

Neuromorphic Controller for Low Power Systems

From Devices to Circuits

by

Saurabh Sinha

A Dissertation Presented in Partial Fulfillment  
of the Requirements for the Degree  
Doctor of Philosophy

Approved November 2011 by the  
Graduate Supervisory Committee:

Yu Cao, Chair

Hongbin Yu

Jennifer Blain Christen

Bertan Bakkaloglu

ARIZONA STATE UNIVERSITY

December 2011

## ABSTRACT

A workload-aware low-power neuromorphic controller for dynamic power and thermal management in VLSI systems is presented. The neuromorphic controller predicts future workload and temperature values based on the past values and CPU performance counters and preemptively regulates supply voltage and frequency. System-level measurements from state-of-the-art commercial microprocessors are used to get workload, temperature and CPU performance counter values. The controller is designed and simulated using circuit-design and synthesis tools. At device-level, on-chip planar inductors suffer from low inductance occupying large chip area. On-chip inductors with integrated magnetic materials are designed, simulated and fabricated to explore performance-efficiency trade offs and explore potential applications such as resonant clocking and on-chip voltage regulation. A system level study is conducted to evaluate the effect of on-chip voltage regulator employing magnetic inductors as the output filter. It is concluded that neuromorphic power controller is beneficial for fine-grained per-core power management in conjunction with on-chip voltage regulators utilizing scaled magnetic inductors.



## DEDICATION

To my parents Pijus Kumar Sinha and Ruma Sinha and my fiance Ripa. Their contribution cannot be mentioned in words.

## ACKNOWLEDGEMENTS

I would like to express my sincere gratitude respect to my advisor and mentor, Dr. Yu Cao, for his continued support and invaluable guidance during the course of this degree, without which this work would not have been possible. His trust in my capabilities and patience in helping me understand new concepts have been the motivating force for this work. His enthusiastic support in academic, research and personal endeavors has been a constant source of inspiration.

I am grateful to Dr. Hongbin Yu and Dr. Bertan Bakkaloglu for their insightful discussions during research meetings and for being on my PhD committee. I would also like to thank Dr. Jennifer Blain Christen for agreeing to be on my committee. I am thankful to fellow graduate students, Asha Balijepalli, Kautilya Patel, Nishith Desai, Anshu Pratik, Parasaran Ramanujam, Dhaval Shah, Sarth Shah, Sudarshan Rangunathan, Amar Velamala and Nirav Desai for invigorating discussions and feedback about my research. My thanks to Dr. Schroder and Dr. Clark for teaching Semiconductor Device Theory 1 and Advanced VLSI Design courses respectively, thus providing a very strong foundation of this field. A special mention for the EEE 525 VLSI Design class of Spring 2011 where I learned some very basic yet essential concepts by 'teaching'.

My internship at ARM, Austin has helped me grow as a researcher. I would like to thank Greg Yeric, Brian Cline, Vikas Chandra and the R&D group for showing me how exciting industrial research can be.

I am indebted to my family for their unconditional love, support and patience. My special thanks to the Apartment 31 gang, the 'Cure for Fridays' late night music group, Chirayu, Parth and Ripa, for making graduate life so much fun.

# TABLE OF CONTENTS

	Page
LIST OF TABLES .....	vii
LIST OF FIGURES .....	viii
CHAPTER	
1 INTRODUCTION .....	1
1.1 Dynamic Power Management .....	3
1.2 Neuromorphic Power Management .....	4
1.3 Voltage Regulators.....	6
1.2 On-chip Spiral Inductor with Integrated Magnetic Materials.....	7
1.1 The Integrated View .....	8
2 NEUROMORPHIC POWER CONTROLLER.....	10
2.1 Background and Previous Work.....	12
Spiking Neuron Architecture.....	12
Dynamic Power Management.....	13
Dynamic Thermal Management .....	14
2.2 CMOS Implementation of Neuromorphic DVS Controller .....	16
Analog Design .....	16
Neurons .....	16
Modulation.....	17
Computation .....	18
Training .....	19
Digital Design.....	20
Error and Training.....	21
Temperature Prediction.....	21

CHAPTER	Page
Implementation.....	22
2.3 Workload Characterization and Power Modeling.....	23
Power Modeling.....	25
2.4 Results and Discussion.....	27
Prediction Accuracy.....	29
Comparison with Hardware-based Schemes.....	31
Power Savings.....	31
2.5 Fine-grained Power Management.....	34
2.6 Conclusion and Future Work.....	35
3 DYNAMIC POWER MANAGEMENT AND REQUIREMENTS FOR ON- CHIP VOLTAGE REGULATORS.....	38
3.1 Switch-mode Converter Basics.....	40
Efficiency Model.....	43
3.2 Review of On-Chip Voltage Regulators.....	45
4 DESIGN OF ON-CHIP SPIRAL INDUCTORS WITH INTEGRATED MAGNETIC MATERIALS.....	48
4.1 Simulation and Analytical Study.....	49
Analytical Model of Stripline Inductors with Magnetic Materials.....	49
Comparison between Analytical Model and 3D EM Solver.....	50
Optimization of Inductor Structure.....	53
4.2 Device Fabrication and Measurement.....	55
Inductors with Magnetic Dots for High Frequency Applications.....	55
Fabrication.....	55
Results and Discussions.....	56

