an array of threshold units without rectification. The numerical results were obtained for a common threshold $U = 1$ and a common Gaussian input signal $x$ with a standard deviation $\sigma = 1$. Results are shown with multiplicative noise, where the input to each threshold unit was given by $v_i = Dx\eta_i$, for various numbers of threshold units from $N = 1$ to $N = 128$ (*solid lines*). Results are also shown for a deterministic channel where the input to each threshold unit was given by $v_i = Dx$ (*dashed line*)

crossings are not possible. Secind, we also note from Fig. 1 that the location of the maxima of the curves are almost identical for different values of $N$. Again this differs from SSR with additive noise where it is observed that the position of the maxima increase to larger $D$ as $N$ is increased.

Figure 2 shows the results obtained by computer simulation of (1) and (2) in which the input signal was rectified ($F(x) = R(x)$). In contrast to the results shown in Fig. 1 for the model without rectification, the plot of mutual information against noise intensity for the model with rectification had the following features. First, a local maxima existed only for sufficiently large value $N$, in these studies this required $N > 8$. Second, the local maxima for different values of $N$ occurred at different noise intensities; the location of the maximum shifted further to the region of weak noise intensity for larger values $N$. Third, the mutual information asymptotically approached 1 from above or below for large $N$ and $D$. We also note that the maximum information attained with rectification is slightly less than the that with no rectification.

The dashed lines in Figs. 1 and 2 show the mutual information for deterministic channels defined by $v_i = Dx$ (Fig. 1: no rectification) and $v_i = DR(x)$ (Fig. 2: rectification); because these channels were deterministic, an increase in the number of channels would not have lead to an increase in mutual

**Fig. 2.** The mutual information $I$ as a function of the noise intensity $D$ for an array of threshold units with rectification. The model parameters were otherwise as shown in Fig. 1. The *solid lines* show results with multiplicative noise, where the input to each threshold unit was given by $v_i = Dx\eta_i$, for various numbers of threshold units from $N = 1$ to $N = 128$. The *dashed line* shows the result for a deterministic channel where the input to each threshold unit was given by $v_i = Dx$

information (because each channel is identical and therefore carries identical information). At each noise intensity, the power of $v_i$ was the same for the deterministic and stochastic channels. A comparison of these plots with those from the stochastic channels with multiplicative noise shows that multiplicative noise can enhance the ability of an array of threshold units to transfer information if there are a sufficient number of threshold units.

We now discuss what happens when the signal is scaled (or eqivalently the threshold) by a factor $\alpha$ so that the scaled signal $\tilde{x}$ is given by $\tilde{x} = \alpha x$. We note that the models with and without rectification have the respective equivalences shown below:

$$v_i = D\tilde{x}\eta_i = D(\alpha x)\eta_i = (D\alpha)x\eta_i = \tilde{D}x\eta_i \ ,$$
$$\tilde{v}_i = DR(\tilde{x})\eta_i = DR(\alpha x)\eta_i = D\alpha R(x)\eta_i = \tilde{D}R(x)\eta_i \ , \qquad (19)$$

Consequently, this analysis shows that there is a simple scaling behaviour between the variance of the Gaussian input signal and the noise level $D$; hence there also exists a similar scaling between signal level and threshold level $U$. This is illustrated in Figs. 3 and 4 where it is shown that the curves of mutual information against noise intensity have identical forms for different values of $U$ but are shifted along noise intensity axis, e.g. increasing the threshold by a

**Fig. 3.** The mutual information $I$ as a function of the noise intensity $D$ for the model without rectification. The numerical results were obtained for a Gaussian input signal with standard deviation $\sigma = 1$



**Fig. 4.** The mutual information $I$ as a function of the noise intensity $D$ for the model with rectification. The numerical results were obtained for a Gaussian input signal with standard deviation $\sigma = 1$

factor of 10 causes the optimum mutual information to occur when the noise intensity is also 10 times larger.

## 4.1 Detailed Analysis of the Model in the Limit of the Large $D$

Perhaps the most striking result of the multiplicative noise is that at large noise intensity $D$, the mutual information tends to a non-zero value (1 for large $N$) when the signal is rectified. This is significantly different to what happens with additive noise. To understand this effect we now undertake a detailed analysis of the large noise limit.

Taking the non-rectified model first we can see that in the limit of high noise intensity, (18) and (15) can be calculated analytically. For a finite threshold $(-\infty < U < +\infty)$, and in the limit of large values of $D$, the instantaneous level of the input signal and the threshold level have no effect on the output state of each threshold unit and the probability that the output is in state 1, $q(x)$ is given by the probability that the Gaussian noise exceeds 0; in the limit of large $D$, $q_x$ is therefore 0.5. The probability that $j$ units are in the state 1 is therefore independent of $x$ and is given by

$$Prob\left\{\sum_{i=1}^{N} y_i = j \middle| x\right\} = \frac{C_j^N}{2^N}, \tag{20}$$

From this indpendence between $j$ and $x$, it follows that

$$Prob\left\{\sum_{i=1}^{N} y_i = j\right\} = Prob\left\{\sum_{i=1}^{N} y_i = j \middle| x\right\}, \tag{21}$$

and the output entropy and output entropy conditional on the input are therefore identical. From (21), it follows that the mutual information in the limit of large $D$ is zero; this analytic result is confirmed by the simulated results in Fig. 1.

We now consider what happens for the model with rectification in the limit of the large noise intensity. For the threshold $U$ less than or equal to zero, $q_x$ is again equal to 0.5 at high noise intensity and the behaviour of the models with and without rectification is therefore the same, i.e $\lim_{D\to+\infty} I = 0$. But when the threshold $U$ is positive and finite $(0 < U < +\infty)$, the probability $q_x = 0$ if $x \leq 0$, and $q_x = 1/2$ if $x > 0$. The conditional probability is therefore given by

$$Prob\left\{\sum_{i=1}^{N} y_i = j \middle| x \leq 0\right\} = \begin{cases} 0 : j \neq 0, \\ 1 : j = 0, \end{cases}$$

$$Prob\left\{\sum_{i=1}^{N} y_i = j \middle| x > 0\right\} = \frac{C_j^N}{2^N}. \tag{22}$$

Therefore, in contrast to the model without rectification, the conditional probability of the model with rectification is dependent on $x$. From (22) it follows that

$$Prob\left\{\sum_{i=1}^{N} y_i = j \Big| j \neq 0\right\} = \int_{-\infty}^{\infty} dx P_x(x) Prob\left\{\sum_{i=1}^{N} y_i = j \Big| x\right\}$$

$$= \int_{-\infty}^{0} dx P_x(x) Prob\left\{\sum_{i=1}^{N} y_i = j \Big| x \leq 0\right\}$$

$$+ \int_{0}^{\infty} dx P_x(x) Prob\left\{\sum_{i=1}^{N} y_i = j \Big| x > 0\right\}$$

$$= \frac{C_j^N}{2^N} \int_{0}^{\infty} P_x(x) dx \ ,$$

$$Prob\left\{\sum_{i=1}^{N} y_i = 0\right\} = \int_{-\infty}^{0} P_x(x) dx \ . \tag{23}$$

Since the probability distribution of the Gaussian input signal is symmetric, i.e. $P_x(x) = P_x(-x)$, and $\int_{-\infty}^{0} P_x(x) dx = \int_{0}^{\infty} P_x(x) dx = 1/2$, the previous expression can be rewritten as

$$Prob\left\{\sum_{i=1}^{N} y_i = j \Big| j \neq 0\right\} = \frac{C_j^N}{2^{N+1}}, \qquad Prob\left\{\sum_{i=1}^{N} y_i = 0\right\} = \frac{1}{2} \ . \tag{24}$$

By using (18), (15) and (24) the following expression for the information entropies is obtained,

$$H(j|x) = -\int_{-\infty}^{0} dx P_x(x) \sum_{j=0}^{N} P\left\{\sum_{i=1}^{N} y_i = j | x \leq 0\right\} \log_2 P\left\{\sum_{i=1}^{N} y_i = j | x \leq 0\right\}$$

$$- \int_{0}^{\infty} dx P_x(x) \sum_{j=0}^{N} P\left\{\sum_{i=1}^{N} y_i = j | x > 0\right\} \log_2 P\left\{\sum_{i=1}^{N} y_i = j | x > 0\right\}$$

$$= \frac{N}{2} - \frac{1}{2^{N+1}} \sum_{j=1}^{N} C_j^N \log_2 C_j^N \ , \tag{25}$$

and

$$H(j) = 1 + \frac{N}{2} - \frac{N+1}{2^{N+1}} - \frac{1}{2^{N+1}} \sum_{j=1}^{N} C_j^N \log_2 C_j^N \ , \tag{26}$$

The mutual information in the limit of large $D$ is therefore

$$\lim_{D \to +\infty} I = 1 - \frac{N+1}{2^{N+1}} \tag{27}$$

The last term of (27) decreases with increasing $N$, and for large $N$, the mutual information can be approximated by $I = 1$. Confirmation of this analytic result is shown in Fig. 2.

This result can, however, be understood in an intuitive manner as follows. As a consequence of the multiplicative structure, increasing the noise intensity is no different than increasing the signal amplitude whilst keeping the noise fixed (see scaling relations derived earlier). Consequently this is the same as fixing the noise and signal levels but reducing the threshold towards zero. If the signal is now half-wave rectified then no threshold crossings can occur during a negatively going portion of the signal but will occur for positive going signal portions (with probability approaching unity because of the very low threshold). The system therefore acts as a polarity detector or equivalently a two-state (i.e. 1-bit) quantiser. Hence the information must tend to 1-bit for sufficiently large noise, or sufficiently low threshold or sufficiently large signal. This effect is not observed in the non-rectified case because threshold crossings are now also possible in the negative portion of the signal. We finally note that the addition of any extra external noise to the signal will reduce the asymptotic value of the information achieved at large $D$.

## 5 Conclusion

In this paper we have investigated a simple system neural model that consisted of an array of threshold units with multiplicative noise. We have shown a relationship between mutual information and noise intensity with multiplicative noise that is similar to the phenomenon of stochastic resonance and suprathreshold stochastic resonance that have previously been observed with additive noise. In particular, we have shown that a non-zero level of multiplicative noise can lead to increased information transfer in threshold systems. We have further shown that the system can transmit more information when the number of the threshold units that receive the common input signal is increased.

Given that sensory neurons have a threshold-like behaviour and that there are multiplicative noise sources in sensory systems, our results suggest that the multiplicative noise may be an essential part of sensory coding. Although the results here concern a Gaussian input noise, they are general in that the analysis required only that the signal be symmetric, i.e. that $P_x(x) = P_x(-x)$. Identical results would therefore be expected in the limit of high noise-intensity for exponential or other symmetric distributions of the signal.

# References

1. L. Gammaitoni, P. Hanggi, P. Jung and F. Marchesoni, "stochastic resonance", *Rev. of Mod. Phys.* **70**, pp. 223–287 (1998).
2. A. Bulsara and L. Gammaitoni, "Tuning into noise", *Phys. Today* **49**, pp. 39–45 (1996).
3. A. Longtin, A. Bulsara and F. Moss, "Time-interval sequences in bistable systems and the noise-induced transmission of information by sensory neurons ", *Phys. Rev. Lett.* **67**, pp. 656–659 (1991).
4. J. K. Douglass, L. Wilkens, E. Pantazelou and F. Moss, "Noise enhancement of information transfer in crayfish mechanoreceptors by stochastic resonance", *Nature* (London) **365**, pp. 337–340 (1993).
5. J. J. Collins, C. C. Chow and T. T. Imhoff, "Stochastic resonance without tuning", *Nature* **376**, pp. 236–238 (1995).
6. N. G. Stocks, D. Allingham and R. P. Morse, "The application of suprathreshold stochastic resonance to cochlear implant coding", *Fluctuation and Noise Letters* **2** (3), pp. L169–L181 (2002).
7. N. G. Stocks and R. Mannella, "Suprathreshold stochastic resonance in a neuronal network model: A possible strategy for sensory coding", in *Future Directions for Intelligent Systems and Information Sciences*, N. Kasabov, ed., Physica-Verlag, Heidelberg, pp. 236–247 (2000).
8. N. G. Stocks and R. Mannella, "Generic noise-enhanced coding in neuronal arrays", *Phys. Rev. E* **64**, p. 030902(R), (2001).
9. N. G. Stocks, "Suprathreshold stochastic resonance in multilevel threshold systems", *Phys. Rev. Lett.* **84**, pp. 2310–2313 (2000).
10. N. G. Stocks, "Information transmission in parallel threshold arrays: Suprathreshold stochastic resonance", *Phys. Rev. E* **63**, pp. 411141–411149 (2001).
11. N. G. Stocks, "Information transmission in parallel threshold networks: Suprathreshold stochastic resonance and coding efficiency", in *Proceedings of the 16th international conference on Noise in Physical Systems and 1/f Fluctuations*, G. Bosman Ed., World Scientific, pp. 594–597 (2001).
12. M. D. McDonnell, N. G. Stocks, C. E. M. Pearce and D. Abbott, "Optimal information transmission in nonlinear arrays through suprathreshold stochastic resonance", (accepted for publication) *Phys. Lett. A.*
13. M. D. McDonnell, N. G. Stocks, C. E. M. Pearce and D. Abbott, "Quantization in the presence of large amplitude threshold noise ", *J. Fluctuation and Noise Lett.* **5**, pp. L457–L468 (2005).
14. D. Rousseau and F. Chapeau-Blondeau, "Constructive role of noise in signal detection from parallel arrays of quantizers", *Signal Processing* **85**, pp. 571–580 (2005).
15. M. D. McDonnell, D. Abbott and C. E. M. Pearce, "An analysis of noise enhanced information transmission in an array of comparators", *Microelectronics Journal* **33**, pp. 1079–1089 (2002).
16. V. C. Anderson, "Digital array phasing", *J. Acoustic. Soc. Am.* **32**, 867 (1960).
17. D. Rousseau, F. Duan and F. Chapeau-Blondeau, "Suprathreshold stochastic resonance and noise-enhanced Fisher information in arrays of threshold devices", *Phys. Rev. E* **68**, pp. 311071–3110710 (2003).

18. D. Rousseau and F. Chapeau-Blondeau, "Suprathreshold stochastic resonance and signal-to-noise ratio improvement in arrays of comparators", *Phys. Lett. A* **321**, pp. 280–290 (2004).
19. Y. Wang and L. Wu, "Stochastic resonance and noise-enhanced Fisher information", *Fluctuation and Noise Letters* **5**, pp. L435–L442 (2005).
20. T. Hoch, G. Wenning, K. Obermayer, "Adaptation using local information for maximizing the global cost", *Neurocomputing* **52– 54**, pp. 541–546 (2003).
21. N. G. Stocks, "Suprathreshold stochastic resonance: an exact result for uniformly distributed signal and noise", *Phys. Lett. A* **279**, pp. 308–312 (2001).
22. A. Manwani and C. Koch, "Detecting and estimating signals in noisy cable structure, I: Neuronal noise sources", *Neural Computation* **11**, pp. 1797–1829 (1999).
23. E. F. Evans, "Cochlear nerve and cochlear nucleus", in *Handbook of sensory physiology*, W. D. Keidel and W. D. Neff, eds., pp. 1–108, Springer, Berlin (1975).
24. T. Furukawa, Y. Hayashida and S. Matsura, "Quantal analysis of the size of excitatory post-synaptic potentials at synapses between hair cells and afferent nerve fibres in goldfish", *J. Physiol.* **276**, pp. 211–226 (1978).
25. C. E. Shannon and W. Weaver, *The mathematical theory of communication*, University of Illinois, Urbana (1949).

# Noise for Health: Phage-Based Rapid Bacterial Identification Method[*]

M.D. King[1], S. Seo[2], J. Kim[2], M. Cheng[2], S. Higgins[2], R. Young[3], D.H. McIntyre, B. Thien[1], A.R. McFarland[1], and L.B. Kish[2,+]

[1] Department of Mechanical Engineering, Texas A&M University, College Station, TX 77843-3123
[2] Department of Electrical Engineering, Texas A&M University, College Station, TX 77843-3128
[3] Department of Biochemistry and Biophysics, Texas A&M University, 2128 TAMU, College Station TX 77843-2128
[4] Integrative Center for Homeland Security, Texas A&M University, College Station, TX 77843-1112

**Abstract.** Recently, the authors have developed and demonstrated a novel sensing technology, named SEPTIC (SEnsing of Phage-Triggered Ion Cascades), for the rapid, inexpensive and specific identification of bacteria. The method combines the specificity and fast response of the bacteriophage ("phages"; viruses that specifically detect and kill bacteria) with the sensitivity of the nano-scale fluctuation-enhanced sensing.

In its prototype form based on a nanowell chip, SEPTIC has already been shown to be capable of unambiguous identification of live bacteria on a time scale of seconds to minutes, many times faster than any other system. The technology is based on using noise analysis to detect the massive ionic fluxes associated with the initial step of bacteriophage infection, the injection of the phage DNA into the cell. Ultimately, sensors based on this new technology would be able to save many lives.

**Keywords:** Bacterial identification; Bacteriophage, Biochip; Nanowell; Stochastic signals.

## 1 Introduction

One of the most imperative needs in clinical and agricultural practice as well as in homeland security applications is the rapid and sensitive identification

---

of bacteria. For example, certain strains of the bacterium *Escherichia coli* (*E. coli*) can cause widespread illness if they get into the food supply. Even though several technologies are available for the identification of bacteria or viruses in humans, veterinary and agricultural diagnostic laboratories, such as culturing and polymerase chain reaction (PCR), these approaches have difficulty such as time required for culturing of bacteria, expensive instrumentation or poor selectivity between living and dead bacteria. A rapid and inexpensive method for detecting and subtyping bacteria suitable for large-scale surveillance efforts, much less employment in the field, is not available.

The new method for the rapid detection and identification of bacteria, SEnsing of Phage-Triggered Ion Cascade (SEPTIC), is detecting and analyzing the electrical field caused by the stochastic emission of ions during phage infection by measuring the microscopic voltage fluctuations in a nano-well device [1–3], in less than 10 minutes. Virtually, a single bacterium can be detected and identified. Because the phage infection is a very selective process, where only bacteria of a specific strain are infected, a SEPTIC-biochip containing an array of sensors where each sensor is sensitized with a different phage, can detect and identify all relevant bacteria with extraordinary speed and selectivity. Though the method works with high reproducibility, presently there are more questions about its biology and physics than answers.

## 2 Background

Presently we will discuss the phage biology behind the phenomenon and also consider some practical physical conditions which may play an important role in determining the observed effects. Unfortunately, the exact details of the ion channels and the dynamics of ion emission, which are important to understand the possible charging effects, are unknown.

Bacteriophages are the most numerous biological entities, estimated at $10^{31}$ in the biosphere, and are unimaginably diverse [4]. Attempts to exploit the specificity of phages in detection and identification of pathogenic bacteria have been burdened by the requirement of culturing the target bacteria, growing the infected culture, and assaying the production of progeny virions, processes which at minimum require hours and also knowledge of the culture conditions is required. However, when we considered the fundamental pathway of the phage infection process, a potential way to avoid these limitations was suggested. The committed step in bacteriophage infection is irreversible adsorption. For double-stranded DNA (dsDNA) phages, this results from interactions between the specific adsorption apparatus, usually tail fibers, with specific receptors on the surface of the host cell [5–7]. For two of the three main morphotypes of dsDNA phages, the myophages with contractile tails and the siphophages with flexible tails, the injection of DNA into the host cell follows rapidly and involves the transitory formation of a channel through which the phage DNA passes into the target cytoplasm [9–11]. Concomitant

**Fig. 1.** The schematic representation of the experimental apparatus. The nanowell with the analyte drople. Side view (*left*) and upper view (*right*)

with injection is a short-lived membrane depolarization and an efflux of ions, including a substantial fraction of the $\sim$0.2M potassium salts present in the cytoplasm, at a rate of $\sim$10$^6$/sec per infected cell [8, 11]. This phenomenon represents an ideal opportunity for bacterial diagnostics, because it not only takes advantage of the well-known specificity available in bacteriophage but it can occur, given sufficient phage concentration, within seconds after admixture of the virions and cells [12]. Moreover, it requires no culturing of the analyte culture and detects only live bacteria.

## 3 Experimental Analysis

The nanowell is a lateral capacitor with a 150 nm gap between the two $4 \times 4$ micron size Ti plates (Fig. 1).