CHAPTER	Page
Equivalent Circuit Model and Parasitic Extraction .....	57
Inductors with Magnetic Rings .....	60
Fabrication .....	61
Results and Discussions.....	62
Permeability Dependent Inductance and Q factor .....	66
Effect of Magnetic Material Patterining .....	66
Inductors using CoZrTa.....	67
Conclusion .....	69
4.3 Potential Application: Resonant Clock Distribution .....	70
4.1 Potential Application: Magnetic Inductors with Voltage Regulators .....	72
Area Savings.....	76
4.5 Conclusion.....	76
5 CONCLUSION .....	77
5.1 Future Work .....	78
REFERENCES.....	79

## LIST OF TABLES

Table		Page
2.1	Processor Specifications .....	23
2.2	Processor Performance Events .....	23
2.3	Summary of Results .....	30
2.4	Calculated Power Savings .....	33
3.1	Related Work .....	46
4.1	Extracted Values from Measurements .....	60
4.2	Simulation Parameters .....	73

## LIST OF FIGURES

Figure	Page
1.1 Power Dissipation of desktop and mobile processors with time .....	2
1.2 Traditional constant supply and frequency vs. DVFS .....	3
1.3 Power Saving due to DVFS .....	4
1.4 Block diagram showing software based power management .....	5
1.5 Simplified schematic of a switch-mode DC-DC converter .....	6
1.6 On-Chip spiral inductor with integrated magnetic material .....	8
1.7 Scope of this work .....	9
2.1 Plot of machine complexity vs. task complexity .....	11
2.2 Workload and DVFS profile and corresponding temperature.....	11
2.3 Schematic of biological neuron, membrane potential and circuit equivalent .....	12
2.4 Block diagram of conventional vs. neuromorphic DVFS .....	13
2.5 Block diagram of spiking neuromorphic controller .....	15
2.6 Schematic of integrate and fire neuron and membrane potential .....	17
2.7 Current to spike rate modulation range .....	18
2.8 Schematic diagram of the plastic synapse .....	18
2.9 Schematic diagram of training circuit and updates in synaptic weight .....	19
2.10 Block level schematic of digital neuromorphic circuit.....	21
2.11 Timing digram showing spike generation, prediction and adaptation.....	22
2.12 Normalized power dissipation vs. operating frequency .....	27
2.13 Workload Prediction and associated P-states.....	28
2.14 Workoad and Temperature Prediction .....	30
2.15 Effect if bias on performance hit and over-estimation .....	32
2.16 Power savings and associated delay for different applications .....	33

Figure	Page
2.17 Effect of changing the ratio of leakage to dynamic power .....	33
2.18 Plot showing percentage of CPU utilization/memory access time .....	34
2.19 Actual vs. predicted workload from neuromorphic controller for <i>ocean</i> .....	36
2.20 Prediction Error vs. update times .....	37
3.1 Trend of SoC clock frequency and number of cores.....	39
3.2 Simple schematic of switch-mode regulator .....	40
3.3 Interleaved DC-DC converter architecture .....	41
3.4 Current ripple cancellation using interleaved architecture .....	42
3.5 Value of filter inductor as a function of frequency .....	43
3.6 Efficiency of voltage regulator for different Q inductors.....	44
3.7 Plot of projection of inductor area with number of cores.....	44
4.1 Structure of stripline inductor to derive analytical equations.....	50
4.2 Comparison of analytical model for varying permeability .....	51
4.3 Comparison of analytical model for varying film thickness.....	52
4.4 Comparison of analytical model for varying conductivity of magnetic film .....	52
4.5 Comparison of analytical model for varying dielectric gap.....	53
4.6 Simulation structures of bare, rings, stripes and magnetic dots.....	54
4.7 The complex permeability spectra of patterned NiFe laminations.....	54
4.8 Bare Inductor .....	56
4.9 Inductor with magnetic film .....	57
4.10 Inductor with patterned NiFe dots .....	57
4.11 Measured and simulated L and Q for different inductor structures.....	58
4.12 Equivalent circuit model to extract inductance.....	59
4.13 3D EM simulation showing current density in magnetic film and dots .....	60