This biochip was used to perform initial experiments with bacteriophages $\lambda$S105 ($\lambda \Delta stf, tfa :: Cam^R$ cI857 S105). As $\lambda^S$ and $\lambda^R$ strains low motility *Escherichia coli* (*E. coli*) W3110 $\Delta fhuA$ and W3110 $\Delta fhuA\Delta lamB$ were used, respectively. The experimental details and results were presented earlier [1]. The basic experimental protocol was to mix the purified phage stock (about $2 \times 10^{10}$ pfu/ml) with the host cells (mid-log phase cells, washed and resuspended in 5 mM MgSO$_4$), incubate at 37$^\circ$C for various times and measure the voltage fluctuations in the nanowell. The isogenic host mutant strain (*E. coli* W3110 $\Delta fhuA\Delta lamB$) was used as negative control, for which we did not anticipate injection leakage of ions.

The power density spectrum of the electrical field fluctuations when there is reaction (phage infection) has $1/f^2$ shape, see Fig. 2, which implies correlation times beyond the time window of observation.

**Fig. 2.** Power density spectrum before the DNA passage and during the injection. Response of sensitive and resistant bacteria with phages. In the case of $\lambda^R$ bacteria (*E. coli* W3110 $\Delta fhuA$ $\Delta lamB$, negative response, blue line), phage infection did not occur, so the spectrum of voltage fluctuations in the nanowell follows approximately $1/f$ shape. In the case of positive response of $\lambda^S$ bacteria (*E. coli* W3110 $\Delta fhuA$, red line), phage infection occurred, so the fluctuations are enhanced, resulting in a steeper spectrum with a $1/f^2$ shape. Comprehensive data are given by [1]

During the experiment, a ∼5 μl droplet of the analyte containing bacteria and bacteriophage was placed in and around the nanowell, as indicated by the dotted circle in Fig. 1. Two probes, providing input to an external voltage amplifier, were firmly pressed on the contact pads. In order to prevent short circuit between the two pads and reduce unwanted noise caused by background ions in the solution, the surface of the chip, except for the nanowell sensing area, was covered with resist AZ5214. A 6 μm*8 μm rectangular window was then opened by photolithography. The voltage fluctuations induced on the nanowell device by the electrical field was amplified by a low-noise *preamplifier* SR560 with high input impedance (100 M$\Omega$) and fed into *signal acquisition unit* ML750 PowerLab/4SP as illustrated in Fig. 1. The power density spectrum of the fluctuations was calculated by a *Dynamic Signal Analyzer* SR785. The nanowell was placed in a *double screening box* (Amuneal Manufacturing Corp.) to prevent electromagnetic disturbance. The double screening box and the preamplifier were placed on an *anti-vibration platform* 100 BM-2 (Nano-K) to avoid potential artifacts caused by vibrations. The time window of the determination of the power density spectrum $S_u(f)$ was 2 minutes.

An important rule is that the analyte droplet should not touch the two probes connected to the external preamplifier, otherwise large unwanted noise would result.

In the mixture of non-adsorbing phages and bacteria (the phage does not infect the bacteria), the voltage fluctuations in the nanowell were small,

displaying PSD of approximately $1/f$ shape. On the other hand, in the solution where the phages invade the bacteria (i.e., the bacteria is sensitive to the phage), we observed large and slow stochastic waves with various time and amplitude scales. These fluctuations had PSD of approximately $1/f^2$ shape in the frequency range of 1–10 Hz, which is conjectured to be due to the ion efflux from the bacteria and possibly the Brownian motion of the charged bacteria. Figure 2 shows an example of PSD plots corresponding to sensitive and resistant bacteria that were incubated for 3 minutes and then mixed with $\lambda$S105 phages.

Note that other sources of electric field noise may also exist, including temperature fluctuations, background ions, etc. However, our experiment showed that these effects were much weaker than the transitory ion leakage due to phage invasion. Therefore, our fluctuation-enhanced nanoscale electric field detection system was successfully applied in all our experiments.

The detection sensitivity can be significantly improved if the external preamplifier is integrated into the same chip so that the noise due to external cable and input stage is eliminated. An estimation of the minimum number of bacteria that can be detected is shown in Fig. 3, based on Linear Response Theory. When a JFET preamplifier with thermal noise reduction technique is integrated into the chip, the SEPTIC technology is projected to detect the presence of 2 bacteria. For certain phages, the sensitivity is even 100 times higher, i.e., SEPTIC can detect a single bacterium even if the number of ions leaking from the bacterium were 100 times less.



**Fig. 3.** Projected sensitivities, based on the assumption of linear response with the bacteriophage lambda $\Delta$UR. Figure is based on data extracted from a recent publication [1]. With some types of phages, the sensitivity limits are about 100 times lower

Taking the above-mentioned facts into account, the following aspects are to be considered during the talk:

(i)     The phage attachment to the specific receptors on the bacterial surface, triggering the opening of ion channels.

(ii)    The difference between the DNA injection mechanisms of the two siphophages ($\lambda$ versus T5).

(iii)   Currents and ion numbers: the net ion charge vs. the current observed during the measurements.

(iv)    Mean distance between bacteria: about 10 micron during the experiments published so far.

(v)     Charging energy and mirror force: if the ion emission is not neutral, the bacteria may get charged. Considerations of the minimal effect (based on the equipartition energy, kT/2) indicate that the charging is of the order of 1000*e. That would yield a considerable mirror force within the Debye screening length.

(vi)    Debye length <1 micron (the deionized water has about 1 micron Debye length); therefore mirror forces seem to be screened and can not serve as a possible explanation of bacteria attraction to the surface.

(vii)   Diffusion coefficients/times of ions. $D$ is in the order of $10^{-9}$ m$^2$/s, thus to diffuse out of the nanowell's vicinity (8 micron) they need less than 0.1 second.

(viii)  Diffusion coefficients/times of bacteria. D is in the order of $10^{-12}$ m$^2$/s thus bacteria need about 10 seconds to diffuse out of the nanowell's vicinity.

(ix)    Debye relaxation time and high-pass filter response. Charge changes within the Debye relaxation time can get through the Debye screening; however, this phenomenon is characterized by strong high-pass filtering (time-derivative) characteristics.

# 4 Earlier Model Considerations

The physical basis of the fluctuations detected using the nano-well device has not been unambiguously determined but a reasonable model can be proposed. During the process of DNA injection of a siphophage or a myophage, each irreversibly-adsorbed virion triggers the opening of a single channel in the cytoplasmic membrane, by as yet unknown components, through which the phage DNA molecule passes. After the DNA passes through into the cytoplasm, the channel is sealed, usually within 1 minute and, again, by unknown processes. Bulk solution measurements have shown that injection is coupled to transient cellular depolarization requiring ion flows on the order of $10^8$ ions per infected cell [11, 13]. The emitted ions will execute a rapid Brownian motion and many will be able to escape from the vicinity of the bacterium. Together with the randomness of both the timing and the spatial orientation

of the ion emission, these ion leakage events would generate stochastic spatiotemporal electrical field fluctuations at the micron or submicron scale, as detected in the nano-well device. In this part, we consider a few possible noise generation mechanisms [14]. Most of them do not explain the experimental results and, at the moment, there are more questions than answers. A valid picture has to explain both the magnitude and the slow dynamics of the noise.

(a) Thermal noise? No, and neither is the Brownian motion of ions. Reason: that would decrease by increasing ion concentration because of the corresponding lower resistance.
(b) Excess motion of bacteria? No, they are not motile.
(c) Spatiotemporal charge density fluctuations caused by the emitted ion clouds? No, that dissipates too fast, see point vii above.
(d) Displacement currents? Due to the high-pass filter effect, to produce a $1/f^2$ noise by fast charging spikes is possible only if there are strong inter-spike correlations (e.g., spike frequency modulation), however the efficiency would be very low. So far, no indication exists that the bacteria would produce a sequence of charging spikes with the proper modulation. However, as this effect may cause mirror forces beyond the Debye length, it should be further investigated.
(e) Concentration cell? A bacterium asymmetrically located in the nanowell and slowly emitting ions may yield a voltage due to the concentration cell effect. Because of the slow diffusion dynamics of bacteria (see point viii), this effect is a possible candidate. Open questions are if, during phage infection, the bacteria leak ions for a prolonged time which would be necessary for this effect to be a proper explanation.

## 5 Conclusions

During a phage infection process, the myophage or siphophage injects its DNA into the host cell, resulting in a transitory, massive ion efflux. This phenomenon provides a perfect opportunity for the rapid and specific detection of bacteria, due to the fact that a certain type of phage can only invade a specific type of bacteria, resulting in an ion efflux within seconds of phage DNA injection. To this end, we have successfully fabricated a biochip whose core element is a nanowell device on $LiNbO_3$ substrate, which comprises of two 4 μm wide Ti electrodes with a 150 nm gap between them. The fabrication was performed by combining electron beam lithography, contact photolithography and RIE without using lift-off process. Detection experiments were conducted by mixing phage with sensitive and resistant bacteria, respectively, on the nanowell region (no voltage is applied to the nanowell) and analyzing the voltage fluctuations across the surface. In mixtures where the bacteria were sensitive to the phage, large $1/f^2$ noise (1–10 Hz) was observed, while with mixtures where bacteria were resistant to the phage, only small $1/f$ noise due

to preamplifier was observed. The large $1/f^2$ noise is assumed to derive from temporary charging fluctuations of the bacteria due to phage-induced ion leakage. Our preliminary experiment showed 100% success rate in identification of bacteria on the scale of minutes, whereas other known technologies require hours to days of bacteria culturing and are often not specific. The sensitivity is expected to be at the single bacterium level. Ultimately, this technology could prove invaluable in clinical, veterinary and agriculture practice, as well as in applications to microbiological threat detection and reduction in biodefense applications, integrated with the IJAG (Inkjet Aerosol Generator), capable of generating clusters of bacterial particles in a 1liter/min airflow controlled by a LASER detection system.

# References

1. M.D. King, S. Seo, J.U. Kim, R.F. Young, M. Cheng and L.B. Kish. 2005. Rapid Detection and Identification of Bacteria: SEnsing of Phage-Triggered Ion Cascade (SEPTIC). Journal of Biological Physics and Chemistry, **5**:(2005).
2. L.B. Kish, M. Cheng, J.U. Kim, S. Seo, M.D. King, R. Young, A. Der, G. Schmera. 2005. Estimation of detection limits of the phage-invasion based identification of bacteria, *Fluctuation and Noise Lett.*, **5** (1): L105 L108 (2005).
3. M.D. King, M. Cheng, S. Seo, J.U. Kong, L.B. Kish and R. Young. SPIE Symposium (Conference on Noise in Biological systems), Austin (2005).
4. R.W. Hendrix, M.C. Smith, R.N. Burns, M.E. Ford, & G.F. Hatfull, "Evolutionary relationships among diverse bacteriophages and prophages: all the world's a phage". *Proc Natl Acad Sci.USA.* **96**, 2192–2197 (1999).
5. E.B. Goldberg, L. Grinius, & L. Letellier, "Recognition, attachment and injection", pp. 347–357. *In* Karam, J.D. et al. (eds.), *Molecular biology of bacteriophage T4*. American Society for Microbiology, Washington D.C., USA, (1994).
6. A.A. Lindberg, "Bacteriophage receptors. *Annu Rev Microbiol.* **27**, 205–241 (1973).
7. U. Henning & S. Hashemolhosseini, Receptor recognition by T-even type coliphages. *In* Karam, J.D. et al. (eds.), *Molecular biology of bacteriophage T4*. American Society for Microbiology, Washington D.C., USA, (1994).
8. L. Letellier & P. Boulanger, "Involvement of ion channels in the transport of phage DNA through the cytoplasmic membrane of *E. coli*". *Biochimie* **71**, 167–174 (1989).
9. S. Silver, E. Levine & P.M. Spielman Cation fluxes and permeability changes accompanying bacteriophage infection of *Escherichia coli. J Virol* **2**, 763–781, (1968).
10. D.J. MacKay & V.C. Bode, "Events in lambda injection between phage adsorption and DNA entry". *Virol* **72**, 154–166 (1976).
11. P. Boulanger & L. Letellier, "Characterization of ion channels involved in the penetration of phage T4 DNA into *Escherichia coli* cells". *J Biol Chem* **263**, 9767–9775 (1988).
12. M. Schwartz, "The adsorption of coliphage lambda to its host: effect of variations in the surface density of receptor and in phage-receptor affinity". *J Mol Biol* **103**, 521-536 (1976).

13. E.V. Kalasauskaite, D.L. Kadisaite, R.J. Daugelavicius, L.L. Grinius & A.A. Jesaitis, "Study on energy supply for genetic processes. Requirement for membrane potential in *Escherichia coli* infection by phage T4". *Eur J Biochem.* **130**, 123–130 (1983).
14. R.J. Biard and L.B. Kish, "Enhancing the the sensitivity of the SEPTIC bacterium detection method by concentrating the phage-infected bacteria via DC electrical current", *Fluctuations and Noise Lett.* **5** (June 2005), in press.

# Contributed Papers

# Parametric Resonance Near Hopf-Turing Instability Boundary

A. Bhatacharyay[1], and J.K. Bhattacharjee[2]

[1] Dipartimento di Fisica "G. Galilei", Università di padova, Italy
   `arijit@pd.infn.it`
[2] Department of Theoretical Physics, Indian Association for the Cultivation of
   Science, India
   `tpjkb@mahendra.iacs.res.in`

Dissipative chemical systems which are having to interact with its environment do not generally see the feeding or removal of species in a uniform or constant manner as is generally taken in simpler form of the theory of an ideal system. There are variations and an expansion of them in Fourier modes can always make some of them vulnerable (or useful if we have some control on it) for the present state of the system. A systematic study of parametric resonance is therefore very important for such systems. Surface wave of fluids generated by vertical oscillation is a well known example where parametric resonance breaks the continuous spatial symmetry [1–3]. In reaction diffusion systems the effect of parametric resonance has been widely studied to see frequency entrainment and multiphase oscillation [4–7]. Existence of multiphase oscillations are theoretically accounted for by showing the stability of phase separated oscillatory orders in complex Ginzburg-Landau (**GL**) equation or in some other reaction diffusion model. Recently observed cluster pattern which does not involve an intrinsic length scales [8, 9] are interesting spatial structures generated by parametric resonance in a chemical system. In this paper we are going to work on a simple reaction diffusion model namely Gierer-Meinhardt (**GM**) model which has two coupled variables or chemical species with varying diffusivity and interact to produce Turing as well as Hopf instability to its basic homogeneous steady state [10]. We will be showing that in the parameter regime where a homogeneous temporal oscillatory mode (Hopf mode) grows, a global oscillatory force in the form of a temporal modulation of one of the parameters of the system can resonate with the existing Hopf mode to generate a new length scale which is different from that characterized by the Turing Instability. A similar linear analysis has been accepted for publication in *Physics Letters A* however in the present work we are explicitly calculating the wave number of the instability in a particular case and extending the theory to see the effect of nonlinearity.

The simplified form of GM model in one dimension looks like [10]

$$\frac{\partial a}{\partial t} = D\frac{\partial^2 a}{\partial x^2} + \frac{a^2}{b} - a + \sigma$$

$$\frac{\partial b}{\partial t} = \frac{\partial^2 b}{\partial x^2} + \mu(a^2 - b) . \tag{1}$$

The variables $a$ and $b$ are local concentrations of the activator and the inhibitor species, $\sigma$ is feeding rate of the activator in the system where the cross reaction coefficient and the removal rate of the inhibitor are the same and is denoted by $\mu$. A linear stability analysis of the above model for its homogeneous steady state characterized by the $a = 1 + \sigma$ and $b = (1 + \sigma)^2$ shows that the basic homogeneous state loses stability to a Turing instability when $\mu D \leq (\sqrt{2/(1 + \sigma)} - 1)^2$ keeping $\mu > (1 - \sigma)/(1 + \sigma)$. Below the horizontal continuous line at $\mu = \mu_{00} = (1 - \sigma)/(1 + \sigma)$ in the phase diagram ($\mu$ vs $D$ plane Fig.1) a Hopf mode grows with $k = 0$ where $k$ is the wave number. There is a codimension 2 point where the boundaries $\mu D = (\sqrt{2/(1 + \sigma)} - 1)^2$ (broken line in Fig.1) and $\mu = (1 - \sigma)/(1 + \sigma)$ meet. The oscillation frequency of the zero-growth Hopf instability on the phase boundary $\mu = \mu_{00}$ is given by $\omega_0 = \sqrt{(1 - \sigma)/(1 + \sigma)}$. In what follows we will develop a perturbation expansion near this boundary $\mu = \mu_{00}$ where the removal rate $\mu$ has an $O(\epsilon)$ temporal modulation of frequency $\omega$.



**Fig. 1.** Phase diagram in $\mu$ vs D plane as obtained from the linear stability analysis of the simplified model

Here we are going to consider the problem which is often treated in hydrodynamic instabilities – the one where the control parameter is given a sinusoidal temporal variation [11–15]. We put a modulation in $\mu$ as $\mu = \mu_0(1 + \epsilon \cos(\omega t))$. The linearized equation now reads

$$L_0 \begin{pmatrix} \delta a \\ \delta b \end{pmatrix} = \epsilon \cos(\omega t) \begin{pmatrix} 0 & 0 \\ 2\mu_0(1 + \sigma) & -\mu_0 \end{pmatrix} \begin{pmatrix} \delta a \\ \delta b \end{pmatrix} . \tag{2}$$

Let us expand $\mu_0$ as $\mu_{00} + \epsilon\mu_{01} + \epsilon^2\mu_{02} + h.o.t.$, where $\mu_{00}$ has the value as has already been mentioned which can also be verified from the $O(1)$ solution

of the (2). In $O(1)$, the eigenvectors for the homogeneous oscillatory state at the Hopf-bifurcation boundary $\mu = \mu_{00}$ are

$$\begin{pmatrix} \delta a_\pm \\ \delta b_\pm \end{pmatrix} e^{\pm i\omega t} = \begin{pmatrix} 1 \\ 2\mu_{00}(1+\sigma) \\ \mu_{00} \pm i\omega_0 \end{pmatrix} e^{\pm i\omega_0 t} \tag{3}$$

where the linear operator $L_0$ has the form

$$L_0 = \begin{pmatrix} \dfrac{\partial}{\partial t} - \dfrac{1-\sigma}{1+\sigma} & \dfrac{1}{(1+\sigma)^2} \\ -2\mu(1+\sigma) & \dfrac{\partial}{\partial t} + \mu \end{pmatrix} \tag{4}$$

Thus we see that $\delta a_\pm$ and $\delta b_\pm$ has a constant phase difference since the production of the inhibitor depends on the concentration of the activator. At this point we consider this Phase difference $\phi(\mu)$ has an additive part which varies on a slower time and larger space scale. So the structure of $\phi(\mu)$ is taken as

$$\phi = \phi_c + \delta\phi_0(X, \tau) \tag{5}$$

where $\phi_c$ is the critical value of and can be easily obtained from (4). Let us expand $\delta A$ and $\delta B$ as

$$\delta a = \delta a_0 + \epsilon \delta a_1 + \epsilon^2 \delta a_2 + h.o.t.$$
$$\delta b = \delta b_0 + \epsilon \delta b_1 + \epsilon^2 \delta b_2 + h.o.t.$$

and introduce the multiple space and time scales as $x = x_0 + \epsilon X$ and $t = t_0 + \epsilon\tau$ respectively. Now, at $O(\epsilon)$ the equation looks like looks like

$$L_0 \begin{pmatrix} \delta a_1 \\ \delta b_1 \end{pmatrix}$$
$$= \begin{pmatrix} \delta a_+^\phi \frac{\partial}{\partial\tau}(\delta\phi_0) \\ -\delta b_+^\phi \frac{\partial}{\partial\tau}(\delta\phi_0) + \mu_{01}[2(1+\sigma)\delta a_+ - \delta b_+^\phi] \pm \frac{\mu_{00}}{2}[2(1+\sigma)\delta a_- - \delta b_-^\phi] \end{pmatrix}$$
$$+ \begin{pmatrix} \delta a_-^\phi \frac{\partial}{\partial\tau}(\delta\phi_0) \\ \mp \delta b_-^\phi \frac{\partial}{\partial\tau}(\delta\phi_0) \pm \mu_{01}[2(1+\sigma)\delta a_- - \delta b_-^\phi] + \frac{\mu_{00}}{2}[2(1+\sigma)\delta a_+ - \delta b_+^\phi] \end{pmatrix} \tag{6}$$

where we have taken $\omega = 2\omega_0$ and on the r.h.s. only secular terms have been considered. The first term on the right hand side correspond to $e^{i\omega_0 t}$ whereas the second one is related to $e^{-i\omega_0 t}$. Note that for $\omega \neq 2\omega_0$ the equation is solvable for $\mu_{01} = 0$ and that will require the temporal part of $\delta\phi$ be a constant. A superscript $\phi$ of $\delta a_\pm$ and $\delta b_\pm$ indicates that it has an additional factor $e^{\delta\phi_{1,2}}$ (say) and the upper(lower) one of $\pm$ or $\mp$ signes in the above expression corresponds to symmetric(antisymmetric) combination of solutions. In the antisymmetric after we make a transformation as $\theta_{1,2} = e^{\delta\phi_{1,2}}$ we get a coupled set of linear equations as

$$-(1 + M)\frac{\partial \theta_1}{\partial \tau} + C_1\theta_1 - C_2\theta_2 = 0$$

$$-(1 - \bar{M})\frac{\partial \theta_2}{\partial \tau} - \bar{C}_1\theta_2 + \bar{C}_2\theta_1 = 0 \tag{7}$$

where $M$, $C_1$ and $C_2$ are different complex numbers and a bar on them indicates complex conjugate. The characteristic equation for (7) is

$$\lambda^2(1 - |M|^2 + M - \bar{M}) + \lambda(M\bar{C}_1 + \bar{M}C_1 - C_1 + \bar{C}_1) + |C_2|^2 - |C_1|^2 = 0$$

Since $\bar{M} - M$ and $C_1 - \bar{C}_1$ are two different pure imaginary numbers the above expression will result in a $\lambda$ which in general is a complex number. So the frequency of oscillation for the system will change. When the solution of $O(1)$ equation considered symmetric characteristic equation can be written in the form

$$\lambda^2(1 + |M|^2 + M + \bar{M}) - \lambda(M\bar{C}_1 + \bar{M}C_1 + C_1 + \bar{C}_1) + |C_1|^2 - |C_2|^2 = 0$$

The above equation requires $|C_1| > |C_2|$ for $\lambda$ to be complex and requires one to go sufficiently below the phase boundary since $|C_1|$ is proportional to $\mu_{01}$. We also see that there is no oscillation caused by the external force on the phase boundary for which $|C_1| = 0$ in the symmetric case whereas in the antisymmetric case it does.

In what follows we will concentrate on the resonance at $O(\epsilon^2)$ at the forcing frequency $\omega \neq 2\omega_0$. In this case at $O(\epsilon)$ for $\mu_{01} = 0$ the phase variation will come out as not time dependent and from now on we will take $\mu_{01} = 0$. The solution of $O(\epsilon)$ equation for $(\omega \neq 2\omega_0)$ will now definitely include $e^{i(\pm\omega_0 \pm \omega)}$, where this term will make secular terms appear in the next higher order for all forcing frequencies. Thus, The $O(\epsilon)$ solution will come as

$$\begin{pmatrix} +\delta\bar{a}_\pm \\ +\delta\bar{b}_\pm \end{pmatrix} = \frac{\mu_{00}}{2\Delta_{(\omega_0 \pm \omega)}}[2(1+\sigma)\delta a_+ - \delta b_+] \begin{pmatrix} \frac{-1}{(1+\sigma)^2} \\ i(\omega_0 \pm \omega) - \mu_{00} \end{pmatrix} e^{i(\omega_0 \pm \omega)t}$$

and

$$\begin{pmatrix} -\delta\bar{a}_\pm \\ -\delta\bar{b}_\pm \end{pmatrix} = \frac{\mu_{00}}{2\Delta_{(-\omega_0 \pm \omega)}}[2(1+\sigma)\delta a_- - \delta b_-] \begin{pmatrix} \frac{-1}{(1+\sigma)^2} \\ i(-\omega_0 \pm \omega) - \mu_{00} \end{pmatrix} e^{i(-\omega_0 \pm \omega)t}$$

The $\Delta_{(\pm\omega_0 \pm \omega)}$ being the determinant of $L_0$ at frequency $(\pm\omega_0 \pm \omega)$. In the case when $O(\epsilon)$ equation proportional to $e^{i(\omega_0 \pm \omega)}$ the solvability condition at $O(\epsilon^2)$ will result in following phase equation

$$-\left[\frac{\partial\delta\phi_0}{\partial X^2} + \left(\frac{\partial\delta\phi_0}{\partial X}\right)^2\right] \times \delta b_+^\phi = \mu_{02}[2(1+\sigma)\delta a_+ - \delta b_+^\phi]$$

$$-\frac{\mu_{00}^2}{2}\frac{\mu_{00}}{2\Delta_{(\omega_0 + \omega)}}(2(1+\sigma)\delta a_+ - \delta b_+^\phi)[1 + i(\omega_0 \pm \omega)]$$

$$\tag{8}$$

In writing the above equation we have associated the phase variable to the inhibitor part only because of the fact that a coupling is not necessary for the result we are after and it would make the analysis simpler. Let us take (for the sake of simplicity) that $\mu_{02} = 0$ and after simplification (10) can be written in the form

$$- \left[ \frac{\partial \delta\phi_0}{\partial X^2} + \left( \frac{\partial \delta\phi_0}{\partial X} \right)^2 \right] = C_1 e^{-\delta\phi_0} + C_2 , \tag{9}$$

where $C_1$ and $C_2$ are two complex numbers. Let us do a Cole-Hopf transformation as $\theta = e^{\delta\phi_0}$ and followed by $\phi = C_1 + C_2\theta$ will ultimately give us the equation

$$-\frac{\partial^2 \phi}{\partial \phi^2} = C_2\phi . \tag{10}$$

From the solution of the above equation and coming back to old variable we arrive at

$$e^{\delta\phi_0} = \frac{e^{iC_3 x}}{C_2} - \frac{C_1}{C_2} \tag{11}$$

where $C_3 = \sqrt{C_2}$ is a complex number whose real part will give the wave number generated which in terms of actual parameters (present case) is given by

$$k_1 = \pm\sqrt{\frac{\mu_{00}^3}{8[\mu_{00} - (\omega_0 - \omega)^2]}} \times [1 \pm (1 + 4(\omega_0 - \omega))^{\frac{1}{2}}] \tag{12}$$

and

$$k_2 = \pm\frac{\mu_{00}^3(\omega_0 + \omega)}{4[\mu_{00} - (\omega_0 + \omega)^2]} \times \sqrt{\frac{\mu_{00}^3}{8[(\omega_0 + \omega)^2 - \mu_{00}]}} \times [1 \pm (1 + 4(\omega_0 + \omega))^{\frac{1}{2}}] \tag{13}$$

Let us extend the theory to include nonlinear terms. By shifting the origin of the concentration scale to the basic homogeneous fixed point we can rewrite the model (1) in the following form.

$$\begin{pmatrix} \frac{\partial}{\partial t} - D\,\nabla^2 - \frac{1-\sigma}{1+\sigma} & \frac{1}{(1+\sigma)^2} \\ -2\mu(1+\sigma) & \frac{\partial}{\partial t} - \nabla^2 + \mu \end{pmatrix} \begin{pmatrix} a \\ b \end{pmatrix} = \begin{pmatrix} \frac{-b\frac{\partial a}{\partial t} + bD\nabla^2 a + a^2 - ab}{(1+\sigma)^2} \\ \mu a^2 \end{pmatrix} \tag{14}$$

and now, if we expand the variables in the way that $a = \epsilon^{\frac{1}{2}} a_0 + \epsilon a_1 + \epsilon^{\frac{3}{2}} a_2 +$ h.o.t. and $b = \epsilon^{\frac{1}{2}} b_0 + \epsilon b_1 + \epsilon^{\frac{3}{2}} b_2 + $ h.o.t., where the parameters are expanded as $\mu = \mu_0(1 + \epsilon^{\frac{1}{2}} \cos\omega t)$, $\mu_0 = \mu_{00} + \epsilon^{\frac{1}{2}}\mu_{01} + \epsilon\mu_{02} + $ h.o.t. Multiple scales are introduced as $x = x + \epsilon^{\frac{1}{2}} X$, $t = t + \epsilon^{\frac{1}{2}}\tau$ and there is an additive phase variation on large space and time scales as mentioned in the linear case.

In $O(\epsilon^{\frac{1}{2}})$ and $O(\epsilon)$ the nonlinear part does not introduce any extra secular term and the theory is the same as that obtained from the previous linear analysis. At $O(\epsilon^{\frac{3}{2}})$ the secular term appearing from the nonlinear part will make the resulting phase equation nonlinear which would look like

$$-\frac{\partial^2 \theta}{\partial X^2} = C_1 + C_2\theta + C_3\theta^2 \tag{15}$$

where $C_1$, $C_2$ and $C_3$ are complex numbers. The $\theta^2$ term originates from nonlinear terms $a_1 b_0$ where $b_0$ is $O(\epsilon^{\frac{1}{2}})$ solution and $a_1$ is $O(\epsilon)$ solution. Other nonlinear terms will change $C_1$ and $C_2$ from its linear analysis expression. Clearly, this equation does not support creating a wave number as found in the linear theory and in the absence of $C_3$ if $C_2$ is nonzero one can get a spatial phase modulation.

Now, we would like to summarize by mentioning that from the linear theory one can easily see that there are two coexistent states one is homogeneous oscillatory and the other is a traveling wave whose wave number have been explicitly given in (12) and (13) under the conditions mentioned. This spatial instability is definitely different from the Turing instability. It is easy to note that the expression of the wave number does not involve the diffusion constant which is always present for Turing instability. The wave numbers are never zero unless $\mu_{00} = 0$ but we know that it requires $\sigma = 1$ which is unacceptable since the basic Hopf instability will vanish. Now, the presence of a force is necessary to have control on these spatial instabilities which are generated close to an instability boundary. Even when there is no parametric modulation by just moving a little away from the instability boundary $\mu = \mu_{00}$ one generates spatial instability and if that is unwanted one has the opportunity to clean the system up with a parametric modulation of any frequency other than $\omega = \omega_0$. The nonlinearity in general suppress the generation of new spatial length scale but close to the instability boundary such spurious instabilities can be quite vulnerable. In the $O(\epsilon)$ part of the theory (linear/nonlinear) we see that the frequency of oscillation of the system can change irrespective of the presence or absence of a forcing term ($|C_2| = 0$) if we move away from the instability boundary. The force is necessary for creation of a temporal instability on the phase boundary $\mu = \mu_{00}$. Here also the role of a particular mode $\omega = 2\omega_0$ is of controlling the situation and only in the presence of the force of particular frequency we can better tune $\mu_{01}$ to control the frequency of oscillation of the system.

# References

1. Miles J.W. (1984) J. Fluid Mech. 148:451–460
2. Milner S.T. (1991) J. Fluid Mech. 225:81–100
3. Staliunas K., Longhi S., Valcárcel G.J. (2002) Phys. Rev. Lett. 89:210406
4. Petrov V., Ouyang Q., Swinney H.L. (1997) Nature (London) 388:655–657
5. Elphick C., Hagberg A., Meron E. (1998) Phys. Rev. Lett. 80:5007–5010
6. Elphick C., Hagberg A., Meron E. (1999) Phys. Rev. E 59:5285–5291
7. Lin A.L., Bertram K., Martinez K., Swinney H.L. (2000) Phys. Rev. Lett. 84:4240–4243
8. Vanag V.K., Zhabotinsky A.M., Epstein I.R. (2001) Phys. Rev. Lett. 86:552–555

9. Vanag V.K., Yang L., Dolnik A.M., Zhabotinsky A.M., Epstein I.R. (2000) Nature (London) 406:389–391
10. Koch A.J., Meinhardt H. (1994) Rev. Mod. Phys. 66:1481–1507
11. Bhattacharjee J.K., Banerjee K. (1983) Phys. Rev. Lett. 51:2248–2251
12. Ahlers G., Hohenberg P.C., Lücke M. (1985) Phys. Rev. A 32:3493–3519
13. Hall P. (1975) J. Fluid Mech. 67:29–63
14. Fauve S., Kumar K., Laroche C., Beysens D., Garrabos Y. (1992) Plys. Rev. Lett. 68:3160–3163
15. Holms H. (1995) Introduction to perturbation methods, Springer-Verlag

# Recurrent Neural Networks in Rainfall–Runoff Modeling at Daily Scale

E.C. Carcano[1], P. Bartolini[2], and M. Muselli[3]

[1] I.M.A.G.E. Padova
[2] D.I.A.M. Genova
[3] C.N.R.-I.E.I.I.T. Genova

**Abstract.** This work aims to simulate potential scenarios in Rainfall-Runoff (R-R) transformation at daily scale, mainly perceived for the control and management of water resources, using feed-forward multilayer perceptrons (MLP) and, subsequently, Jordan Recurrent Neural Networks (JNN). R-R transformation is one of the most complex issue in hydrological environment due to high temporal and spatial variability, very strong and non linear interconnections among variables: a good challenge for Artificial Neural Networks (ANN). Abilities and limitations of MLP and JNN models have been investigated, especially focusing on drought periods where water resources management and control are particulary needed. The study compares the results of the two networks typologies to outputs from a conceptual linear model and then to physical context of two small Ligurian catchments. It also demonstrates the remarkable improvement obtained with the JNN approach especially when rainfall memory effect is employed as an additional input.

## 1 Introduction

The simulation of rainfall-runoff (R-R) relationships has been an unavoidable issue of hydrological research for several decades and has resulted in plenty of models proposed in literature.

Following Beck (1991), these models can merely be divided in: *metric, conceptual, and physics-based.* Metric models are deeply observation-oriented, pursuing the system response by extrapolating information from the available data. Conceptual models, on the other hand, describe all the relevant components of hydrological processes as simplified conceptualizations, whereas physics-based models aim to reproduce the hydrological behavior of a basin by using the concepts of classical continuum mechanics. Another distinction proposed in literature deals with different levels of prior knowledge available which lead to three different color-coded types of model: white, grey and black box. In the first case the model is perfectly known, in the second one some physical insight is allowed, but several parameters still need to be determined from data. In black-box models, unfortunately, no physical insight is possible

and the structure of the model is chosen inside families which show good flexibility and 'have been successfully' employed in the past [9]. Artificial Neural Networks (ANN) represent one of these families and have been widely investigated in Hydrology since the middle 1990's. They are especially appreciated for their abilities to treat difficult issues such as: high non linearity and huge amount of variables involved in R-R transformation. Nowadays, in fact, in our enviroment, the current trend seems to be to focus on the existing data rather than the physical process. Moreover, in situations where the information is available only at specific sites in a basin or when only rainfall-runoff data sequences are known, it becomes very difficult to develop a conceptual model of complex processes (Kokkonen, Jakeman, 20001; Thirumalaiah, Deo, 2000). In this scenario, ANN have been proved to provide better solutions when applied to: 1) complex systems that, otherwise, may be poorly mimiced or understood, 2) problems tainted by noise or that involve pattern recognition, diagnosis and generalization, 3) circumstances where input is incomplete or ambiguous by nature. An ANN is able of modeling the R-R relationships due to its ability to generalize patterns even when noisy and ambiguous input data are considered and to synthesize a reliable model without needing any prior knowledge on functional relationship between dependent and independent variables. Mathematically, it can be considered as a universal approximator having the ability to learn from examples without knowing explicitly the physics of the problems considered [2].

ANN hydrological applications have been widely devoted to *forecasting* rather than *data generation* context. *Data generation* and *forecasting* are the two major tasks in hydrological researches, and deal, respectively, with prediction of future values and simulation of potential scenarios. *Data generation* is anyway of vital importance in ordinary applications such as: management and optimum control of system or supervision of drought periods.

In the R-R prediction environment pioneer researches are due to [12] and [8], who used a feed-forward ANN to forecast runoff values only from rainfall data. Their purpose was to demonstrate the learning ability of the network, rather than applying it to a practical situation affected by noise. Moreover, several efforts encouraged the comparison between ANN and other models, such as classical ARMA [11] and ARMAX, coupled with Kalman Filter, in order to perform daily flow forecasting. Laio (2003) compared ANN to a nonlinear prediction model, based on deterministic chaos for the prediction of flood events at hourly scale. Very recent works have been done by [1,7] and [9], who coupled a feed-forward ANN to a simple black-box linear model (auxiliary model) to be able to forecast daily streamflows starting from daily measured rainfall, runoff and evapotraspiration series of several storm events of two large Indian catchments. Far away from *forecasting* problem is the work of [4], who used ANN to interpolate the management rules of Pozzillo reservoir, given from dynamic programming.

Still, no deep attention has been paid for attempting the simulation of streamflow data particulary at daily scale, and moreover, when rainfall rates

are the unique measures available. Hence, this work aims to investigate the performances of feed-forward and recurrent ANN in streamflow simulation context at daily scale, and their ability of managing the presence of a high number of zero values in rainfall input data.

# 2 Problem Description

The process of transformation of inflow into streamflow data requires a substantial number of informations (precipitation, evapotranspiration, temperature, soil moisture, . . . ), but when some input series are rather limited or even totally absents, conceptual models normally fail to represent the complex and non linear dynamics of R-R processes and empirical models become much more reliable. The worst situation occurs, like in our case, when the system relies only on observed precipitations and on its previous calculated streamflows. Despite a consequent very poor modeling, reproducing daily runoff values when only rainfall rates are provided is of practical use even when no water-gauges are available or also when time registration discrepancies occur between rainfall and runoff measures, thus offering misleading interpretations of the catchment's results. Moreover, R-R modelling, particulary at daily scale, entails a very challenging benchmark: long sequences of zeroes in rainfall input series, that, apart from "contradictory" information to process, are a unique and remarkable security in measures since they are, with a high probability, not affected by accidental errors. Our crucial question consists in finding an ANN procedure able to process long sequences of non rainy days and, at the same time, to reproduce the "system memory" over a catchment when minimal information data occur [3,12]. In this context, simulations start asking different ANN approaches to be able to reproduce, at first, modeled streamflows from a conceptual linear model ("controlled experiments") and then observed streamflows of two Ligurian catchments: Argentina and Impero.

# 3 Procedure Description

## 3.1 Analysis of the Referenced Conceptual Model

To test their ability of modeling R-R mean daily data, ANN have been questioned to reproduce, at least, outputs from a simple conceptual linear model with one reservoir. Conceptual streamflow values ensue from controlled experiments wherein input precipitations are considered already as effective rainfall, directly able to produce calculated daily streamflows (through the functional dependence (1) and skipping, therefore, the very complex and non linear trasformation from observed (total) rainfall into effective rainfall. Mostly for this reason target calculated streamflow are more simple to be reproduced,

since they are free from data errors. Instantaneous streamflow value is given by:

$$Q(t) = Q(0) \cdot e^{-\beta \cdot t} + q_{e,j} \left(1 - e^{-\beta \cdot t}\right) \tag{1}$$

where:

- $Q(t)$ is the istantaneous calculated streamflow value;
- $Q(0)$ is the initial istantaneous streamflow value calculated at the beginning of the interval considered;
- $q_{e,j}$ is the constant daily discharge, given by the product of mean daily rainfall intensity and catchment's area;
- $\beta$ is the model's parameter;
- $t$ is the time, $t \in (0, \Delta t)$, $\Delta t$ is the time interval considered;
- $j$ is the index of the day, $j \in (1, N)$, $N$ is the total number of days considered.

The mean daily streamflow value $\langle Q \rangle$ is then given by:

$$\langle Q \rangle = \frac{Q(0)}{\beta \cdot \Delta t} \left(1 - e^{-\beta \cdot \Delta t}\right) + q_{e,j} - \frac{q_{e,j}}{\beta \cdot \Delta t} \left(1 - e^{-\beta \cdot \Delta t}\right) \tag{2}$$

## 4 Network Methodologies

Two network typologies (MLP and JNN) have been challenged to reproduce the target runoff datum and particulary the beginning of the rise in the flow hydrograph and the decay curve where "contradictory informations" occur, namely, when precipitation has ended and flow is yet decreasing.

### 4.1 MLP and Tapped Dealyed Lines (TDL)

A standard two layered feed-forward neural network with sigmoidal activation function has been considered for modeling the R-R data; conjugate gradient optimization has been adopted to train the ANN. The prime focus regards the input informations, or rather, the number of antecedent daily precipitations necessary for the hydrograph's reconstruction. It is evident that a single rainfall value is not sufficient to reproduce the runoff measure, since the same (null) input value would be associated with different output values. Therefore, a collection of past rainfall measures has to be fed to the ANN to catch the variability of the physical system. This can be realized by employing the Tapped Delayed Line (TDL) approach by which the last $m$ values $x(t), x(t-\tau), \ldots, x(t-(m-1)\tau)$ of a signal $x(t)$ are simultaneously presented at the input of a network, Fig. 1.

Herein a single input serie $x(t)$ is shown. $x(t), x(t-\tau), \ldots, x(t-(m-1)\tau)$ are fully connected to a hidden layer, where $\tau$ stands for the temporal elementary unit backwards. In the first row input information is loaded, in the second one (hidden layer) three processing elements occur to produce a single output: $y(t)$.