Figure	Page
4.14 Optical images of inductors with magnetic films and rings.....	61
4.15 B-H hystereis loop and process flow .....	62
4.16 Measured and simulation of L and Q vs. frequency for magnetic rings.....	63
4.17 Simulation plots of normalized Q vs. L density .....	64
4.18 Measurement result showing relative gain in L and Q at high frequencies.....	65
4.19 Dependence of inductance and Q factor with width of magnetic bar .....	65
4.20 Optical images of rectangular and stripline inductors with CoZrTa .....	67
4.21 Measurement of L and Q vs. frequency for 2 turn inductor with CoZrTa .....	68
4.22 Measurement of La dn Q vs. frequency for stripline inductor with CoZrTa.....	69
4.23 Plot showing inductance requirement of on-chip DC-DC converter.....	71
4.24 Circuit schematic of 3-level DC-DC converter from [6] .....	72
4.25 Chip micrograph of 3-level DC-DC converter from [6].....	73
4.26 Top view of single turn inductor with magnetic rings .....	74
4.27 Inductance and Q factor vs. frequency of magnetic inductor .....	74
4.28 Real impedance of magnetic inductor with frequency.....	75

## Chapter 1

### INTRODUCTION

The semiconductor industry has seen tremendous growth owing to the consistent scaling of device dimensions leading to lower costs, better performance and higher density. However, power dissipation has increased due to higher clock speeds and very high leakage currents at advanced technology nodes (sub-90nm). Power or energy consumption have become one of the major design considerations in today's complex digital integrated circuits. Low power design techniques right from system/software level to microarchitecture, circuit design and device level are required to combat the increasing demands from modern system on a chip (SoC) based integrated circuits found in popular consumer electronics such as netbooks, tablets and smartphones.

Listed below are some of the issues that motivate low power design of integrated circuits

1. In the case of large data centers and server farms, power dissipation is typically in the order of MWs, which is transformed into heat. This requires extensive design to cool the room and the system racks, which costs additional energy.
2. On-chip power dissipation results in temperature increase of the die. Elaborate methods such as water-cooled heat-sinks and on-board fans are required to control the temperature and keep the chip operational and reliable.
3. With the advent of mobile devices such as laptops, tablets and smartphones that are battery-operated, the available energy is fixed and it is imperative to design low power circuits to increase the battery lifetime or the time between recharges.

Fig 1.1 shows the increase in power consumption of desktop and mobile processors manufactured by Intel and AMD in the last 20 years. As can be clearly seen microprocessors manufactured by both companies hit the power wall at around 100+ watts, beyond which traditional cooling methods stopped working. In order to keep up with Moore's law, keep power dissipation within limits and cater to a larger growing mobile market, processors with

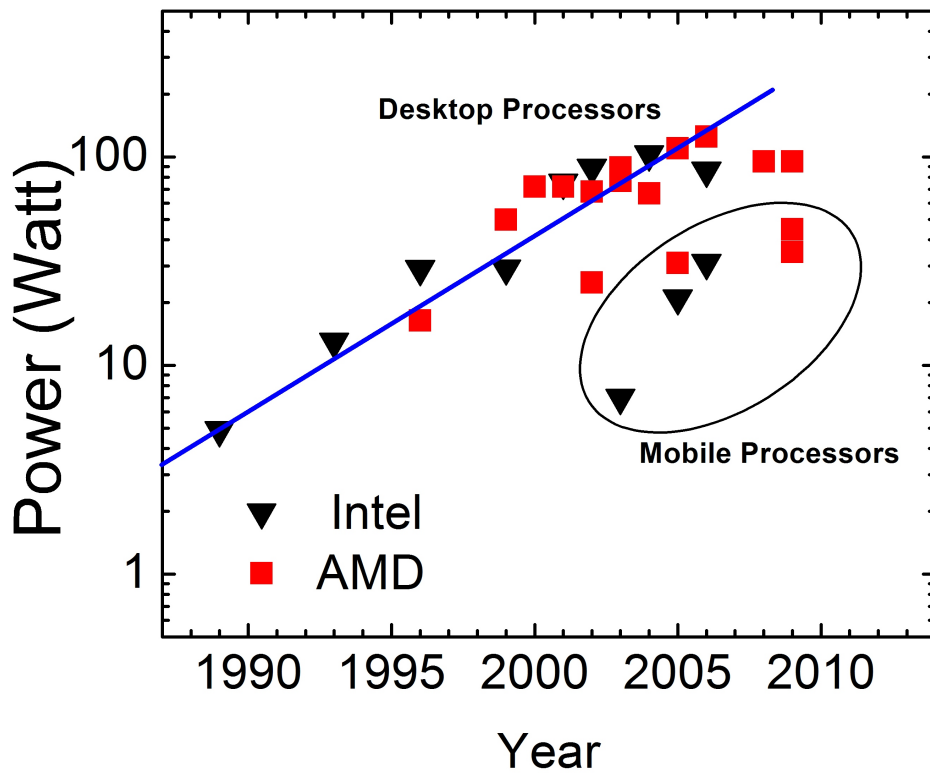


Figure 1.1: Power dissipation of desktop and mobile processor with time. (Source: wikipedia)

multiple cores at lower clock frequency, lower pipeline depth, aggressive clock gating incorporating dynamic voltage and frequency scaling were manufactured. The maximum power dissipation of desktop and server processors is limited to about 130 watts. A new segment of processors for the mobile computing market such as laptops are designed for significantly lower power dissipation limited to 25-35 watts.