**Fig. 1.** A time-delay neural network

## 4.2 MLP Results: Controlled Experiments

An undesired insuperable discrepancy arises between the two compared signals, showed by the appearance of continuous and persistent plateaus in network's modeled daily runoff, which make MLP unable to reproduce the behavior of target, even if a large number of previous rainfall and network's calculated runoff values have been included in the TDL. Figure 2 shows a yearly extract of the training set.

To carry out simulation a feed-forward MLP with six input informations (features) was selected. The input features include: current precipitation value, two tapped delayed lines of three previous precipitations and two network's calculated runoff data. Three nodes in the hidden layer have been introduced, thus obtaining a mean square error ($MSE$) of the residuals ($Q_i - \hat{Q}_i$), (through the functional dependence (3)) equals to: 1.068 e-03.

$$MSE = \frac{1}{N} \sum_{i=1}^{N} \left( Q_i - \hat{Q}_i \right)^2 \tag{3}$$

## 4.3 Jordan Neural Networks

A more remarkable way to recognize and reproduce temporal sequences consists in using recurrent neural networks (RNN) herein invoked for their ability to model a sort of system memory or inertia by considering cues of a recent past back into the input layer and enable, furthermore, a simple and more incisive performance than traditional tapped-delayed lines. Since training a

**Fig. 2.** Bad reconstructions of low values and recession curves

complete RNN involves a high computational cost and can lead to stability problems, simplified models have been proposed in the literature, which can be trained by conventional back-propagation. Elman (1990), Kohonen (1989) and Jordan (1986, 1989) introduced three models of this kind. In our work, Jordan neural networks (JNN) have been considered; here, recurrency is accomplished through the following updating rule [6]:

$$X_i(t) = \alpha \cdot X_i(t-1) + Y_i(t-1) \tag{4}$$

where $X_i$ is the *i-th* component of the additional vector $X$, built at given time $t$ from previous outputs $Y_i$ and $\alpha$ is the strength coefficient to be chosen by the user. Figure 3 schematizes the followed procedure.

## 4.4 JNN Results: Controlled Experiments

Our procedure entails two different approaches: one with "rainfall memory effect" where $Y_i$ (through the functional dependence (4)) are rainfall data and $\alpha_1$ is the strength coefficient, and another one with runoff memory effect with $\alpha_2$ coefficient, where $Y_i$ are network's calculated discharges; and traditional recurrent procedure occurs. Figure 4 shows a yearly extract of the training set.

A JNN with three feautures has been introduced. Inputs include: rainfall and runoff memory effects and previous rainfall datum. Three nodes in the hidden layer have been introduced, thus obtaining a $MSE$ equals to: 4.752 e-05.

**Fig. 3.** Continuous arrows represent connections from the *i-th* unit of the output layer to the *i-th* unit of the input layer, while *dashed arrows* represent fully connected layers. Alpha stands for the strength's memory connection to be searched



**Fig. 4.** Good reproductions of "decay curves" obtained with rainfall memory effect; $\alpha_1$ equals to 0.75 and $\alpha_2$ equals to 0.72

## 5 Analysis of Physical Context

An harder task was to reproduce real data, which are affected by noise and require a more scrupulous exploration of suitable rainfall strength coefficient ($\alpha_1$), expecially for the intermittent behavior of low values signal. Alphas have been searched with standard gradient algorithm during training.

### 5.1 Results: Physical Context

In both simulations "rainfall memory effect" resulted of noticeble major importance, compared to the remaining input informations. Therefore, calculated memory effect was set aside and $\alpha_2$ coefficient was considered constant for the two different hydrographs reconstructions, Figs. 5 and 6.



**Fig. 5.** A three years simulation sequence obtained from five years training data: Argentina River

Some spurious peaks arise, both in Figs. 5 and 6, due to sporadic rainfall events that lead to excessive rises in modeled streamflow data.

Simulations have been performed with rainfall, runoff memory effects, a tapped delayed line of two previous rainfall data and seasonality, that is, the current day in Julian calendar. For both rivers four nodes have been introduced in the hidden layer, thus obtaining, respectively for Argentina

**Fig. 6.** Idem: Impero River

River: $\alpha_1$ equals to 0.82, $\alpha_2 0.70$ and $MSE$ equals to 7.489 e-05; for Impero River: $\alpha_1$ equals to $0.75, \alpha_2 0.70$ and $MSE$ equals to 5.120 e-05.

## 6 Comments

The results pointed out, in both cases, the major importance played by the "rainfall memory effect" that allow the network to remember cues from the recent past, thus reproducing the "runoff decay curves" with high accuracy.

This is a consequence of the fact that rainfall values are provided externally and, therefore, are not affected by errors due to the reconstruction process (like the runoff value). Good results have been obtained when streamflow targets to be reproduced ensue from controlled experiments. Similar performances, moreover, can be achieved when dealing with measured target streamflows, provided that the rainfall data at some previous day are given as input to the JNN, besides weighted sums for memory effects.

## References

1. Anctil F., Michel C., Perrin C., A.V., 2004. A soil moisture index as an auxiliary ANN input for streamflow forecasting, J. Hydrol., 286,155–167.
2. ASCE Task Committee on Application of the Artificial Neural Networks in Hydrology, 2000a. ANN in Hydrology I: preliminary concepts. J. Hydrol. Engng., ASCE 5 (2) 115–123.

3. Campolo M., Soldati A., Andreussi P., 1999. River flood forecasting with a Neural Network model, Water Resources Research, 35, 1191–1197.
4. Cancelliere A., Ancarani A., Giuliano G., Rossi G., 2000. Determinazione delle regole di esercizio di un serbatoio tramite reti neurali, XXVII Convegno di Idraulica e Costruzioni Idrauliche.
5. Laio F., Porporato A., Revelli R., Ridolfi L., 2003. A comparison of nonlinear flood forecasting methods, Water Resources Research, Vol. 39.
6. Hertz J., Krogh A., Palmer G., Introduction to the theory of Neural Computation, 1991. Addison-Wesley Publishing Company, Redwood City, California.
7. Lin G.-F., Chen L.H., 2004. A non linear rainfall-runoff model using radial basis function network, J. Hydrol., 289, 1–8.
8. Minns A., Hall. M., 1996. Artificial Neural Networks as rainfall-runoff models, Hydrol. Sci. J. 41(3), 399–417.
9. Rajurkar M., Kothyari U., Chaube U., 2004. Modeling of the daily rainfall-runoff relationship with artificial neural network, J. Hydrol., 285, 96–113.
10. Sjöberg J., Zhang Q., Lennart L., Benveniste A., Deylon B., 1995. Nonlinear Black-box in system Identification: a Unified Overview, Automatica, Vol. 31, 1691–1724.
11. Thirumalaiah K., Makarand Deo, 2000. Hydrological forecasting using Neural Networks, Journal of Hydrologic Engineering, 180–189.
12. Zhu M.L., Fujita M., 1994. Comparisons between fuzzy reasoning and neural network methods to forecast runoff discharge. Journal of Hydroscience and Hydraulic Engennering 12(2), 131–141.

# Distributed Data Acquisition System for Environment Monitoring Nonlinear Processes

G. Costache

University "POLITEHNICA" of Bucharest, Faculty of Automatic Control and Computers

**Abstract.** The main objective of this paper is to describe a solution to realize on-line connections between remote workstations placed into different locations. Unidirectional or bidirectional informational transfer has to be guaranteed. Local PCs network or Internet are considered as transmission support. The soft component represented by a LabView virtual instrument will coordinate and control the target communication.

This paper aims to proof the concept of remote control of sensors, for automated monitoring of ecological processes. Using National Instrument LabVIEW framework, we realized an application that monitor/command using the same server, different ecological processes.

Real-world data logging applications are typically more involved than just acquiring and recording signals, typically involving some combination of online analysis, offline analysis, display, report generation, and data sharing. In this paper we present a virtual instrument (VI) that converts a PC into a data logger and also the same instruments are used to offline analysis of measured and recording data and to sharing data with another application like Excel.

**Keywords:** Environmental processes, data transfer, monitoring, remote command, remote control, DataSocket communication.

## 1 Introduction

Remote communication between the place from where the parameters are monitoring to a control and command center represent a modern aspect in measurement field. Data transmission configurations are depending by the application requirements. The presence of a Reception Centre is no obligatory. If the LAN is able to geographical manage the application, we can identify each PC through IP address allocated by network server. Internet will be used if the application geographical surface is growing. In this case, a PC must have a real IP address publish on Internet.

Technological process supervision assumes acquisition, processing and returning the commands for a large number of parameters. At the same time, it is necessary to monitor the investigated parameters, because the history is needed in guaranteeing the performances, but also establishing the causes which lead to malfunctions.

Communication technologies and we refer here to Internet, which is mainly spread, allow data, acquired from long distances, to be concentrated in high performance servers. In the same time, following the temporal evolution of acquired parameters can be done from any computer, connected to Internet. This issue is suited both in technological processes, but also in education process, where a laboratory with high performance acquisition equipments and an Internet connection can be used by other laboratories less equipped.

Basic elements of a data logging and analyzing system are presented in Fig. 1. Acquiring is the process of actually measuring the physical parameters and bringing them into your logging system. Online analysis consists of any processing done to the data while are acquiring. It includes alarms, data scaling, and sometimes control, among others.

Logging is an obvious requirement of every data logging system. Offline analysis is everything done with the data after it has been acquired in order to extract useful information from it. The final functional block is made up display, reporting, and data sharing. These are all the "miscellaneous" requirements that fill out the functionality of a data logging system [6].



**Fig. 1.** Basic elements of a data logging and analysis system

## 2 Why Data Acquisition Through Internet

In many of industrial processes, it is preferred the control at a distance of the equipments based on the received parameters. Sometimes, this process is absolutely necessary because of some high risk elements, or because of the conditions that do not allow the user to access it. As an example, one can mention the atomic piles where because of the high degree of radiations, the watch and control of the parameters must be done at a distance. Another example can be the watch of the furnaces, especially in the foundry works

where because of the high temperatures, the user cannot physically stay in those areas.

More than this, an advantage of the watch/control at a distance is the possibility to watch more over than one equipment with a single acquisition server. This aspect allows the taking over of the data from more sensors connected to various industrial equipments and also the transmission of the controls to them. The process can be developed from more than one host in the same time. Thus, there can be a client host that commands through Internet the parameters' modification of an installation, relying on a pressure sensor, while another client modifies the working characteristics of a thermal system, both sensors being connected to the same server.

To sum up, one can say that the advantages of control at a distance of the ecology processes, relying on the watch/control of the parameters are: people should not have to work in the risk areas; the possibility to control more than one equipment in the same time from different hosts; the possibility to watch from a larger number of stations, without extra investments. It should be mentioned that the control of a sensor can be done by a single terminal, while the watch can be developed from many such terminals; the possibility of making some statistics in real time of the monitoring equipment's working.

## 2.1 Virtual Instrument Function

Communication module can be realized as application thought interconnecting of data control with subVIs that allow communication of PC with the target monitored and controlled technological process. Remote data transfer is necessary on long term and low frequency applications, i.e., environment parameters control and monitoring, or dangerous remote operations. In Fig. 2. is presented a remote control and monitoring system for a technological process.

Sensors and traducers $T$ convert electrical or no electrical analogical process data into electrical signals. For PC processing, after passing through Conditioning Block $CC$, the signal must be converted from time continuous variation into a digital time representation, under binary form. This operation is performed by analog-numerical converter $CAN$. Data are taken by PCprocess and processed by virtual instrumentation software. This has two functions: to establish a connection with process via GPIB communication and to realize the communication with remote monitoring and control unit (PCcontrol) via Internet or LAN.

The virtual instrumentation software placed on PCcontrol realize communication with nearby process, but contains processing, recording and visualization modules too. After data interpretation, the commands are send to process by network. In this case, the commands are send by numeric-analogical converter $CAN$ and conditioning circuit $CC$ to actuators $A$.

**Fig. 2.** Transfer, monitoring and control system architecture

# 3 Virtual Instrument

Software is of critical importance in PC-based data logging systems, because well-written logging software determines how data is stored, how quickly data can be written to disk, and how efficiently disk space is used and gives different data management capabilities.

The software use [8], in this paper, to create a PC-based data logger is LabVIEW and also the same LabVIEW programs, that are called virtual instruments or VIs, are used for offline analysis.

Every VI uses functions that manipulate input from the user interface or other sources and display that information and every VI has two components front panel respectively diagram bloc.

The VI that perform both basics function that are data logging and data analysis has two mains front panels ONLINE ACQUISITION respectively OFFLINE ANALYSIS selected by means of the Tab Control.

## 3.1 Acquisition Virtual Instrument

The front panel corresponding to ONLINE ACQUISITION is presented in Fig. 3 and this front panel is used to control acquisition process and also to

**Fig. 3.** Online Acquisition Front Panel

display the important parameters of the acquisition process like acquisition state or measured values and also record date and time. There are three components that are waveform graph in stack plot and two blocks Acquisition Control that control acquisition parameters like parallel port address number of acquisitions, number of points per acquisition respectively Acquisition Progress that show in different form current value acquired. Waveform graph display the filtered and non filtered form of the acquired signal.

To ensure the read bits on data section on the parallel port are used a Sequence Structure where the first 3 sequences are used for parallel port settings [8]. All other program elements are included into a While Loop that repeats the subdiagram inside it until the conditional terminal, an input terminal, receives a particular Boolean value so that acquisition can be interrupt by user.

After port settings is select the file where are storage the measured values and for that are used File Dialog respectively New File. The program can assure one ore more acquisition; number of these is selected by Fig. 3. ON-LINE ACQUISITION Front Panel properly control and number of points of every this acquisition that mean number of values is also selected by properly control, both on the front panel.

## 3.2 Analysis Virtual Instrument

LabVIEW offers many of built-in analysis functions that cover different areas and methods of extracting information from acquired data. With these functions it can be make a virtual instrument (VI) that front panel has two components switching by a proper Tab Control between SAMPLE RECORD and DATA RECORD that mean analysis of every 2000 acquiring points (one sample) or all data representing acquiring points.

The VI corresponding to one sample and to all data is presented in Fig. 4. On that panel are arranged the indicators elements that assure the visual-

**Fig. 4.** Offline Analysis Front Panel (N. Patrascoiu al all, 2005)

ization the results of the data analysis about measured and recorded values. Those elements indicate, at one time with record time the limits values and statistical values. For the limit values are indicate also the time when these limits are accomplished. The statistical values represent the statistic parameters described above that is mean, standard deviation and variance computed Two Waveform Graphs are used to the graphical analysis.

One of these display the graphical evolution of record data i.e. the measured signal and for this exist possibly to setup the lower and upper limits. Another graph display time period for that data is in excess of lower or upper limits. In this way is possible to determine which of the acquired data is into domain that is limited by lower and upper setting limits. It can be possible to use LabVIEW measurement data files into another application. To save measurement data files is use the Write LabVIEW Measurement File Express.vi that generates the proper output files.

The LabVIEW data file is a tab-delimited text file that can open with a spreadsheet application like Microsoft Excel or a text-editing application. The recorded data with graph representation in Excel for the same acquisitions data sequences are presented in Fig. 8. In addition to the data an Express.vi generates, the *.lvm* file includes information about the data, such as the date and time the data was generated.

## 4 The Communication Module

The application has two main parts: the part of data acquisition and data processing; the part of communication.

The module responsible for communication was developed using NI DataSocket Server-Client communication. In order for the application to work, the Data Server application must be running.

The server module of the application writes data in the server, and then every client can read the data. The client which have rights of issuing commands to the server, can also write data in the server. Also, the server application can read command data from the server and apply it to the sensors. Below is presented the module responsible for writing data into the DataSocket server (Fig. 6ab).

# 5 The Software Applications

The development of Internet applications and the possibility to use data acquisition instruments connected to a PC allows the users to make a central system of measurements near to the process and use it to transfer data towards every place where users need [7].

Using a TCP/IP connection between the computers (that includes a Server part and a Client part), it is possible to develop a system of subVIs using LabVIEW environment. This works like a Server-Client structure, allowing the users to perform data acquisition and transfer. Usually, data are transferred to a different place to be processed by a monitoring system. Since the amount of data to be transmitted to the monitoring system can be important, it is suitable to compress them before the transfer is carried on. This allows a faster communication between the monitoring system and the process.

The software application described hereafter integrates two parts. The first part is represented by the LabVIEW environment that helps to perform data acquisition, saving and transfer. The second part represented by the software application responsible with data compression. Both software applications are user friendly.

## 5.1 The LabVIEW Application

The TCP/IP connection uses as server application the diagram depicted in Fig. 6a. The server gives the result of a request on its ports and then a client that reads the port takes the result.

The server blocks are as follows: *String* is the numeric field to send; *Port* is the communication port where the data are received; *TCP Listen* checks for client requests; *TCP Write* writes the data on the communication port; *TCP Close* closes the communication.

The client diagram is drawn in Fig. 7.b)

The blocks are here as follows [7]: *Address* is the address of server, in numeric format; *Port2* is the port of communication where the data are stored (if the data are available); *TCP Read* reads the data if they arrived to the port of communication.

In Fig. 11a it is illustrated how the capture of data in the *\*.txt* files type is realized. The file resides in the same place where as the instrumentation for data acquisition. The compression application can be initiated to compress

**Fig. 5.** Excel Data Representation



**Fig. 6.** (**a**) The DataSocket Server; (**b**) The module the writes data in the DataSocket Server



**Fig. 7.** (**a**) Data server structure; (**b**) Data Client structure

the files before performing the remote transfer to the monitoring system. The received data have to be first decompressed and then processed. The results are shown in graphical mode.

The next Fig. 8b shows an example of *\*.txt* file, created with the *Write To Spreadsheet File.vi* block. The numerical values are encoded by digits, using the ASCII table. Each ASCII code represents a symbol within the original data string to compress.

## 6 Conclusions

The computer has become a friendly instrument that offers more and more features for processing data, display and print the results, anticipatory buffering and saving, connection to the network and sending the data to another

**Fig. 8. (a)** The structure of Write block charged to capture the data; **(b)** An example of data file in ASCII format

computer. All these operations are equally important but nowadays, when many applications use remote data from and towards processes, the efficiency of methods employed to operate data transfers constitutes the milestone of the communication between processing units.

The applications was designed and developed to prove a couple of concepts about the data acquisition in general and some notions about the possibility of adding remote controlling/monitoring. From one point of view one can process the experimental data gathered from a real process, but one can also see the result of one remote command sent to industrial equipment in the real time. The main part is, as we mentioned earlier, the server with the data acquisition board. As main topic of this paper we presented the process of remote controlling of the ecological parameters, but the server can deal with more than one application in the same time. Thus it can monitor parameters provided by 16 sensors and it can lead till four industrial equipments in the same time, using NI Field Point.

The presented bidirectional transfer module allows to data remote transfer via Internet or LAN. This module can be used as subVI into many remote communication applications: environment parameters control and monitoring, or dangerous remote operations.

## References

1. C. Donciu, A. Trandabat, M. Cretu, Servere Tcp-Ip Destinate Teletransmisiilor De Date In Labview, Revista De Instrumentatie Virtuala, Vol. 3, Nr. 4, pp. 83–87, 2001.
2. C. Donciu, C. Fosalau, M. Cretu, Modul de transfer bidirectional al datelor la distanta, Revista De Instrumentatie Virtuala, vol. 4, nr. 3, pp. 68–70, 2001.
3. C. Nitu, F. Krapivin, A. Andro, Sisteme inteligente in ecologie, Ed. Printech, Bucuresti, 2000.
4. R. Munteanu, Aplicatii ale instrumentatiei virtuale in teletransmiterea datelor, Revista De Instrumentatie Virtuala, vol. 3, nr. 3(7), pp. 130–134, 1999.
5. M. Ghercioiu, Conectarea la Web din panoul instrumentului virtual, Revista de Instrumentatie Virtuala, vol. 1, nr. 2, pp. 60–61, 1998.

6.  M. Vlad, I. Rancea, M. Trufas, V. Sgarciu, Acquisition And Monitoring of Process Parameters Using Internet, the 15th International Conference on Control Systems and Computer Science, 2005, Bucharest, Romania.
7.  A. Dumitrascu, H. Mihai, D. Stefanoiu, Data Compression And Transfer Within Environmental Processes, the 15th International Conference on Control Systems and Computer Science, 2005, Bucharest, Romania.
8.  N. Patrascoiu, A. Tomus, Data Logging And Analysis By Virtual Instrumentation, the 15th International Conference on Control Systems and Computer Science, 2005, Bucharest, Romania.