In the mobile/hand-held device landscape the power envelope is limited to about 1 watt for processor speeds up to 800MHz to 1.5GHz. Future SoCs are slated to run at speeds of 2GHz, have multiple cores with prolonged battery life. To keep up with these specifications, it is crucial to develop low power schemes at every level of hierarchy in the design of consumer electronic products beginning with novel devices and structures, cutting-edge circuit design techniques, power-aware micro-architecture with tightly integrated software development.

The goal of this dissertation is to explore multiple techniques for low power design beginning with system-level dynamic power management to device level research to improve silicon real-estate. Additionally, we study and predict the feasibility of fine-grained power

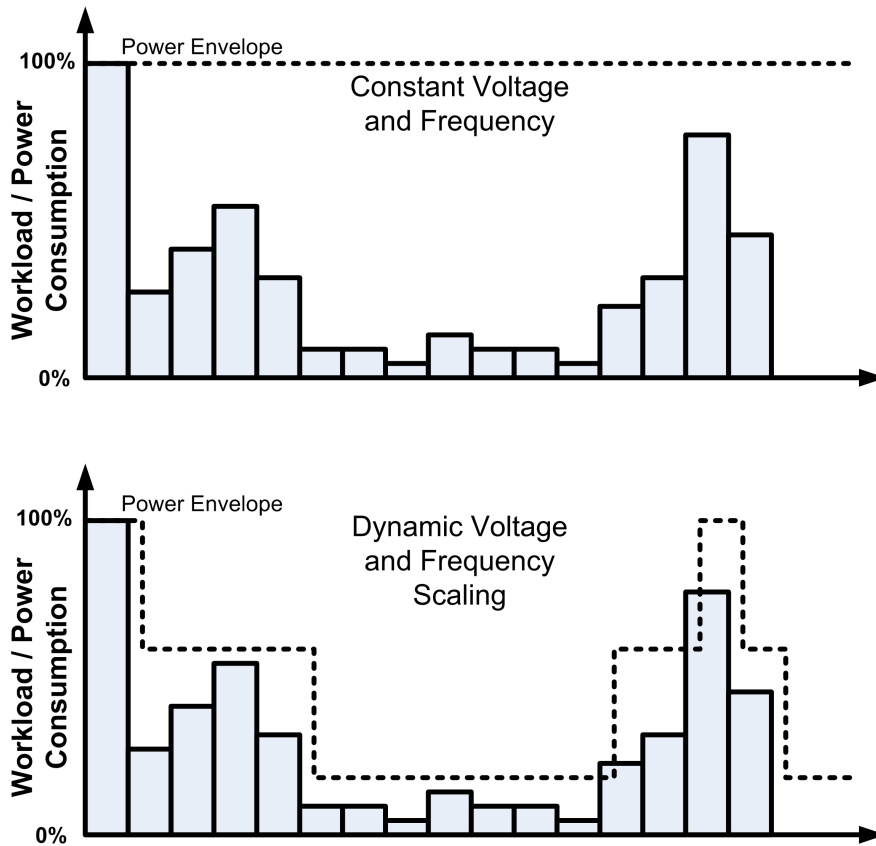


Figure 1.2: Traditional constant supply and frequency power scheme vs. dynamic voltage and frequency scaling.

management techniques for multi-core architectures for the future. The next few sections will give an introduction to the topics covered in the dissertation starting with dynamic voltage and frequency scaling for reducing dynamic power.

### 1.1 Dynamic Power Management

Microprocessors are typically designed to operate at a target frequency ( $f$ ) and supply voltage ( $V_{DD}$ ) for peak performance. Active power dissipation in a CMOS circuit is given by

$$P = \alpha C V_{DD}^2 f \quad (1.1)$$

where  $C$  is the total switched capacitance and  $\alpha$  is the switching activity factor. However, peak performance is not always required and significant energy savings can be achieved by dynamically varying the operating voltage and frequency of the processor to meet the instantaneous workload demand. This method to dynamically reduce power dissipation is

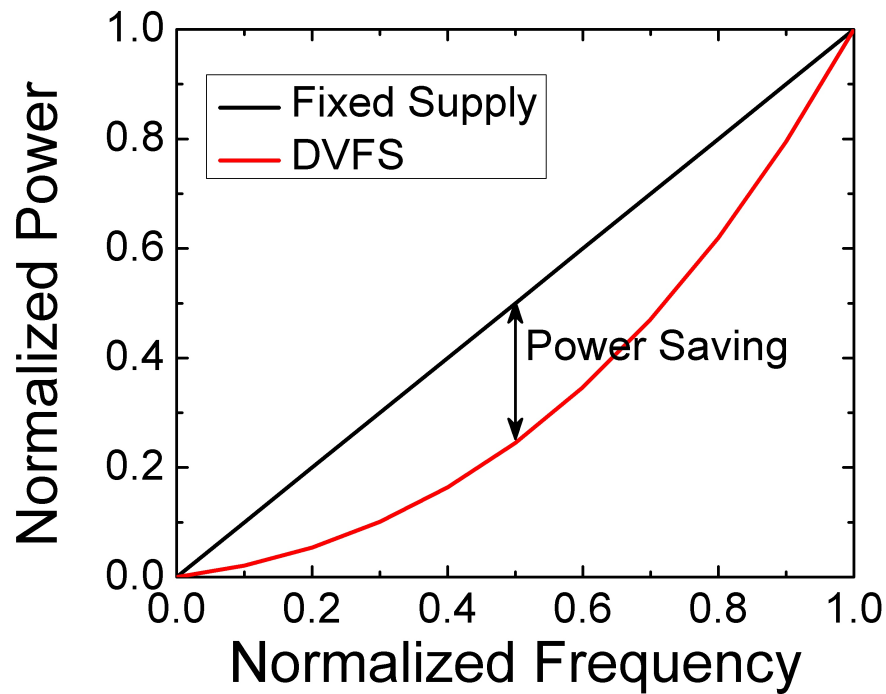


Figure 1.3: Power saving due to DVFS scheme.

known as dynamic voltage and frequency scaling (DVFS). Most DVFS schemes proposed ‘react’ to the changing workload and appropriately scale  $V$  and  $f$  and are implemented in software. Fig. 1.2 and 1.3 shows a qualitative plot of implementing DVFS with workload requirement and the corresponding power savings.

Modern microprocessors are complex systems with multiple cores, non-trivial architectures, employing millions of aggressively scaled sub-100nm transistors designed under extreme power, speed and reliability constraints. An efficient control scheme to manage dynamic power, thermal limits, reliability, error detection-correction and workload scheduling is required.

## 1.2 Neuromorphic Power Management

For efficient power savings we present a hardware based neuromorphic controller for ‘predictive’ dynamic voltage and frequency scaling. Neuromorphic engineering is an emerging interdisciplinary research field that takes inspiration from biology, physics, mathematics, computer science and integrated circuit design with the goal of emulating neural computational architectures. However, unlike conventional approaches of implementing

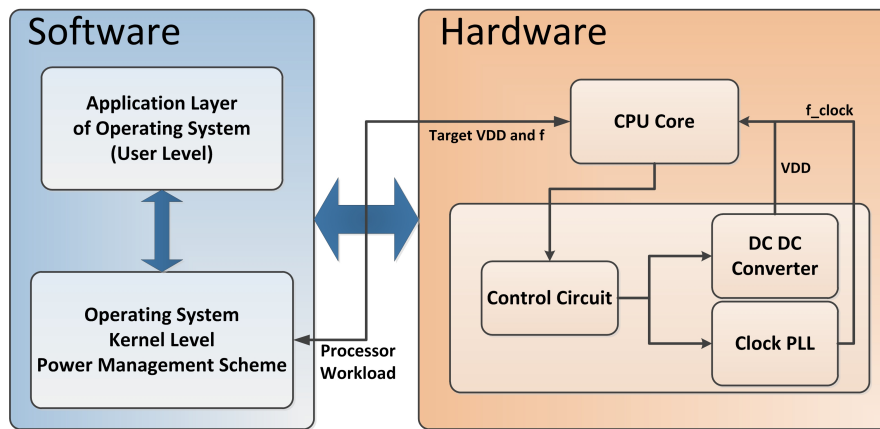


Figure 1.4: Block diagram showing a software based power management scheme.

neuromorphic systems to emulate biological counterparts, we have attempted to explore the potential of the neuromorphic architecture for processor power and thermal management.

The main features of the controller are:

1. The controller uses core workload, frequency, temperature and CPU performance counter values to dynamically manage power and core temperature.
2. The controller employs a single-neuron recurrent network to predict workload and temperature of the microprocessor based on the past workload/temperature profile.
3. Predicted workload and temperature allows pre-emptive scaling of voltage and frequency with high accuracy. Worst case accuracy of the neuromorphic controller is as good as an OS/software based DVFS scheme.
4. Analog and digital versions of the controller are implemented to compare performance metrics.
5. A complete hardware implementation reduces computational load on the processor.
6. The neuromorphic controller has been found to be suitable for fast transitioning workloads.

The details of the neuromorphic controller, its digital and analog implementation, accuracy and impact on power savings are presented in Chapter 2.