# Automatic Safety Control in Food Processing

R. Furlanetto[1], F. Tassan[1], and M. Toppano[2]

[1] Technology Development Center, Electrolux Professional S.p.A., viale Treviso, 15-33170 Pordenone, Italy

[2] Artificial Perception Laboratories – DEEI, University of Trieste, via Valerio, 10-34127 Trieste, Italy

## 1 Introduction

In food processing, the hygienic issues are the most important issues for the health of the consumer; in particular the microbiological safety of cooked food is also regulated by hygienic legislation.

The undercooked meat, for example, may be a vehicle for many pathogens, consequently it is recommended to cook the foods to reduce the microbial concentration. During thermal processing of foods, the population of microorganisms present in the food decreases depending on the temperature of the product. The population of vegetative cells, such as *E.Coli*, *Salmonella* and *Lysteria monocytogenes*, will decrease in a logarithmic manner. The population of microbial spores will decrease according to a similar pattern but only after an initial lag time.

The automatic supervision of food in terms of hygiene represents a huge advantage for persons who have no longer to take care of safety aspects during the cooking.

With this type of system is behind this idea: an oven that automatically controls the safety of the food in absolutely respect of HACCP recommendations (Hazard Analysis & Critical Control Point, Directive CEE 93/43 which contains the guidelines regarding the hygiene in food preparation, production, storage and distribution).

This is a complex system that provides the automatic validation on the achievement of the microbiological safety during the cooking process in professional and domestic ovens. Our idea consists of novel type of control able to solve the problem of safety with a high quality of cooked products.

## 2 Method Description

The residual bacterial content in a cooked food depends to a substantial extent on the actual time during which a given minimum temperature level is

allowed to persist in the food, i.e. on the same food being allowed to remain at a minimum temperature level for a definite period of time, and such an information is not automatically and readily available in a prior-art cooking appliances.

Since the bacterial content in a food is generally known to vary as a function of both the temperature to which the food itself is exposed and the length of time of such exposure to said temperature, the possibility is given for a function to be plotted that links the reduction in the bacterial content with both these parameters and that, as a result, is representative of the evolution pattern of the same bacterial content in the food.

Figure 1 shows a typical plot of microbial population, over time during cooking process. It is important to define the decimal reduction time $D$ as the time necessary for 90% reduction in the microbial population. The initial microbial population has no influence on the $D$ value since the magnitude is directly related to the difference of values.



**Fig. 1.** Semi-logarithmic plot of microbial population over time for several temperatures

Exposure of the microbial population to higher temperatures results in a decrease in the $D$ value, as shown in Fig. 1 that represents the relation of time, $D$ and population.

$$t = D \cdot (\log N_0 - \log N) \tag{1}$$

The thermal death time $F$ is the time required to cause a reduction in a population of microorganisms or spores and can be expressed as a multiple of $D$ value. Considering $D_{T0}$ value at reference temperature and $n$ number of decimal reduction required it is possible to write (2).

$$F = D_{T0} \cdot (\log N_0 - \log N) = n \cdot D_{T0} \tag{2}$$

Since $D$ changes with temperature, the thermal resistance constant $z$ is a unique factor describing thermal resistance of the bacteria spores and can be defined as the number of degrees below reference temperature at which $t$ increases by a factor 10 ($T_0$ is the reference temperature).

$$z = \frac{T - T_0}{\log F - \log t} \tag{3}$$

Having therefore so defined the function of bacterial content reduction, the possibility arises for a plurality of degrees of known reduction thereof to be defined as well, which correspond to respective values $F_0, F_1, F_2, \ldots F_n$ that such a function $F$ can take. As this will be exemplified further on, to these values there can be associated pre-established periods of preservability under storage conditions (i.e. shelf life) of the food items having been processed. In a few words, use is made of a conventionally established and acknowledged function of the evolution pattern of the bacterial count, i.e. content in a food, and some characteristics of the hygienic state of the same food are identified experimentally along with the aptitude thereof to be kept in store, i.e. preserved before its bacterial content rises again to an unacceptable level.

Knowing the necessary value for commercial sterilization, corrects parameters time and temperature can be calculated by (2) and (3) [4].

If it is assumed that $dS$ is the fraction of the process towards reaching thermal death and this is accomplished in time $dt$, the two entities are correlated by (4), where $t_1$ is the thermal death time at temperature $T_1$, assuming that the destruction is additive.

$$dS = \frac{1}{t_1} dt \tag{4}$$

When the thermal death time has been reached, that means that effective sterilization has been achieved, and for $dS$ can be written the following relationship:

$$\int dS = 1 \tag{5}$$

Thus it is possible to calculate the $F$ parameter using (3), (4) and (5).

$$F = \int_{t_1}^{t_2} 10^{\frac{(T-T0)}{z}} dt \tag{6}$$

Where $t_2$ and $t_1$ represents a general time slot, $T_0$ the reference temperature and $T$ the minimum internal temperature of the product. Usually $t_1$ represents the instant to which there corresponds a temperature in excess of a specified value (e.g. 50°C), $t_2$ instant of final reading and $T_0$ reference temperature (e.g. 71°C).

In our approach, during the cooking process, the value of $F$ is calculated in a continuous manner, and is further compared with values $F_0, F_1, F_2, \ldots F_n$

that are contained in a pre-defined table, in which said values are associated to and characteristic of the kind or category of the food being each time handled. When the integral is greater than to $F_0$ implies that the sterilization process is complete, as the necessary fraction of the bacteria has been destroyed. In this way, the factors $F$ and $z$ can be combined with time-temperature curve and integrated to evaluate a sterilizing process [3].

For example, Fig. 2 shows the core temperature and the $F$ parameter for roast-beef cooking with a set point of $58°C$.



**Fig. 2.** Core Temperature and thermal death time $F$ for roast-beef cooking. Reference pathogen is Lysteria monocytogenes with $z$ equals to 10

## 3 Validation Algorithm

At the beginning of the cooking process, the processing and control device at the same time starts calculating the integral of the formula (6) in view of delivering at each single instant the value of $F$.

At each such instant, said processing and control device automatically compares the value of $F$, as this has so been just calculated, with a set of values $F_0, F_1, F_2, \ldots F_n$ that will have been appropriately pre-defined and stored in the unit's memory.

These values are determined experimentally for each category of risk of food, **A** or **B** (high risk or low risk), in view of identifying the corresponding maximum allowable shelf-life, i.e. the longest storage period before the bacterial content of the food increases again to an unacceptable value. For instance, a determined code, e.g. "SAFE 1", may be associated to a given value of $F$. The same applies to all other pre-determined values of $F$, as this is better exemplified in the table below:

To the code SAFE-0 there will for instance correspond, for each selected category of food, a maximum allowable time of 5 hours before said food is

**Table 1.** Safety conditions

| Value of $F$ | Display | Max. Allowable shelf-life |
|---|---|---|
| $F < F_0$ | UNSAFE | |
| $F > F_0$ | SAFE-0 | 5 hours |
| $F > F_1$ | SAFE-1 | 1 day (refrigerated storage) |
| $F > F_2$ | SAFE-2 | 5 days (refrigerated storage) |
| . . . | . . . | |
| $F < F_n$ | SAFE-n | . . . |

eventually served for consumption; again, to the code SAFE-1 there may correspond a shelf-life of one day if the food is kept under correct refrigerated storage conditions. In this way, at the end of the cooking process is it possible to be immediately informed about the cooked food having a more or less acceptable bacterial content.

Under the above cited circumstance (i.e. $F_{meas} < F_0$), the cooking cycle is able to go on automatically at least until said value of $F$ eventually reaches the value of $F_0$.

## 4 Numerical Analysis

The computation of the formula with a fixed precision architecture requires the quantization of the basic function exponential using a look-up table. The integration is made by a continuous addition of the different values coming from the exponential look-up table.

A probe insertion identification algorithm is required during the cooking start-up and the result of this special software is a warning indicating the probe out of the food. To implement this functionality is necessary to have a multi-point sensor that also gives the possibility to estimate the lowest temperature into the food. It is possible to identify two main not insert probe situations:

1. *probe correctly inserted into the holding support*;
2. *probe free in the cavity*;

The second situation is the worst case because it is difficult distinguishes between probe fully inserted in the food and probe free in the cavity. To solve these problems a knowledge-based algorithm is used, implementing a particular temperature growing pattern recognition [8].

## 5 Conclusions

Modeling a cooking process is a very hard task, in particular for health and safety problems. In our work using the relation between $F$ index and Core

Temperature, we carried out a novel method implemented in embedded system, able to detect in real time the safety of foods during the cooking.

# References

1. M A Hague, K E Warren, M C Hunt, D H Kropf, C L Kastner, S L Stroda, D E Johnson (1994) Endpoint temperature, internal cooked color, and expressible juice color relationships in ground beef patties, Journal of Food Science, vol. 59, no. 3, pp. 465–470
2. Joseph Kerry, John Kerry, D Ledward (2002) Meat Processing, CRC Press, Woodhead Publishing
3. R L Earle (2004) Unit operations in food processing, The New Zealand Institute of Food Science & Technology (Inc.), Web Edition
4. R P Singh, D R Heldman (1993) Introduction to Food Engineering, Academic Press
5. B Kosko (1992) Neural networks and fuzzy system: a dynamical system approach to machine intelligence, Prentice-Hall, Englewood Cliffs
6. R J Mammone, Y Zeevi (1991) Neural Networks theory and applications, Academic Press
7. A Boscolo, C Mangiavacchi, O Tuzzi (1994) Data and models fusion techniques in measurement, XIII IMEKO World Congress, Torino
8. J M C Sousa, U Kaymak (2002) Fuzzy decision making modeling and control, World Scientific Series in Robotics and Intelligent Systems, vol. 27
9. K Hiroda (1993) Industrial Applications of Fuzzy Technology, Springer-Verlag

# Using a TI C6701 DSP Rapid Prototyping System for Nonlinear Adaptive Filtering to Mitigate Interference

R. Goshorn and D. Goshorn

Space and Naval Warfare Systems Center, Code 2373, 53560 Hull Street, San Diego, CA 92152-5001

## 1 Introduction

Designing and implementing nonlinear systems, onto hardware devices experiences a paradigm shift with the innovative rapid prototyping system (RPS). The RPS is low cost, commercially off-the-shelf (COTS), and ingeniously establishes a direct path from initial design to a hardware implementation operating in real-time, eliminating several levels of tedious programming for hardware, with automatic code generation. In this paper, the hardware device is a digital signal processor (DSP) embedded in a prototype board. The RPS extends to real-world hardware testing specific to application, from which the nonlinear system design can be revised and optimized. From the RPS, final nonlinear system hardware can be designed with a high level of confidence, and the prototype board can continue to simulate other nonlinear devices.

This paper explicates an overview of the RPS for nonlinear system hardware development. Section 2 addresses the nonlinear adaptive filter on the RPS for narrowband interference mitigation in communications channels. Section 3 conveys the overall RPS structure, procedure, and challenges overcome. Section 4 yields ideas for RPS expansion and concludes the paper.

## 2 Nonlinear Applications on the Rapid Prototyping System

The RPS can implement any nonlinear system, such as nonlinear control systems and coupled nonlinear differential equations. This section describes the innovative implementation of nonlinear adaptive filter theory using the RPS for interference mitigation in real-world radio systems. In this communications application, the radio fails to detect the desired signal due to narrowband interference. Yet with the nonlinear adaptive filter hardware from the RPS, the radio captures the desired signal at an excellent bit error rate (BER). Thus,

the RPS is used not only to design the nonlinear filter, but also to prove its validity when the hardware implants in real-world applications.

The input, $\mathbf{y}$, to the nonlinear adaptive filter is the received shaped offset quadrature phase-shift keying (SOQPSK) modulated signal, $\mathbf{d}$ [1], with linear chirp narrowband interference $\mathbf{i}$, and additive white Gaussian noise, $\eta$, such that

$$\mathbf{y} = \mathbf{d} + \mathbf{i} + \eta \tag{1}$$

$$\mathbf{i}\,[n] = \sigma_{\mathbf{i}} e^{(j2\pi(f_{\mathbf{i}_o} + \frac{\psi}{2}n)n + \theta)} \tag{2}$$

with $f_{\mathbf{i}_o} = -200\,\mathrm{Hz}$ and chirp rate $\psi = 16\,\mathrm{Hz}$. Interference $\mathbf{i}$ has a uniform random initial phase, $\theta$ from a uniform distribution on the interval [0, 1]. Interference $\mathbf{i}$ is assumed to have zero-mean and

$$\sigma_{\mathbf{i}}^2 = \sigma_{\mathbf{d}}^2 \cdot 10^{(-SIR_{dB}/10)} \tag{3}$$

where the signal-to-interference ratio (SIR) is set to $-10\,\mathrm{dB}$. The BER of the received signal is 0.4947. Therefore, the received signal has no information and thus needs to be filtered. The normalized least mean square (NLMS) adaptive predictor structure is used to mitigate the interference, with the weight adaptation [2]

$$\mathbf{w}_{n+1} = \mathbf{w}_n + \mu \frac{\mathbf{y}_n}{\mathbf{y}_n^H \mathbf{y}_n} \hat{\mathbf{d}}^* \tag{4}$$

$$\hat{\mathbf{d}} = \mathbf{y} - \mathbf{w}_n^H \mathbf{y}_n \tag{5}$$

where $\mathbf{y}_n = [y[n-L]y[n-2L]\dots y[n-M\cdot L]]^T$. The NLMS adaptive filter uses the last M symbols to compute the dynamic weights to predict the interference, where M is the filter order and set to the value of ten. The NLMS step-size $\mu$, is assumed a value of 0.007195. The values for parameters M and $\mu$ were optimized for both bit error rate and convergence time, to yield a near 0.01 BER for chirp interference. Discussion of the method of optimization is found in [3]. The M values of $\mathbf{y}$ used to predict are separated by a minimum tap delay, L with a value of 14, a function of the sampling and symbol frequencies. The NLMS adaptive filter has been shown to exploit nonlinear behavior, surpassing the traditional Wiener filter solution in several applications, as explained in [4], which cites further references.

The NLMS filter weights behave nonlinearly, as the weights dynamically oscillate as a result of the interference contaminated signal. As the real and imaginary parts of the weights nonlinearly oscillate, the corresponding magnitude of the NLMS adaptive filter output converges.

## 3 Rapid Prototyping System Structure

The RPS aptly designs, tests, and validates any nonlinear system on hardware in real-world applications, entailing a high-level design only. The overall software and hardware system structure can be seen in Fig. 1.

**Fig. 1.** RPS software and hardware structures (**A**) RPS Steps 1, 2 and (**B**) RPS Step 3

The required COTS software is MathWorks' MATLAB®, Simulink®, Real Time Workshop®, Link to Code Composer Studio™, Data Acquisition Toolbox, and Embedded Target for Texas Instrument (TI) C6000™ DSP. MATLAB® is software for design, analysis, and data visualizations. Simulink®, an add-on to MATLAB® , provides the environment for high-level system modeling, simulation, and validation. Real Time Workshop® provides the capability to generate automatic C code equivalent to the high-level Simulink® model. The Link to Composer Studio™ links MATLAB® 's generated C code to TI's Code Composer Studio™, where it is then compiled into assembly code for the DSP [5]. The Embedded Target for TI C6000™ provides the interface between the hardware and MATLAB®/Simulink® software [6]. Finally, the Data Acquisition Toolbox allows for data exchange between MATLAB® and the RPS board [7].

The required RPS prototype board for the applications in this paper is the Texas Instruments TMS320C6701 EVMC6701 Evaluation Module (EVM) board which embeds the C6701 DSP. This board is one full-length peripheral component interface (PCI) slot, allowing the use of the C6701 DSP internally in your PC or externally on real-world hardware applications or on PC laptops. If used externally with a PC laptop, as in the nonlinear filter application, the EVM board must connect and communicate with the laptop in analog form through audio cables, or digitally through RTDX (Real-Time Digital Exchange) cables. Both the PC and EVM board embed analog-to-digital and digital-to-analog converters supporting a range of sampling frequencies. The prototype hardware has the Crystal CS4231A-KL codec while the PC utilizes Analog Devices AD1981A AC'97 SoundMAX codec [7,8].

The following encapsulates the RPS in three steps, labeled in Fig. 1. The first step is to create the high-level, block-diagram model of a nonlinear system in Simulink® software. Test and verify the model in the Simulink® environment. The second step happens with only a click of a button in Simulink®, in which then the equivalent ANSI/ISO C-code equivalent to the model is automatically created and compiled into assembly code which is assembled into machine language for the C6701 DSP. Upon completion of code generation, the executable file downloads onto the DSP and executes the nonlinear systems which were originally designed in Simulink® [5, 6]. The third step is the signal exchange between the PC and EVM using MathWork's Data Acquisition Toolbox. The digital signal generated in MATLAB® is modulated, converted to analog, and transmitted to the prototype hardware, where the prototype's codec converts the signal back to digital for demodulation, filtering, and modulation back by the C6701 DSP. The prototype's codec converts the filtered, modulated signal back to analog and is transmitted to the PC for post-processing and analyzing. These three steps of the nonlinear adaptive filter RPS are automated through a custom-made graphical user interface (GUI), written in MATLAB®.

During the RPS development, there were challenges to overcome. For one, Simulink® models targeted to the EVM board must operate in frame-based processing, defaulted to 64 samples per frame [6]. However, the automatically generated C code would only recognize the default. This was solved by editing two lines of C code. This showed the flexibility of changing the automatically generated code if needed to. Another issue occurred when the RPS was unable to handle a signal of 25 kHz bandwidth, because the PC codec, with maximum sampling rate of 48 kHz, is incapable of sampling at the Nyquist rate of 50 kHz [1]. However, appropriate hardware will be purchased to upgrade the RPS for faster sampling frequencies. Another technicality arose when trying to use the sampling frequency of 48 kHz on the PC's codec, because the Interrupt Service Routine (ISR) would overrun its current application. In general, there was some pioneering required with MathWorks and TI, as it seemed this was one of the earliest rapid prototyping systems implementations. Hopefully, future customers will not have this pioneering.

# 4 Conclusions

The revolutionary RPS is low-cost, COTS, and can be used to test and validate any nonlinear system implanted in the desired real-world application, with no low-level programming (C & below) required. The nonlinear system algorithms can be revised in Simulink® based upon the real-world performance and the RPS run-time software is automatically updated and revised. At any time, interactive debugging and testing of the nonlinear adaptive algorithms can be carried out and immediately executed on the RPS hardware. The RPS can be expanded to include field programmable gate arrays (FPGAs)

which can be automatically updated with Simulink®. The nonlinear models for implementation can also be partitioned and exported onto a DSP and FPGAs [9]. After validating the nonlinear algorithm's performance on the hardware when entrenched in the real world, the automatically generated C or VHDL/Verilog code can be extracted and further optimized for follow-on final hardware designs at any time. In conclusion, this pioneering RPS implements nonlinear systems in the three steps shown, closing the costly and timely gap between traditional software development and implementation on a variety of hardware devices for real-world applications and validation.

# References

1. Proakis, John G. *Digital Communications, $3^{rd}$ ed.* NJ: McGraw-Hill, 1995.
2. Haykin, Simon. *Adaptive Filter Theory, $3^{rd}$ ed*, NJ: Prentice Hall, 1996.
3. Goshorn, Rachel. "Syntactical classification of extracted sequential spectral features adapted to priming selected interference cancellers." PhD Thesis, University of California, San Diego, 2005.
4. Haykin, Simon and Bernard Widrow. *Least-Mean-Square Adaptive Filters, Chap. 9: Steady-State Dynamic Weight Behavior in (N)LMS Adaptive Filters.* NJ: John Wiley & Sons Inc., 2003.
5. *Real Time Workshop® User's Guide V6*: Natick, MA. MathWorks Inc 2005
6. *Embedded Target for the TI TMS320C6000$^{TM}$ DSP Platform User's Guide, V2*: Natick, MA. MathWorks, Inc., 2005.
7. *Data Acquisition Toolbox User's Guide V2*: Natick, MA. MathWorks, Inc 2005.
8. *IBM ThinkPad T41 Access Help User's Guide*, <http://www.ibm.com>.
9. *Link for ModelSim® User's Guide, V1*: Natick, MA. MathWorks, Inc 2005.

# Gunn Oscillations Described by the MEP Hydrodynamical Model of Semiconductors

G. Mascali[1], V. Romano[2], and J.M. Sellier[2]

[1] Dipartimento di Matematica, Università della Calabria, Via Ponte Bucci, cubo 30 b, 87036 Arcavacata di Rende (Cs), Italy, and INFN-Gruppo G. di Cosenza
`g.mascali@unical.it`

[2] Dipartimento di Matematica ed Informatica, Università di Catania, Viale A. Doria 6, 95125 Catania, Italy
`romano@dmi.unict.it, sellier@dmi.unict.it`

## 1 The Model

High-field phenomena in submicron electron devices cannot be described satisfactorily within the framework of the *drift-diffusion models* that do not include energy as a dynamical variable and are valid only in the quasi-stationary limit, while most hydrodynamical models suffer from serious theoretical drawbacks due to the *ad hoc* treatment of the closure problem [1]. Here we employ a moment approach, previously introduced in [2,3] (see also [4] for a complete review) in which the closure procedure is based on the maximum entropy principle while the conduction bands are described by the Kane dispersion relation. The electrons in GaAs are considered as a mixture of two fluids, one representing the electrons in the $\Gamma$-valley and the other the electrons in the four equivalent $L$-valleys. The model comprises the balance equations of electron density, energy density, velocity and energy flux for both populations, coupled to the Poisson equation for the electric potential.

We will give only a brief sketch of the model. For more details the interested reader is referred to [5].

One assumes that the conduction band is described in the neighborhood of each minimum (valley) by the Kane dispersion relation approximation

$$\mathcal{E}_A(k_A)\left[1 + \alpha_A \mathcal{E}_A(k_A)\right] = \frac{\hbar^2 k_A^2}{2m_A^*}, \quad \mathbf{k}_A \in \mathbb{R}^3, \quad A = \Gamma, L \,, \tag{1}$$

involving a parameter $\alpha_A$, called the non-parabolicity factor. $\mathcal{E}_A$ is the electron energy in the $A$-valley measured from the bottom of the valley $\overline{\mathcal{E}}_A$, $m_A^*$ the electron effective mass, $\hbar \mathbf{k}_A$ the *crystal momentum*, $k_A$ its modulus, and $\hbar$ the reduced Planck constant. The values of $\alpha_A$ and the other physical parameters are reported in Table 1 of [5].

At a kinetic level, the system is described by two Boltzmann equations, one for the $\Gamma$-valley and the other for one $L$-valley. The macroscopic balance equations are deduced by taking the moments of the Boltzmann transport equations, that is by multiplying each transport equation by suitable weight functions $\psi(\mathbf{k}_A)$ and integrating over $\mathbb{R}^3$. If we consider the set of weight functions $1$, $\hbar\mathbf{k}_A$, $\mathcal{E}_A$ and $\mathcal{E}_A\,\mathbf{v}_A$, where $\mathbf{v}_A(\mathbf{k}_A) = \frac{1}{\hbar}\nabla_{\mathbf{k}_A}\mathcal{E}_A(\mathbf{k}_A)$ is the electron group velocity, we get the following macroscopic balance equations

$$\frac{\partial n_A}{\partial t} + \frac{\partial(n_A V_A^i)}{\partial x^i} = n_A\,C_{n_A}\,, \tag{2}$$

$$\frac{\partial(n_A P_A^i)}{\partial t} + \frac{\partial(n_A U_A^{ij})}{\partial x^j} + n_A q E^i = n_A C_{P_A^i}\,, \tag{3}$$

$$\frac{\partial(n_A W_A)}{\partial t} + \frac{\partial(n_A S_A^j)}{\partial x^j} + n_A q V_A^j E_j = n_A C_{W_A}\,, \tag{4}$$

$$\frac{\partial(n_A S_A^i)}{\partial t} + \frac{\partial(n_A F_A^{ij})}{\partial x^j} + n_A q E_j G_A^{ij} = n_A C_{S_A^i}\,, \tag{5}$$

where $n_A$ is the electron density, $V_A^i$ the average electron velocity, $W_A$ the average electron energy, $S_A^i$ the average energy flux, $P_A^i$ the average crystal momentum, $U_A^{ij}$ the average crystal momentum flux, $G_A^{ij}$ a coupling term with the electric field $E^i$, $F_A^{ij}$ the average flux of energy flux, $C_{n_A}$ the density production, $C_{P_A^i}$ the crystal momentum production, $C_{W_A}$ the energy production, $C_{S_A^i}$ the energy flux production. All these terms have a clear definition in kinetic theory (see [5]) and refer to electrons in the A-valley, $A = \Gamma, L$. $q$ is the absolute value of the elementary charge.

All the most relevant scatterings for GaAs, that is those between electrons and intervalley non-polar optical phonons and intravalley polar and acoustic phonons, are taken into account. The electron-electron scattering and scattering of electrons with ionized impurities are neglected.

The above system is coupled to the Poisson equation

$$\mathbf{E} = -\nabla\Phi, \quad \epsilon\Delta\Phi = -q(N_+ - N_- - n_\Gamma - 4n_L) \tag{6}$$

where $\Phi$ is the electric potential, $N_+$ and $N_-$ the donor and acceptor densities respectively, and $\epsilon$ the dielectric constant.

These moment equations do not constitute a set of closed relations because of the fluxes and production terms. Therefore constitutive assumptions must be prescribed.

If we assume as fundamental variables $n_A$, $V_A^i$, $W_A$ and $S_A^i$, $A = \Gamma, L$, which have a direct physical interpretation, the closure problem consists of expressing $U_A^{ij}$, $F_A^{ij}$ and $G_A^{ij}$ and the moments of the collision terms $C_{n_A}, C_{P_A}^i$, $C_{W_A}$ and $C_{S_A^i}$ as functions of $n_A$, $V_A^i$, $W_A$ and $S_A^i$.

The Maximum Entropy Principle (hereafter MEP) leads to a systematic way of obtaining constitutive relations on the basis of information theory. According to the MEP, the distribution functions $f_A^{ME}$ which can be used to

evaluate the unknown moments of $f_A$, are stationary points of the electron entropy functional under the constraint of fixed fundamental variables.

This procedure has been used in [5] upon the ansatz of small anisotropy for the $f_A^{ME}$. Formally a *small* anisotropy parameter $\delta$ has been introduced and explicit constitutive equations have been obtained for fluxes and production terms up to the first order in $\delta$. Their explicit expressions are given in [5].

## 2 Simulations of Gunn Oscillations

One can prove that the system (2)–(5) closed with the MEP is hyperbolic in the physically relevant region of the dependent variables [6] and it is well known that the solutions of quasi-linear hyperbolic systems suffer loss of regularity (e.g. formation of shocks). In the last decades several accurate high-order shock capturing numerical schemes have been developed. Most schemes are based on upwind methods and require the solution to the Riemann problem. Unfortunately no analytical solution to the Riemann problem for the model under investigation is available at the present time and an approach based on the full numerical evaluation of the Roe matrix is not practical. Therefore we have resorted to a central differencing scheme. The central schemes known in the literature deal almost exclusively with homogeneous systems. In [7, 8] a suitable extension for one-dimensional balance laws with (possibly *stiff*) source terms has been developed on the basis of the Nessyhau and Tadmor scheme [9] for homogeneous hyperbolic systems.

The complete method is based on a second-order splitting technique which separately solves the system with the source put equal to zero (convection step) and the one with the flux put equal to zero (relaxation step).

Each convective step has the form of a predictor-corrector scheme on a staggered grid. The scheme is second order accurate both in time and space for homogeneous systems. The interested reader is referred to [10] for a complete review of the numerical aspects related to the present paper.

We consider a GaAs diode coupled to an RLC tank circuit which stimulates Gunn oscillatory effects, used for the generation of microwaves. The one-dimensional diode has length $L_d = 2\,\mu\text{m}$ and its doping profile is

$$N_+(x) = \begin{cases} 10^{17} & \text{for} & x < 0.125\,\mu\text{m} , \\ 10^{16} & \text{for} \quad 0.125\,\mu\text{m} < x < 0.15\,\mu\text{m} , \\ 0.5 \times 10^{16} & \text{for} \quad 0.15\,\mu\text{m} < x < 0.1875\,\mu\text{m} , \quad \text{(donors/cm}^3) \\ 10^{16} & \text{for } 0.1875\,\mu\text{m} < x < 1.875\,\mu\text{m}, \\ 10^{17} & \text{for} \quad 1.875\,\mu\text{m} < x . \end{cases}$$

$$(7)$$

with abrupt junctions. As initial conditions we take the equilibrium state while the potential $\Phi(L_d)$ is set either equal to 2V or it is determined by coupling the device to the system of ODE which models the circuit

$$\frac{dV_d}{dt} = \frac{1}{C}\left(I - I_d - \frac{V_d}{R}\right) , \quad \frac{dI}{dt} = \frac{1}{\Lambda}\left(V_B - V_d\right) , \tag{8}$$

where $V_d$ is the voltage through the device, $V_B$ the bias voltage of the circuit, and $I_d$, the particle current in the device, calculated as

$$I_d = -\frac{q\,A}{L_d}\int_0^{L_d}\left(n_\Gamma v_\Gamma + 4\,n_L v_L\right)dx . \tag{9}$$

Finite difference discretization of (8) and (9) allows the diode voltage to be updated at each simulation time step. The values used for the capacitance, $C$, resistance, $R$, and inductance, $\Lambda$, of the circuit are

$$C = \left(\epsilon A/L_d + 0.82 \times 10^{-12}\right)F , \quad R = 25\,\mathrm{ohm} , \quad \Lambda = 3.5 \times 10^{-12}\,\mathrm{henry} ,$$

where the cross-sectional area, $A$, of the diode is $A = 1.0 \times 10^{-3}\,\mathrm{cm}^2$.

The oscillator (8) are given the initial state

$$V_d(t_0) = 2\mathrm{V} , \quad I(t_0) = 0 . \tag{10}$$

The circuit is engaged at time $t_0 = 75\mathrm{ps}$ when the GaAs diode is judged to have reached the steady state illustrated in Figs. 1, 2. One observes, Fig. 3,



**Fig. 1.** Electron velocity vs position. *** MC simulation, continuous line hydrodynamical model

**Fig. 2.** Electron energy vs position. The notation is as in Fig. 1



**Fig. 3.** The potential $V_d$ versus time for the Gunn diode. The initial time coincides with the instant when the circuit is engaged. The notation is as in Fig. 1

that there are some initial oscillations in the electric potential $V_d$ that smooth out and become negligible after about 200 ps.

In order to test the accuracy of the model and the robustness of the numerical scheme, we have also simulated the same case with the Monte Carlo method by using the code ARCHIMEDES [11]. There is a good agreement between the MC simulation and that of the MEP hydrodynamical model and in particular the behavior of the electric potential is excellent.

## Acknowledgments

## References

1. A. M. Anile and O. Muscato, *Improved hydrodynamical model for carrier transport in semiconductors*, Phys. Rev. B 51 (1995), pp. 16728–16740.
2. A. M. Anile and V. Romano, *Non parabolic band transport in semiconductors: closure of the moment equations*, Cont. Mech. Thermodyn., 11 (1999) pp. 307–325.
3. V. Romano, *Non parabolic band transport in semiconductors: closure of the production terms in the moment equations*, Cont. Mech. Thermodyn., 12 (2000) pp. 31–51.
4. A. M. Anile and V. Romano, *Hydrodynamical modeling of charge transport in semiconductors*, Meccanica 35 (2000) pp. 249–296.
5. G. Mascali and V. Romano, *Hydrodynamical model of charge transport in GaAs based on the maximum entropy principle*, Cont. Meh. Thermodyn. 14 (2002) pp. 405–423.
6. G. Mascali and V. Romano, *Simulation of Gunn oscillations with a non-pabolic hydrodynamical model based on the maximum entropy principle*, COMPEL 25 (2005) pp. 35–54.
7. F. Liotta, V. Romano and G. Russo, *Central schemes for systems of balance laws*, International Series of Numerical Mathematics, 130 (1999) pp. 651–660.
8. F. Liotta, V. Romano and G. Russo, *Central schemes for balance laws of relaxation type*, SIAM J. Num. Analysis 38 (2000) pp. 1337–1356.
9. H. Nessyahu and E. Tadmor, *Non-oscillatory central differencing for hyperbolic conservation law*, J. Comp. Physics 87 (1990) pp. 408–463.
10. A. M. Anile, N. Nikiforakis, G. Russo and V. Romano, *Discretization of Semiconductor Device Problems*, Chap. 5 of **Handbbok of Numerical Analysis** Vol. XIII, Elsevier (2005).
11. GNU package Archimedes, freely available at the web-site: www.gnu.org/software/archimedes

# Dynamic Test Data Generation for the Nonlinear Models with Genetic Algorithms

A. Dobrescu

"Politehnica" University of Bucharest, Str. Spl. Independenei 313, Bucharest Code
74206 Romania
sabda23@yahoo.com

**Abstract.** For nonlinear dynamic systems, the classical models are not sufficiently accurate, because the parameters are poorly known and are in general time-variants. So, it is important to develop control systems that incorporate learning capabilities in a way that their control systems automatically improve accuracy in real time and become more autonomous. This paper presents different technique used in dynamic nonlinear applications like dynamic test data generation and genetic algorithms. One example is given: an automatic pilot.

## 1 Introduction

The software programs become an increasingly critical component in all domains. The biggest problem is that the smallest bug in software program can affect exceedingly all the system. Therefore the test of the software is decisive for the success of the entire system. Software testing accounts for 50% of the total cost of software development. This cost could be reduced if the process of testing is automated. Because of the complexity of the software programs, the rules are almost impossible to be found. In software testing, it is desirable to find test inputs that exercise specific program features. To find these inputs by hand is extremely time-consuming, especially when the software is complex. One method to make that automatically is test data generation paradigm commonly known as dynamic test data generation. Dynamic test data generation was originally proposed by Miller and Spooner in 1976 [8] and then investigated further with the TESTGEN system of Korel (1990) – [6], (1996) – [7], the QUEST/Ada system of Chang et al. (1996) [5] and the ADTEST system of Gallagher and Narasimhan (1997) [9]. This paradigm treats parts of a program as functions that can be evaluated by executing the program and whose value is minimal for those inputs that satisfy test adequacy criterion such as code coverage. In this way the problem of generating test data reduces to the better understood problem of function minimization. But for more complicated software the performance of random test data generation

deteriorates. In this situation, genetic algorithms perform considerably better [4, 10].

The same situation is in the automatic control, where a controller has to be identified. The purpose of a controller is to force, in a meritorious way, the actual response of a system conventionally called the *plant* to match a desired response called the *reference signal*. There are classical techniques for designing the parameter values and the topology of controllers. The ubiquitous type of controller is PID controller. The PID controller was patented in 1939 by Albert Callender and Allan Stevenson of Imperial Chemical Limited of Northwich, England [11]. The design process for controllers today is generally channeled along lines established by existing analytical techniques (notably those that lead to a PID-type controller). It would be desirable to have an automatic system for synthesizing creating the design of a controller that was open-ended in the sense that it did not require the human user to prespecify the topology of the controller whether PID or other, but, instead, automatically produced both the overall topology and parameter values directly from a high-level statement of the requirements of the controller.

However, there is no preexisting general-purpose analytic method for automatically creating a controller for arbitrary linear and nonlinear plants that can simultaneously optimize prespecified performance metrics such as minimizing the time required to bring the plant output to the desired value as measured by, say, the integral of the time-weighted absolute error, satisfy time-domain constraints involving, say, overshoot and disturbance rejection, satisfy frequency domain, constraints e.g., bandwidth, and satisfy additional constraints, such as constraints on the magnitude of the control variable and the plant's internal state variables. One method for creating an automatic controller using genetic algorithm will be shown further for the system.

## 2 Genetic Algorithms

Evolutionary computation is the name given to a collection of algorithms based on the evolution of a population toward a solution of a certain problem. These algorithms can be used successfully in many applications requiring the optimization of a certain multi-dimensional function. The population of possible solutions evolves from one generation to the next, ultimately arriving at a satisfactory solution to the problem. These algorithms differ in the way a new population is generated from the present one, and in the way the members are represented within the algorithm. Three types of evolutionary computing techniques have been widely reported recently. These are Genetic Algorithms (GAs), Genetic Programming (GP) and Evolutionary Algorithms (EAs).

Genetics Algorithms were envisaged by Holland [2] in the 1970s as an algorithmic concept based on a Darwinian-type survival-of-the-fittest strategy, where stronger individuals in the population have a higher chance of creating

an offspring. A genetic algorithm is implemented as a computerized search and optimization procedure that uses principles of natural genetics and natural selection. The basic approach is to model the possible solutions to the search problem as strings of ones and zeros. Various portions of these bit-strings represent parameters in the search problem. If a problem-solving mechanism can be represented in a reasonably compact form, then GA techniques can be applied using procedures to maintain a population of knowledge structure that represent candidate solutions, and then let that population evolve over time through competition (survival of the fittest and controlled variation). However, too strong a bias towards the best individuals will result in them dominating future generations, thus reducing diversity and increasing the chance of premature convergence on one area of the search space. Conversely, too weak a strategy will result in too much exploration, and not enough evolution for the search to make substantial progress [1].

Genetic algorithms are most appropriate for optimization type problems, and have been applied successfully in a number of automation applications including job shop scheduling, proportional integral derivative (PID) control loops, and the automated design of fuzzy logic controllers.

Every solution candidate has to be evaluated because the GP algorithm maximizes or minimizes some objective function. In the context of structure identification this function should be an appropriate measure of the level of agreement between the model and system response. One example for automatic control is the sum of squared error function:

$$J = \sum_{i=1}^{N} e_i^2 \tag{1}$$

where $e_i$ is the error between experimental data and structure output for each of $N$ data points. Of course many other fitness functions could be used instead. Better model structures evolve as the GP algorithm minimizes the fitness function. Another possible fitness function could be:

$$J = 500^* \left( 8 - \log_{10} \left( \sum_{i=1}^{N} e_i^2 \right) \right) \tag{2}$$

This function was presented in [3] and it was successfully used to calculate the fitness of solution candidates for identifying a coupled water tank system using Genetic Programming.

## 3 An Automatic Pilot System

**Autopilots** mechanically guide a vehicle without assistance from a human being. An autopilot doesn't refer specifically to aircraft. Autopilots for boats and ships are called by the same name and serve the same purpose. They also

usually use similar processes. The aircraft autopilots are the most complex, and the most critical. The design of aircraft autopilots involves the flight dynamics that implies the orientation of air and space vehicles and how to control the critical flight parameters, typically named pitch, roll and yaw. The automatic pilot system for an aircraft normally involves control for various properties of the aircraft, depending of air properties like velocity, pressure, density, and temperature.

Very important factors in pursuance of the reference are shock waves that form in front of the nose of aircraft. The vehicle obtains thrust by the reaction to the ejection of fast moving exhaust from within the rocket engine.

In all rockets the exhaust is formed from propellant which is carried within the rocket prior to its release. Rocket thrust is due to the exhaust gases applying pressure on the inside surfaces of the rocket engine as they accelerate.

Rockets are also used for deceleration, to transfer to a lower-energy orbit, for example to enter into a circular orbit from outside, to de-orbit for landing, for the whole landing if there is no atmosphere.

Rockets must be used when there is no other substance (land, water, or air) that an aircraft may push against, such as in space. In these circumstances, it is necessary to carry all the propellant within the vehicle.

Rockets are particularly useful when very high speeds are required, such as orbital speed (mach 25 or so). The speeds that a rocket vehicle can reach can be calculated by the rocket equation; which gives the speed difference ("delta-v") in terms of the exhaust speed and ratio of initial mass to final mass ("mass ratio").

Common mass ratios for vehicles are 20/1 for dense propellants such as liquid oxygen and kerosene, 25/1 for dense monopropellants such as hydrogen peroxide, and 10/1 for liquid oxygen and liquid hydrogen. However, mass ratio is highly dependent on many factors such as the type of engine the vehicle uses and structural safety margins.

Sometimes, particularly in launch scenarios, the required velocity (delta-v) for a mission is unattainable because the propellant, structure, guidance and engines weigh so much as prevent the mass ratio from being high enough. This problem is frequently solved by staging – the rocket sheds excess weight (usually tankage and engines) to attain a higher effective mass ratio thus permitting a higher delta-v.

Typically the acceleration of a rocket increases with time due to applying the same thrust to a decreasing mass, with discontinuities when stages burn out, and starting at a lower acceleration with the new stage firing.

Tsiolkovsky's rocket equation, named after Konstantin Tsiolkovsky who first derived it, considers the principle of a rocket: a device that can apply acceleration to it by expelling part of its mass with high speed in the opposite direction, due to the conservation of momentum.

It says that for any maneuver or any journey involving a number of maneuvers:

$$\Delta v = v_e \ln \frac{m_0}{m_1} \tag{3}$$

or equivalently

$$m_1 = m_0 e^{-\Delta v/v_e}, m_0 = m_1 e^{\Delta v/v_e} \tag{4}$$

where $m_0$ is the initial total mass, and $m_1$ the final total mass and $v_e$ the velocity of the rocket exhaust with respect to the rocket.

$$1 - \frac{m_1}{m_0} = 1 - e^{-\Delta v/v_e} \tag{5}$$

is the mass fraction (the part of the initial total mass that is spent as reaction mass).