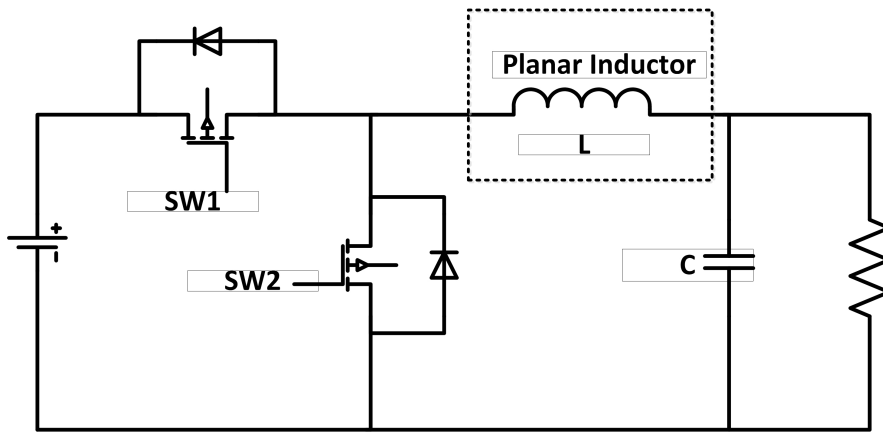


Figure 1.5: Simplified schematic of a switch-mode DC-DC converter.

Fig. 1.4 shows a high-level block diagram of a typical software based dynamic voltage and frequency scaling scheme. Based on the processor workload a target voltage and frequency value is chosen by the software (or hardware). The control circuit signals the PLL to lock to the target frequency. Subsequently the DC-DC converter is signalled to change the supply voltage. In order to implement DVFS efficiently, the frequency and voltage transition needs to be smooth with minimal timing overhead and voltage drooping. With the advent of multi-core architecture finer granularity DVFS and per-core on-chip voltage regulator implementation has become an area of active research. The next section gives a short introduction to switch-mode DC-DC converters with emphasis on on-chip implementation.

### 1.3 Voltage Regulators

With increased scaling and aggressive integration of multiple cores and SOCs for mobile devices, it is common to have multiple supply and frequency domains on the same chip. On-chip distributed power management is essential for reliable low power operation, dynamic voltage scaling and increasing battery life of mobile devices[3],[7]. DC-DC converters can be realized as linear, switch-cap or switched mode inductive converters. Typically switched mode inductive converters are used because of their high efficiency. However, the inductor is generally implemented off-chip or on-package because of their prohibitively large size[8]. Using interleaved topologies, the inductor requirements can be relaxed allowing on-chip implementation[9],[10]. The number of inductors increases as the number of interleaved stages and there is no real area advantage.

On-chip inductors occupy areas larger than a few hundred microns in diameter with inductance values in the range of a few nanohenries. In order to realize distributed power delivery and management for SOCs, it is necessary to explore methods to reduce the area of on-chip inductors with similar or better performance. One of the most common approach involves integrating magnetic materials with on-chip spiral inductors. Switch-mode DC-DC converters with integrated magnetic core inductors have been reported and studied in [11, 12, 13]. Our goal is to undertake a quantitative study of DC-DC converters utilizing scaled magnetic inductors in the neuromorphic DVS loop. In the next section, an introduction to scaling on-chip spiral inductors using integrated magnetic materials is presented.

#### 1.4 On-chip Spiral Inductors with Integrated Magnetic Materials

The electronics industry has seen tremendous growth in terms of performance, scalability and market owing to aggressive scaling of on-chip active devices, namely CMOS transistors at an exponential rate. However, passive elements, especially on-chip spiral inductors have not scaled at the same rate. The inability to effectively scale on-chip spiral inductors is one of the major hurdles preventing realization of monolithic system-on-a-chip.

On-chip spiral inductors with integrated magnetic materials have been explored for decades as a viable option to scale the inductor dimensions with similar or better performance metrics. Numerous structures such as single or double magnetic thin film above and below the inductor metal layer have been explored and implemented. Other structures including solenoidal, toroidal and MEMS based inductors with magnetic cores have also been presented. Fig. 1.6 shows a possible structure of a CMOS compatible spiral inductor with integrated magnetic material surrounding the metal lines.

Inductance increase in the range of 2X to about 20X is possible when the magnetic material completely surrounds the inductor metal wire, thus providing a continuous high permeability path to the magnetic flux. Detailed simulations in 3D EM solvers such as Ansys HFSS have been conducted to identify the optimum material properties and structure that gives maximum increase in inductance and quality factor. Based on the target application different structures and materials have been identified to provide the best results. Standalone spiral inductors with and without magnetic materials were fabricated and measured to validate the