$\Delta v$ (delta v) is the integration over time of the magnitude of the acceleration produced by using the rocket engine (not the acceleration due to other sources such as gravity or drag). For the typical case of acceleration in the direction of the velocity, this is the increase of the speed. In the case of acceleration in opposite direction (deceleration) it is the decrease of the speed. Note that gravity or drag also changes velocity, but they are not part of the quantity delta-v. Hence delta-v is not simply the change in speed or velocity. However, thrust is often applied in short bursts, and during these short periods the other sources of acceleration may be negligible, and the delta-v of one burst may be simply approximated by the speed change. The total delta-v can simply be found by addition, even though between bursts the magnitude and direction of the velocity changes due to gravity, e.g. in an elliptic orbit.

Note that, as mentioned, at any time the *magnitude* of the acceleration contributes to the delta-v, hence always a non-negative value, regardless of whether the rocket is used for acceleration or deceleration. This again demonstrates that delta-v is not simply the change in speed or velocity: the latter may be zero if we first accelerate and than decelerate, but the delta-v accumulates.

The equation is obtained by integrating the conservation of momentum equation:

$$m^* dv = v^* dm \tag{6}$$

for a simple rocket that emits mass at a constant velocity ($dm$ is here the reaction mass; if it is the change of the rocket mass then there is a minus sign in the latter equation).

Although an extreme simplification, the rocket equation captures the essentials of rocket flight physics in a single short equation. It happens that delta-v is one of the most important quantities in orbital mechanics that quantifies how difficult it is to get from one trajectory to another.

Clearly, to achieve a large delta-v, either $m_0$ must be huge (growing exponentially as delta-v rises), or $m_1$ must be tiny, or $v$ must be very high, or some combination of all of these.

A Simple Genetic Algorithm has been used with an initial population of 20 chromosomes. Each chromosome is a sequence of 100 variables. We have

**Fig. 1.** The response with a normal controller

established the maximum number of generations as 1200 and a generation gap of 0.9.

In the Figs. 1 and 2 there are presented the results obtained with an PID algorithm and another with genetic algoritms.

# References

1. C. Jain Lakhmi, N.M. Martin (1998). Fusion of Neural Networks, Fuzzy Systems and Genetic Algorithms: Industrial Applications
2. J. H. Holland (1975). Adaptation in Natural and Artificial Systems, MIT Press, Cambridge
3. G. J. Gray., D. J. Murray, Li Y. Smith, K. C. Sharman, T. Weinbrenner (1998). Nonlinear model structure identification using genetic programming. Control Engineering Practice
4. C. Michael, G. McGraw. Automated Software Test Data Generation for Complex Programs
5. K. Chang, J. Cross, W. Carlisle and S. Liao (1996). A performance evaluation of heuristics based test case generation methods for software branch coverage. International Journal of Software Engineering and Knowledge Engineering
6. B. Korel (1990). Automated software test data generation. IEEE Transactions on Software Engineering

**Fig. 2.** The response with genetic algorithms controller

7. B. Korel (1996). Automated test data generation for programs with procedures. In Proceedings of the International Symposium on Software Testing and Analysis

8. W. Miller and D. L. Spooner (1976). Automatic generation of floating point test data. IEEE TSE

9. M. J. Gallagher and V. L. Narasimhan (1997). Adtest: A test data generation suite for ada software systems. IEEE TSE

10. Christopher C. Michael Gary E. McGraw Michael A. Schatz Curtis C. Waltony. Genetic algorithms for Dynamic Test Data Generation

11. J. R. Koza, M. A. Keane, F. H. Bennett W. Mydlowec (2000). Automatic Creation of Human-Competitive Programs and Controllers by Means of Genetic Programming

# Neuro-Fuzzy Based Nonlinear Models

C. Nitu and A. Dobrescu

"Politehnica" University of Bucharest, Str. Spl. Independenei 313, Bucharest Code
74206 Romania
cnitu@ecosys.pub.ro

**Abstract.** This paper presents a fuzzyfication method that consists in determining
the parameters of the membership functions by using control error intervals and
functions which cross these intervals. The proposed method can use constant or
variable control error intervals and linear or nonlinear functions which cross these
intervals and which determine the variable parameters of the neuro-fuzzy based
nonlinear model with the designed membership functions. First, nonlinear models
are presented, then a concrete system is analyzed in both cases: classic and adaptive.

**Keywords:** Fuzzy control systems, automatic control, human-centered de-
sign, human factors, human perception, human reliability, human supervisory
control.

## 1 Introduction

The behavior of the human operator in control systems of the technological
processes is defined by two models: the model of the human operator as a
transmitter of the command values and the mental model that describes the
decision process making to solve the problem defined by control law of the
controller. Regarding the deduction processes for mental data manipulation,
the human operator can modify the initial conceptual model depending on the
effects of his actions. After the mental command calculation and the mental
comparison of that result with the values of his own experience, the human
operator can decide the modification of mental model parameters, which can
lead to the improvement of the control process performances. This means that
the conceptual model is an adaptive mental model. The mathematical mental
model can also modify itself by the results of some logical functions that were
mentally run by the human operator [1].

In the design of the fuzzy controller the fuzzyfication of the controller in-
put and output variables consist in the choice of the membership functions
shape and parameters by using their graphs. By using this method of design

the parameters of the membership function are derived in accordance with the chosen graphs, and the equations are defined for the designed fuzzy controller. This fuzzyfication method has disadvantage that it doesn't use any information about the dependence of the parameters in function of the input and output controller variables or of other external signals. This dependence is necessary in the design of the adaptive fuzzy control system and in the modeling of the human operator behavior in complex distributed systems. This disadvantage is eliminated with the proposed method, where the determination of the membership function parameters is made with constant or variable control error intervals and a linear or nonlinear function which cross these intervals.

The main nonlinear models applications in control systems include the very complex processes as follows:

– The fuzzy model of the human operator behavior used in the designing of the digital control algorithms when the mathematical process models can not be derived and the human manual process control is simulated. By the human operator fuzzy model analysis, a lot of rules tables can be experimented by computer simulation, so finally it can be found an optimal procedure for manual command of the process. By this procedure can be trained the human operators for the control of the processes so that a good transitory process of the controlled value at the variation of the reference magnitude can be obtained, and so an optimization of the manually control process. This training method by using of the human operator fuzzy model is suited when the technological processes have nonlinearities that make the empiric training very difficult. The fuzzy methods to study the human operator behavior model allow obtaining a nonlinear mathematical model of the human control actions in normal control of the processes. The fuzzy model FM of the human operator behavior in the designing of the digital fuzzy controller contains the fuzzy mental model FMM that represents the control law in manual control and sensor-motor model SMM of the human operator

– The fuzzy model of the human operator behavior is included in the designing of the adaptive control system. In the designing of the complex applications when it is difficult to obtain precise process models the human operator knowledge is used to design the adaptive control algorithms, Fig. 2. The sensor-motor model SMM is defined by the transfer function:

$$H_{(s)} = \frac{Ke^{-\tau s}}{(1 + sT_1)(1 + sT_{2)}} \tag{1}$$

Where $\tau$ is time delay and $T_1$, $T_2$ the time constant, which depend of human operator characteristics [2].

– The neuro-fuzzy methods are used to improve the fuzzy mental models when the process control needs nonlinear control law, for example for variable structure controls system. The neuro-fuzzy techniques bring many ad-

vantages in mental modeling. It can be designed nonlinear membership functions because the neural networks are adequate for modeling methods for certain nonlinearities. The data acquisition can be made on the fuzzy based rules in according with the human operator experience in process data acquisition.

The calculation structure of a neuro-fuzzy mental model contains the certain following layers (Fig. 1.):



**Fig. 1.** The calculation structure of a neuro-fuzzy mental model

Layer 1 to transmits the input signals values from sensors and human operator directly to the next layer

Layer 2 to create the activation functions for a certain simple neural networks.

Layer 3 to obtain logical fuzzy operations

Layer 4 for the implementation of the inference rules

Layer 5 for the computation of the command values for actuators

The use of the neuro-fuzzy methods in the human operator behavior modeling and for the process modeling allows adopting new process control algorithms for processes that do not have precise mathematical models and claims long experience from the human operator. These new techniques allow process automation where the operator is obliged to intense his efforts and where conventional control algorithms did not allow obtaining the expected performances.

On the bases of the reality regarding the use of neural networks and fuzzy method in the human operator behavior modeling when the operator commands manually technological processes technological parameter, the following conclusions reveal:

The neural networks represent a new design technology for certain process control systems and allow improvement of the indicators that define human operator activities and the efficient functioning of the process.

The neural networks are in fact software products that are designed and used in control systems, but they need training based on data sets chosen and checked by special test data sets [3].

Due to the fact that many processes modify in time their dynamics, a static neuro-fuzzy mental model cannot be used because of the unknown parameters variation in time. The hybrid mental model can be used for control system, being composed by a neural network and a fuzzy model

## 2 The Determination of Membership Functions Paramaters

As it already has said the proposed method for determining the parameters of the membership functions uses constant or variable intervals for the definition of the equations and a linear or a non-linear function which will cross these intervals. The variable parameters of the linear or non-linear functions which will cross the variables of the intervals will represent the variable parameters of the adaptive fuzzy control system with the designed membership functions. For example the parameters of the membership functions with the triangle shapes are shown in Fig. 2.



**Fig. 2.** The membership function parameters design for linear function

For triangle membership functions and for a linear function $y = \propto \varepsilon$, (Fig. 2.) if the variables of the intervals are chosen as being $L_{ii} L_{mi}, L_{si}$, where then the equations of the triangle membership functions are:

$$\mu_{i,\text{tri}}(\alpha) = \mu_{i,\text{tri}}(\varepsilon, a_i, c_i) = \begin{cases} 0, & \varepsilon \leq \frac{L_{ii}}{\alpha} \\ \mu_{1i} = \frac{\varepsilon - a_i}{b_i - a_i} = \frac{\alpha\varepsilon - L_{ii}}{L_{mi} - L_{ii}}, & \frac{L_{ii}}{\alpha} \leq \varepsilon \leq \frac{L_{mi}}{\alpha} \\ \mu_{2i} = \frac{c_{i-n}}{b_i - a_i} = \frac{L_{si} - \alpha\varepsilon}{L_{si} - L_{mi}}, & \frac{L_{mi}}{\alpha} \leq \varepsilon \leq \frac{L_{si}}{\alpha} \\ 0, & \varepsilon \geq c_i \end{cases} \quad (2)$$

**Fig. 3.** The membership function parameters design for a nonlinear function

In the Fig. 2., the value of $\propto$ is 60 and the data obtained are:

$$a_i = 0.3333b_i = 0.6667c_i = 1$$

Similarly, this method can be used to design the membership functions with trapeze and bell shape, for linear or nonlinear function which will cross the chosen intervals of the controller input variables – Figs. 3, 4.

In the Fig. 3. the function has the value:

$$y = -c * e^{\varepsilon/t} \tag{3}$$

The values of c and T are 1 and the data obtained are:

$$a_i = 0.4122b_i = 0.6635c_i = 1$$

In the Fig. 2. the function has the value $y = c * e^{\varepsilon}$, the value of $c$ is 1 and the data obtained are:

$$a_i = 0.0837b_i = 0.5945c_i = 1$$

The described method for the determining the parameters of the controller membership functions allows introducing one or more variable parameters in the equations of the membership functions which are used in the adaptive fuzzy model to obtain the imposed system performances.

## 3 Adaptive Fuzzy Conrol Algorithm

The described adaptive fuzzy model is based on the modification of the membership functions in accordance with chosen criteria which depends of the

**Fig. 4.** The membership function parameters design for a nonlinear function

control system error value. Initially the values of the membership functions parameters are computed knowing the mathematical model of the controlled process. The proposed method to modify the membership functions parameters is presented for the case when chosen criteria is the value of the control system error, when two membership functions are used and when the parameter $\alpha$, can be modified between large limits. In Fig. 5. is presented how to design the membership functions for adaptive fuzzy control system for the described case.

The proposed adaptive fuzzy model was applied for a temperature control system when the plant transfer function was known. The control algorithm includes the following operations : the measurement of the controlled temperature $\theta$; the computing of the control system error $\varepsilon$; the calculation of the membership functions value $\mu_{jI}$ ( $\varepsilon;\theta$ ) in accordance with the chosen value of the parameter $\alpha_i$; and the computing of the command value by using the defuzzyfication method namely the center of area.

The algorithm to adapt the membership function parameters was chosen:

$$\alpha = \begin{cases} \alpha_1 & \text{for} \quad |\varepsilon| < \varepsilon_1 \\ \alpha_2 & \text{for} \quad |\varepsilon_1| \leq |\varepsilon| < |\varepsilon_2| \\ \alpha_3 & \text{for} \quad |\varepsilon| > |\varepsilon_2| \end{cases} \qquad (4)$$

The described algorithm for adaptive fuzzy models is based on the modification of parameter $\alpha$ of the membership functions in function of the control system error. The computer simulations for different transfer functions of the plant proved that presented algorithms offered good results for all studied cases. The proposed adaptive fuzzy model can be extended for the case when membership functions have two or more variable parameters. In this case the nonlinear function Y ($\varepsilon$ ) is described by one or more parameters which are adjusted by adaptive algorithm by using a criteria to improve the system performances.

**Fig. 5.** The adaptive membership functions

For the experienced temperature control system for the control of the temperature in a building was used. The input of the fuzzy control is the error $(\varepsilon = r - \theta)$ and the output is the command.

The response of the system at with classic fuzzy algorithm is presented in Fig. 6.(a) and the response of the system at with an adaptive fuzzy algorithm is presented in Fig. 5(b).



(a)                                                    (b)

**Fig. 6.** The response of the system

# 4 Conclusions

The proposed nonlinear fuzzy models applications are illustrated for two cases, in the designing of the adaptive control system to adapt the control law, a nonlinear model, in function of the system performances and in the modeling of the human operator behavior in complex distributed systems. The human operator modeling requires the using of the neuro-fuzzy method to obtain mental and a sensor-motor model. The mathematical model of the human operator behavior in control systems is nonlinear because the parameters of its transfer functions are modifying in time and can be considered that they have constant value only for a small time interval. It has a complex serial-parallel structure with elements of commutation in function of the mental calculation results and of the process conditions and because in many applications the thinking process can be modified in time.

# References

1. B. Kirwan, (1994) *A Guide to Practical Human Reliability Assessment*, Taylor & Francis.
2. D.L. Day and D.K. Kovacs, *Computers, Communication and Mental Models,* Taylor & Francis.
3. C. Nitu, (2005) *Human Operator in Digital control Systems*, (in Romanian), Ed. MATRIX, Bucuresti.

# Reconfigurable Pattern Generators Using Nonlinear Electronic Circuits

J. Neff[1], V. In[1], A. Kho[1], A. Bulsara[1], B. Meadows[1], A. Palacios[2],
S. Hampton[3], L. Nguyen[3], D. Chi[3], and N. Koussa[3]

[1] Space and Naval Warfare Systems Center, Code 2373, 49590 Lassing Road A341, San Diego, CA 92152-5001, USA
[2] Mathematics Dept., San Diego State University, San Diego, CA 92182, USA
[3] University of California, San Diego, 9500 Gilman Dr., La Jolla, CA 92093 (858) 534-2230

## 1 Introduction

In these experiments we explore the practical implementation of pattern generating electronic circuits. The project takes advantage of fundamental mathematical arguments, based on symmetry, in order to define the circuits and coupling topologies that are used. By constructing networks of low-order dynamical neuron models, and by considering symmetries in the way the neurons are coupled, including time-shift invariant symmetries, the patterns generated by the system are reduced to a predictable few. Using the resulting patterns, the relative phases between the synchronized neurons within the network are used to define gait patterns that are similar to those found in living quadrupeds. Electronic circuits based on these observations are engineered to generate the needed analog signals for driving locomotion in an N-legged robot.

The experiments demonstrate an example central-pattern-generator (CPG) and its use in guiding a walking robot (in this case a quadruped). The project begins by adopting the Fitzhugh-Nagumo neuron model as the fundamental CPG element. In the CPG, the neurons are created using discreet electronic components, but the full design is modular and extensible, consisting of a single motherboard (perforated circuit board using discreet components) and a number of daughter boards that are nearly identical. The motherboard is used to translate the CPG signals into signals required by servos that drive the limbs of the robot while the daughter-boards include two neurons each with parameters and coupling strengths that are programmable via a serial bus. The mother and daughter-boards together include all the coupling circuitry needed to make generic and extensible central-pattern generating networks.

## 2 Background

Central Pattern Generators (CPGs) are networks of neurons that generate self-sustaining patterns of behavior for controlling various physiological functions such as locomotion, mastication and respiration [3,4]. CPGs are typically located in the Central Nervous System (CNS). To initiate a particular function, the CNS first translates the CPG rhythm into a coordinated pattern of activity, and it then sends that coordinated pattern to motorneurons innervating muscle fibers. In many cases, the coordinated pattern is just a faithful image of the CPG rhythm. In invertebrates and primitive vertebrates [5], in particular, it has been established that the fictive locomotion produced by the CPG is very similar to the actual rhythmic motor output. Golubitsky et al. [1, 2, 6, 7] used this fact as a modeling assumption for constructing realistic, biologically-inspired, mathematical models that can reproduce the spatio-temporal patterns of animal locomotion, called *animal gaits*. Central to their works is the postulate that the natural symmetries that are observed throughout the animal kingdom must be present, to some degree, in the actual CPG architecture. This postulate leads them to characterize the phase relations in the gaits of legged animals through periodic solutions, arising via symmetry-breaking Hopf bifurcations [8] of a coupled system of differential equations with symmetry. In the special case of quadrupeds, a critical contribution of the work by Golubitsky et al. is a theorem that states that a network of eight cells with $\mathbf{Z}_4 \times \mathbf{Z}_2$ symmetry is the smallest network that can generate *all* primary gaits: walk, jump, trot, bound, pace, and pronk. Furthermore, their network configuration avoids undesired conjugacies between trot and pace, which would imply inconsistencies with actual experimental observations.

In this work, we present the first circuit realization of an animal (quadruped) robot controlled by a CPG network whose design is based on the work of Golubitsky et al. This Paper is organized as follows. We first present a brief review of the network model equations for readers not familiar with the literature. Then we focus on the circuit realization of an individual neuron as well as on the coupling topology that leads to the actual network implementation. We demonstrate, through hardware simulations of the CPG network, that our electronic CPG and our animal robot can indeed generate all primary gaits–walk, jump, trot, bound, pace, and pronk–just as they were originally described by theory. (The paper so far lacks a time series or example of a complete pattern).

## 3 CPG Network

Following Buono and related works, the minimal CPG network consists of eight neurons arranged into two bilateral arrays. Within each array, the neurons are interconnected in a directed ring fashion. Across the arrays, however, the neurons are bidirectionally coupled, pairwise. Figure 1 depicts the actual

**Fig. 1.** A minimum of eight neurons is required to generate all the patterns for locomotion. In the network, the unidirectional coupling within the *top* and *bottom* arrays result in the $\mathbf{Z}_4$ symmetry. The bidirectional coupling between the arrays results in the $\mathbf{Z}_2$ symmetry

network. This interconnection scheme leads to $\mathbf{Z}_4 \times \mathbf{Z}_2$ symmetry as follows. $\mathbf{Z}_4$ symmetry implies that the network remains unchanged under, simultaneous, cyclic permutations of the neurons on each array. $\mathbf{Z}_2$ symmetry, on the other hand, implies invariance under bilateral exchange between neurons of both arrays. In Fig. 1 the numbered nodes represent the neurons in the CPG and the arrows represent the direction of the coupling. The responses of the neurons one through four are used to drive the left-rear, right-rear, left-front and right-front legs of the robot respectively.

### 3.1 Neuron Circuit Realization

The internal dynamics of each individual neuron in the network is governed by the two-dimensional Fitzhugh-Nagumo equations:

$$\dot{x} = c\left(x + y - \frac{1}{3}x^3\right) \equiv f_1(x, y)$$

$$\dot{y} = -\frac{1}{c}(x - a + by) \quad \equiv f_2(x, y) \,,$$

(1)

where $a$, $b$, and $c$ are parameters. The circuit realization of Fitzhugh-Nagumo model (1) was carried out as a traditional analog computer, where state variables are represented as voltages on nodes of an electronic circuit. A variety of simple operational-amplifier based circuits are used to construct the analog computer. Typically a completed analog computer will include sum and difference amplifier circuits as well as voltage integrating circuits [9]. Special function circuits, such as those that provide a nonlinear input-output response, are constructed using piece-wise linear functions.

Construction of an analog computer typically proceeds in the following way. The system that is being modeled is written in standard form, that is, as a set of coupled 1st-order differential equations as in (1). As an arbitrary

starting point, we use the voltage on the first node to represent the variable $\dot{x}$ and we use this node as the input to the first integrator to obtain the state-variable $x$ as a function of time. The output of the integrator, which is typically scaled by a gain factor, can be used as an input to other op-amp based circuits for computing functions of $x$. The outputs of these functions are then summed using a summing amplifier, the output of which is the initial starting point $\dot{x}$. The equivalency between the left-had side and the right hand side of the equation representing the system is obtained via feedback in the analog computer. By starting at any point on the circuit loop described above, and by carefully accounting for all the gains and sign changes through the loop of the circuit, an electronic circuit for modeling the time evolution of the state variable $x$ can be obtained. Circuits for the remaining sate variables are constructed similarly. Coupling terms between the variables are easily implemented using the summing amplifiers. For brevity, only a simplified and



**Fig. 2.** Circuit for $\dot{x}$, plus piec-wise nonlinear function for $x^3$

partial schematic for the Fitzhugh-Nagumo model is given here. Figure 2 illustrates a circuit loop used to model the state variable $x$ given in (1). The inset is the symbol for representing a single neuron in the CPG network as it appears in Fig. 1. Starting at $\dot{x}$ and moving counter-clockwise around the loop, the circuit involves: an inverting integrator, a non-inverting buffer, and a summing-inverting amplifier, which feeds the input to the integrator. A piece-wise approximation to $x^3$ is also computed and provided as an input to the summer, as well as the state variable $y$. The variable $y$ is computed using a similarly constructed circuit (not shown). The system parameters $a$, $b$, and $c$ and the time scale are all determined by the values of the resistors and capacitor shown. Coupling circuits between the neurons allow us to set both the strength and the signs of the coupling terms $X_{Dir}$ and $X_{Bi}$ (not shown).

A drawback of this design approach is that the system parameters are set by physical properties, such as device capacitances and resistances, which are typically fixed when those devices are fabricated. To circumvent this problem we use solid-state programmable resistors. This allows us to program the parameters of the neuron model using a serial buss. We use the same approach

for programming the coupling strength between neurons within the CPG network. This ability to program the system parameters is necessary for us to demonstrate all the locomotion patterns using a single network.

## 3.2 Network Circuit Realization

Using (1), we write the CPG model proposed by Golubitsky et al. in the following form

$$
\begin{aligned}
\dot{x}_i &= f_1(x_i, y_i) + X_{Dir}(x_{i-2} - x_i) + X_{Bi}(x_{i+\varepsilon} - x_i) \\
\dot{y}_i &= f_2(x_i, y_i) + Y_{Dir}(y_{i-2} - y_i) + Y_{Bi}(y_{i+\varepsilon} - y_i) \,,
\end{aligned}
\tag{2}
$$

where $i = 1, \ldots, 8$ mod 8, $\varepsilon = +1$ if $i$ is odd, and $\varepsilon = -1$ when $i$ is even, and $X_{Dir}$, $X_{Bi}$, $Y_{Dir}$, and $Y_{Bi}$, are the coupling strengths that appear in Fig. 1. The network in Fig. 1 is created using the neuron cirtcuit described above as the unit cell. The system of coupled neurons that make up the CPG are divided into an array of stackable circuit boards, each containing two bilaterally coupled Fitzhugh-Nagumo circuits (such as cells one and two in Fig. 1), and a mother-board that interfaces the CPG to the servos (legs) of the robot. For example, the network shown in Fig. 1 is implemented using 4 circuit boards. Each neuron board attaches itself to a common serial buss, making all the system parameters associated with its two neurons available to be set or read via the buss. Similar to the individual neuron parameters, the coupling strengths in (2) are also set via the serial buss. The system is can be extend to robots with more than four legs, such as hexapods, by adding more neuron boards to the stack.

# 4 Patterns and Locomotion

To select a particular pattern, the appropriate values for the coupling resistors are set using the serial bus. These values are then stored in non-volatile memory on the programmable resistor chips so that, once programmed, the CPG will generate the walking pattern without any digital support. In this regard, the CPG system we've developed is decentralized, meaning that the computation needed for locomotion is performed using low-level circuits that are decentralized from higher level circuits. In this system the high-level functions (which this paper does not address) can be performed using a micro-controller. Currently a micro-controller is only used to program the system parameters. A simple look-up table is used to store the appropriate resistor values, based on the desired mode of locomotion.

Identifying the appropriate system parameters that will result in a desired pattern is a significant challenge to creating useful CPGs. Symmetry arguments which take advantage of the inherent symmetries of a system can quickly help to to identify possible solutions and can give some insight into

the strengths and signs of the coupling terms. For this system, the $\mathbf{Z}_4 \times \mathbf{Z}_2$ symmetry allows the single network to support all the known quadruped patterns, pronk, pace, bound, trot, jump and walk, so long as certain certain sign conventions for the bidirectional and coupling terms are followed. These conventions are given in Table 1. The table shows what set of coupling signs are associated with a particular pattern. These methods do not necessarily provide precise values for the coupling parameters nor do they dictate the stability of the solution. For this work we use values that are similar to those suggested in previous work by Golubitsky et al. In practice we find that, as long as the coupling signs are correct, there is a large parameter space that can produce any particular stable pattern.

**Table 1.** Signs of coupling strengths for each of the animal gaits generated by the electronic CPG

| Gait | $X_{Dir}$ | $X_{Bi}$ | $Y_{Dir}$ | $Y_{Bi}$ |
|------|-----------|----------|-----------|----------|
| Pronk | + | + | + | + |
| Pace | + | − | + | − |
| Bound | − | + | − | + |
| Trot | − | − | − | + |
| Jump | − | + | + | + |
| Walk | − | − | + | + |

## 4.1 Leg Motion

Given the functioning analog CPG circuit described above, additional circuitry is needed to translate the time-vary voltages from the CPG into the pulse-width modulated signals needed for the servo motors. These simple translation circuits are implemented on the the mother board along with the micro-controller described above. In this translation there is a convientient one-to-one comparison between the phases of the syncronized $x$ and $y$ values of a single neuron and the elevation and the forward (or backward) position of a leg. Figure 3a is an example single cycle of the $x$ and $y$ variables for a single neuron as a function of time. For a stable pattern, the $x$ and $y$ variables are synchronized to each-other, oscillating at a common frequency, with a constant phase difference between the two. Since one variable always leads or lags the other, a comparison can be made between the dynamics of the neuron and the motion of a leg. Specifically, each quadruped's leg has two axis of motion, one axis for moving the leg forwards and backwards, and one axis for raising and lowering the leg. We associate the $x$ variable with forward backward motion and the $y$ variable with raising and lowering the leg. Figure 3a shows how the proper sequence of movements, associated with the phase of the neuron, results in coordinated motion of the leg Fig. 3b.

**Fig. 3.** The phases of a neuron oscillator (**a**) and the associated movement of a leg (**b**)

## 5 Conclusions and Future Work

We've developed the first circuit realization of an animal (quadruped) robot controlled by a CPG network whose design is based on the work of Golubitsky et al. The system is extensible and so can support a variety of N-legged robots. Since the system is programmable and uses non-volatile memory, a wide variety of behaviors can be explored and is a good example of a decentralized pattern generator. For the quadruped, the system produces several gait patterns that are functional, in particular the walk and trot patterns. Other patterns, such as jump, are not well supported by the servo driven legs of our robot. This is because, although hobby servos can provide the relatively high torque needed to lift and move the robot, they tend to be too slow for such movements. Robots with more than four legs can be accommodated simply by adding more daughter boards to the CPG and by programming the system appropriately. Using this system, a single network is used to generate all the needed patterns for locomotion for a particular robot.

Although the resulting CPG is physically large, a similar system is being developed as a single application-specific integrated circuit and will similarly demonstrate the concepts behind creating practical GPGs using nonlinear electronic circuits. The design of that system will be minimalist, and will benefit from our observation that CPG systems that possess inherent symmetries are robust to the neuron cells being implemented and to the specific form of the coupling between the cells.

We would like to acknowledge the Naval Research Enterprise Program for sponsoring this work.

## References

1. M. Golubitsky, I. Stewart, P.-L. Buono, and J.J. Collins. *Physica D* **115**, 56 (1998).

2. M. Golubitsky, I. Stewart, P.-L. Buono, and J.J. Collins. *Nature* **401**, 693 (1999).
3. A. Cohen, S. Rossignol and S. Grillner (Eds). *Neural control of rhythmic movements in vertebrates.* Jon Wiley & Sons, New York (1988).
4. D.W. Morton and H.J. Chiel. *J. Comp. Physiol.* **A 173** 519 (1993).
5. S. Grillner, J. Buchanan, P. Wallen, and L. Brodin. In: *Neural Control of Rhythmic Movements in Vertebrates. (A.H. Cohen, S. Rossignol, and S. Grillner eds.)* New-York, Wiley, 129 (1988).
6. P.-L. Buono and M. Golubitsky. *J. Math. Biol.* **42** 291 (2001).
7. P.-L. Buono. *J. Math. Biol.* **42** 327 (2001).
8. M. Golubitsky, I.N. Stewart, and D.G. Schaeffer. *Singularities and Groups in Bifurcation Theory: Vol. II.* Appl. Math. Sci. **69**, Springer-Verlag, New York, 1988.
9. Horrowitz and Hill, "The Art of Electronics."

# Configuring A Non-Linear Process Control System Using Virtual Instrumentation

PhD student A. Enescu, PhD student G. Costache

"Politehnica" University of Bucharest, 313 Splaiul Independentei, Bucharest 060042, Romania
gabrielc@ecosys.pub.ro
alexe@ecosys.pub.ro

**Abstract.** It was in 1896 when S. Arrhenius first noticed the potential effect of human activities on the carbon cycle and the implications for climate change. He put forward the theory that $CO_2$ in the atmosphere was an important greenhouse gas and that it was a by-product of burning fossil fuels. In 1958, Charles Keeling began the observations at Mauna Loa Observatory, 3650 m up a mountain in Hawaii, regarded as far enough away from any carbon dioxide source to be a reliable measuring point. Measurements of $CO_2$ in the atmosphere have been continuous for almost 50 years. In recent decades, $CO_2$ increased on average by 1.4 parts per million (ppm) a year because of the amount of fossil fuels burnt.

## 1 Introduction and History

Carbon is unquestionably one of the most important elements on Earth. It is the principal building block for the organic compounds that make up life. Carbon's electron structure gives it a plus 4 charge, which means that it can readily form bonds with itself, leading to a great diversity in the chemical compounds that can be formed around carbon; hence the diversity and complexity of life. Carbon occurs in many other forms and places on Earth; it is a major constituent of limestone, occurring as calcium carbonate; it is dissolved in ocean water and fresh water; and it is present in the atmosphere as carbon dioxide, the second most important greenhouse gas.

The flow of carbon throughout the biosphere, atmosphere, hydrosphere, and geosphere is one of the most complex, interesting, and important of the global cycles. More than any other global cycle, the carbon cycle challenges us to draw together information from biology, chemistry, oceanography, and geology in order to understand how it works and what causes it to change.

## 2 Kyoto Protocol

There is growing concern about the impact that increased emissions of certain gases, known as "greenhouse gases", are having on the global climate. Because of this, the Kyoto Protocol has been established to limit emissions of greenhouse gases.

Under the Kyoto Protocol, industrialised countries and those in transition to a market economy (the so-called "Annex I countries") have agreed to limit or reduce their emissions of these greenhouse gases. The "Annex I Countries" are those that have taken on emission reduction or limitation targets under the Kyoto Protocol.

The Protocol sets quantified emission limitations and reduction obligations with respect to a basket of six gases. Of these, carbon dioxide ($CO_2$), which derives from the burning of fossil fuels such as coal, oil and gas, is the most important. Methane ($CH_4$) and nitrous oxide ($N_2O$) emissions are also substantial contributors to the problem. The targets define the amount of greenhouse gases that the countries are allowed to emit in the "commitment period" of 2008 to 2012, relative to the amount emitted in 1990. These targets represent either a cut in emissions or a lower rate of increase in emissions.

## 3 Mauna Loa

In the late 1950's, Roger Revelle, an American oceanographer, and a colleague, Charles Keeling, began monitoring atmospheric $CO_2$ at an observatory on Mauna Loa, on the big island of Hawaii. The record from Mauna Loa, shown in Fig. 1 below, is a dramatic sign of global change that captured the attention of the whole world because it shows that this "experiment" we are conducting is apparently having a significant effect on the global carbon cycle. The climatologically consequences of this change are potentially of great importance to the future of the global population.

The Mauna Loa atmospheric $CO_2$ measurements constitute the longest continuous record of atmospheric $CO_2$ concentrations available in the world. The Mauna Loa site is considered one of the most favourable locations for measuring undisturbed air because possible local influences of vegetation or human activities on atmospheric $CO_2$ concentrations are minimal and any influences from volcanic vents may be excluded from the records. The methods and equipment used to obtain these measurements have remained essentially unchanged during the 47-year monitoring program.

Because of the favourable site location, continuous monitoring, and careful selection and scrutiny of the data, the Mauna Loa record is considered to be a precise record and a reliable indicator of the regional trend in the concentrations of atmospheric $CO_2$ in the middle layers of the troposphere. The Mauna Loa record shows a 19.4% increase in the mean annual concentration, from 315.98 parts per million by volume (ppmv) of dry air in 1959 to 377.38 ppmv

**Fig. 1.** Mauna Loa historical records of $CO_2$

in 2004. The 1997–1998 increase in the annual growth rate of 2.87 ppmv represents the largest single yearly jump since the Mauna Loa record began in 1958. This represents an average annual increase of 1.4 ppmv per year.

## 4 Carbon Cycle Modelling

We modified the C.free model of carbon cycle (Fiddaman, 1997), which is an eddy diffusion model with stocks of carbon in the atmosphere, biosphere, mixed ocean layer, and 10 deep ocean layers.

In the C.free model, all emissions initially accumulate in the atmosphere. As the atmospheric concentration of $CO_2$ rises, the uptake of $CO_2$ by the ocean and biosphere increases, and carbon is gradually stored. The atmospheric flux to the biosphere consists of net primary production (NPP). Net primary production grows logarithmically as the atmospheric concentration of $CO_2$ increases according to:

$$NPP = NPP_0 \left( 1 + \beta_b \ln \left( \frac{C_a}{C_{a,0}} \right) \right) \tag{1}$$

$NPP$ = net primary production
$NPP_0$ = reference net primary production
$\beta_b$ = biostimulation coefficient
$C_a$ = $CO_2$ in atmosphere
$C_{a,0}$ = reference $CO_2$ in atmosphere

Because the relationship is logarithmic, the uptake of $CO_2$ by the biosphere is less than proportional to the increase in atmospheric $CO_2$ concentration. Effects of the current biomass stock, temperature, and human disturbance are neglected.

It is worth noting that this formulation is not robust to large deviations in the atmospheric concentration of $CO_2$. As the atmospheric concentration of $CO_2$ approaches zero, net primary production approaches minus infinity, which is not possible given the finite positive stock of biomass. As the concentration of $CO_2$ becomes very high, net primary production can grow arbitrarily large, which is also not possible in reality. Neither of these constraints is a problem for reasonable model trajectories, though.

To simplify the model, detailed biospheres are aggregated into stocks of biomass (leaves, branches, stems, roots) and humus (litter, humus).

$$C_b(t) = \int NPP(t) - \frac{C_b(t)}{\tau_b} dt \tag{2}$$

$C_b$ = carbon in biomass
$\tau_b$ = biomass residence time

$$C_h(t) = \int \frac{\Phi C_b(t)}{\tau_b} - \frac{C_h(t)}{\tau_h} dt \tag{3}$$

$C_h$ = carbon in humus
$\tau_h$ = humus residence time
$\Phi$ = humification fraction

The interaction between the atmosphere and mixed ocean layer involves a shift in chemical equilibrium. $CO_2$ in the ocean reacts to produce $HCO_3^-$ and $CO_3^{-2}$. In equilibrium,

$$C_m = C_{m,0} \left( \frac{C_a}{C_{a,0}} \right)^{\left( \frac{1}{\xi} \right)} \tag{4}$$

$C_m$ = $CO_2$ in mixed ocean layer
$C_{m,0}$ = reference $CO_2$ in mixed ocean layer
$C_a$ = $CO_2$ in atmosphere
$C_{a,0}$ = reference $CO_2$ in atmosphere
$\zeta$ = buffer factor

The atmosphere and mixed ocean adjust to this equilibrium with a time constant of 9.5 years.

The buffer or Revelle factor, $z$, is typically about 10. As a result, the partial pressure of $CO_2$ in the ocean rises about 10 times faster than the total concentration of carbon.

This means that the ocean, while it initially contains about 60 times as much carbon as the preindustrial atmosphere, behaves as if it were only 6 times as large.

The buffer factor itself rises with the atmospheric concentration of $CO_2$ and temperature. This means that the ocean's capacity to absorb $CO_2$ diminishes as the atmospheric concentration rises. The temperature effect (which is omitted in this model) is one of several possible feedback mechanisms between the climate and carbon cycle.

$$\xi = \xi_0 + \delta_b \ln\left(\frac{C_a}{C_{a,0}}\right) \tag{5}$$

$\zeta = $ *buffer factor*
$\zeta_0 = $ *reference buffer factor*
$\delta_b = $ *buffer $CO_2$ coefficient*
$C_a = $ *$CO_2$ in atmosphere*
$C_{a,0} = $ *reference $CO_2$ in atmosphere*

The deep ocean is represented by a simple eddy-diffusion structure similar to that in the Oeschger model, but with fewer layers. Effects of ocean circulation and carbon precipitation, present in more complex models, are neglected. Within the ocean, transport of carbon among ocean layers operates linearly. The flux of carbon between two layers of identical thickness is expressed by:

$$F_{m,n} = \frac{(C_m - C_n)^e}{d^2} \tag{6}$$

$F_{m,n} = $ *carbon flux from layer m to layer n*
$C_k = $ *carbon in layer k*
$e = $ *eddy diffusion coefficient*
$d = $ *depth of layers*

The effective time constant for this interaction, $e/d2$, varies with d, the thickness of the ocean layers. This model employs a 75 meter mixed layer, five 200 meter middle layers, and five 560 meter deep ocean layers with the time constants of 1.4 years, 10.0 years and, respectively, 78.4 years. Models with fewer ocean layers underestimate the short term participation of the ocean in carbon uptake and must increase uptake by other means to compensate.

## 5 Model Results

Rockets The historical carbon emissions (Figs. 2–4) summarize two main components: from fossil fuels burnt and from land use, since 1850. This data provides estimates of global net carbon fluxes, on a year-by-year basis from 1850 through 2000, resulting from fossil fuels burnt and changes in land use (such as harvesting of forest products and clearing for agriculture), taking into account not only the initial removal and oxidation of the carbon in the vegetation, but also subsequent regrowth and changes in soil carbon. We used

**Fig. 2.** Historical carbon emissions from 1850



**Fig. 3.** Different scenarios for carbon emissions

a five years mean of the total historical carbon emissions as an input lookup for the C.free model.

After 2004, we proposed different scenarios for carbon emissions, according with some IPCC emissions scenarios, with Kyoto Protocol and even the utopist scenario of complete emissions cut-off after 2005.

After running the model with these different scenarios, the results were in accordance with IPCC models results. We can clearly see that even if we will stabilize the emissions, the amount of $CO_2$ in atmosphere will continue to rise.

## 6 Conclusions

At the 1997 Kyoto conference, 38 industrialized nations agreed to reduce emissions to about 95% of 1990 rates by 2012. While the agreement is better than

**Fig. 4.** Projected $CO_2$ levels in atmosphere until 2300

business as usual, rapidly developing nations like China are not signatories, and their emissions continue to grow. The policy debate has become a fight over whether to stabilize the emission rate, not the stocks of greenhouse gases that drive the climate. Even if Kyoto were fully implemented, emissions would continue to exceed removal and GHG concentrations would continue to rise. The fight over implementation of the Kyoto Protocol, therefore, has become a debate about how much more GHG concentrations in the atmosphere will rise, and how much faster the global climate will warm. Halting warming, much less reversing it, is not even on the table.

# References

1. Thomas S. Fiddaman (1990) – Feedback Complexity in Integrated Climate-Economy Models, A.B., Engineering Sciences, Dartmouth College
2. *** (2005) – DTI: Department of Trade and Industry [http://www.dti.gov.uk]
3. *** (2005) – Trends: Atmospheric carbon dioxide [http://cdiac.esd.ornl.gov/tre-nds/trends.htm]
4. John D. Sterman, Linda Booth Sweeney (2002) – Cloudy Skies: Assessing Public Understanding of Global Warming

# Understanding Complex Systems

**Edited by J.A. Scott Kelso**