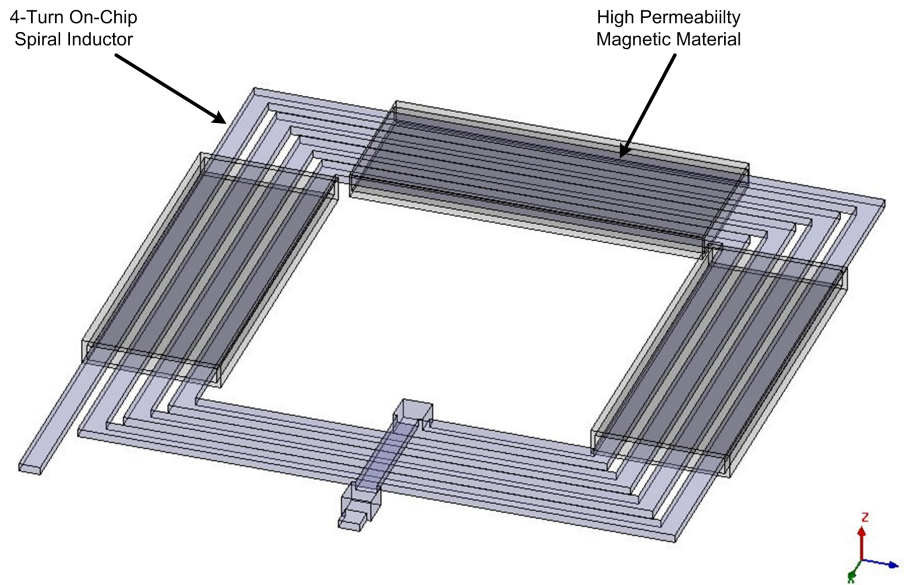


Figure 1.6: On-Chip Spiral Inductor with integrated magnetic material.

simulation study. In addition to design, simulation and measurement, a closed form analytical model for stripline inductors with magnetic materials is developed to match simulation results. Inductors with integrated magnetic materials can be used to improve numerous application utilizing on-chip inductors. Some typical applications include on-chip DC-DC converters (frequency of operation in the sub 300 MHz range) and resonant clock distribution for microprocessors (frequency of operation greater than 1 GHz). Further details of magnetic inductor characterization and potential applications are provided in Chapter 3.

### 1.5 The Integrated View

The goal of this dissertation is to integrate research and results at multiple-levels of design hierarchy from on-chip spiral inductors and their applications in voltage regulators to circuit-design and block-level simulation of neuromorphic DVFS controllers and provide a comprehensive view of how changes at each level impacts power-performance metrics at a system level. Fig.1.7 shows an integrated view/scope of this work which is as follows:

1. From a system-level perspective dynamic voltage and frequency scaling and thermal management is explored in battery operated laptop with a state-of-the-art commercial microprocessor as discussed in Chapter 2.

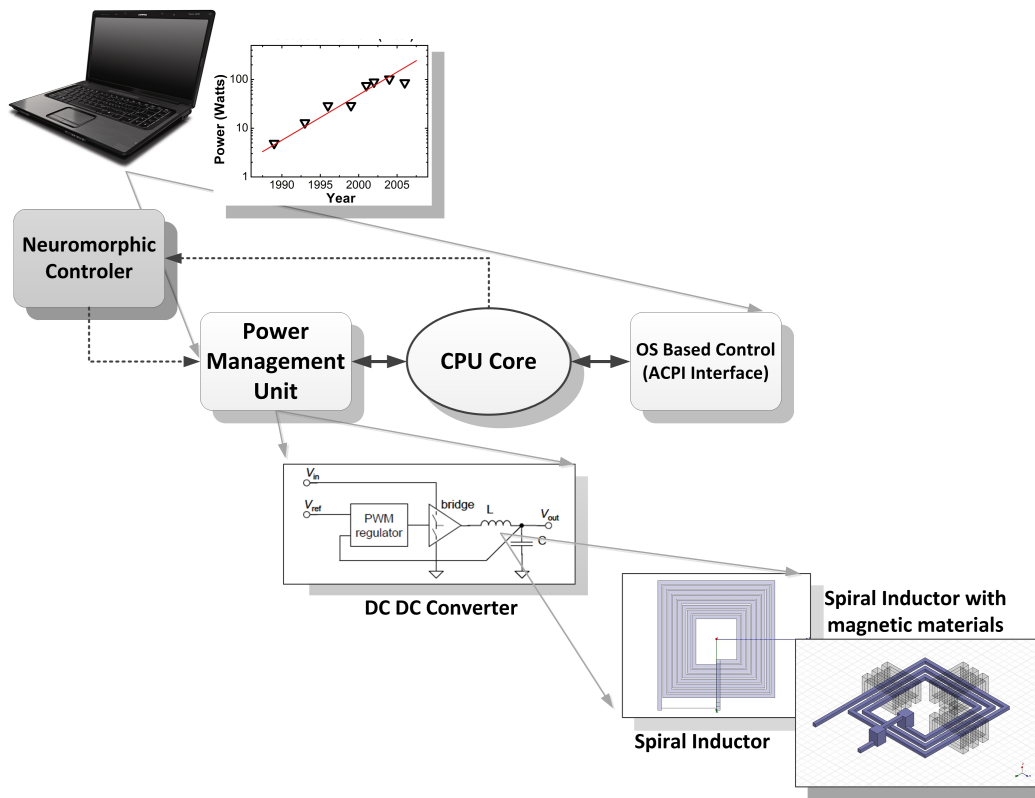


Figure 1.7: Scope of this work.

2. A system level study is conducted to understand the operational requirements of on-chip dc-dc converters for a neuromorphic DVFS and thermal predictive controller, especially in terms of inductance parameters of the output filter is presented in Chapter 3.
3. At a device-level on-chip spiral inductors with integrated magnetic materials are characterized and potential applications such as resonant clock distribution are explored and performance improvements and trade-offs are evaluated. The potential of on-chip spiral inductor with integrated magnetic materials for on-chip voltage regulators used for per-core fine-grained DVFS is explored. Further details are presented in Chapter 4.
4. The final chapter concludes the thesis with remarks on the design trade-offs and opportunities for future work.

### NEUROMORPHIC POWER CONTROLLER

Neuromorphic engineering is an emerging interdisciplinary research field that takes inspiration from biology, physics, mathematics, computer science and integrated circuit design with the goal of emulating neural computational architectures[14]. Accurate modeling and simulation of cognitive processes is a major challenge eluding scientists since the advent of digital computers. In contrast to precise and serial architecture of Boolean logic computers, the nervous system implements robust, parallel, event-driven and reliable computation with unreliable components and uncertain interconnections. Additionally, there is extensive use of adaptation, learning and the ability to handle non-linear and complex tasks[15]. Hence, neuromorphic engineering typically aims to design electronic neural systems that perform tasks similar to their biological counterparts such as silicon neurons and synapses[16, 17], retinas[18] and cochleae[19].

Modern day microprocessors are complex systems with multiple cores, non-trivial architecture, employing millions of aggressively scaled sub-100nm transistors designed under extreme power, speed and reliability constraints[20]. An efficient control scheme for dynamic power, thermal management, reliability, error detection-correction and workload scheduling is required. From Fig. 2.1 it is apparent that neuromorphic architectures are more suitable for complex tasks such as classification and prediction of noisy nonlinear data as compared to Boolean logic architecture.

In this work, we extend the design of our previous neuromorphic DVS controller[21] for temperature prediction. The controller uses core workload, frequency, temperature and CPU performance counters to dynamically manage power and core temperature. Fig.2.2 shows the workload and temperature profile with corresponding frequency levels from a state-of-the-art microprocessor. The controller employs a single-neuron recurrent network to predict the workload and temperature of the microprocessor based on the past workload/temperature profile and regulate the supply voltage and frequency accordingly.

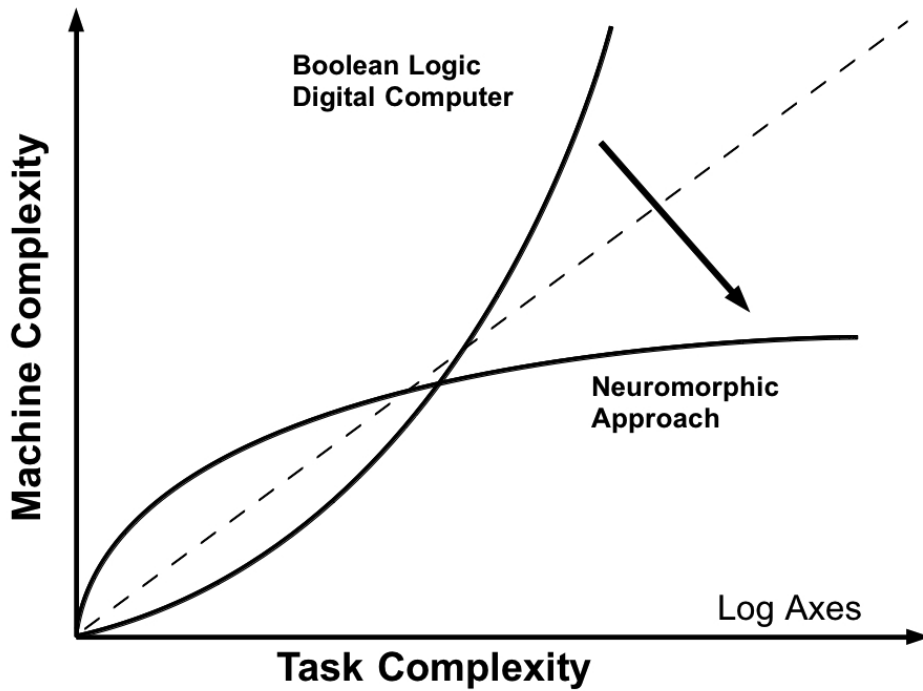


Figure 2.1: Plot of Machine complexity vs. task complexity for Boolean logic computer vs. neuromorphic computing architecture [1].

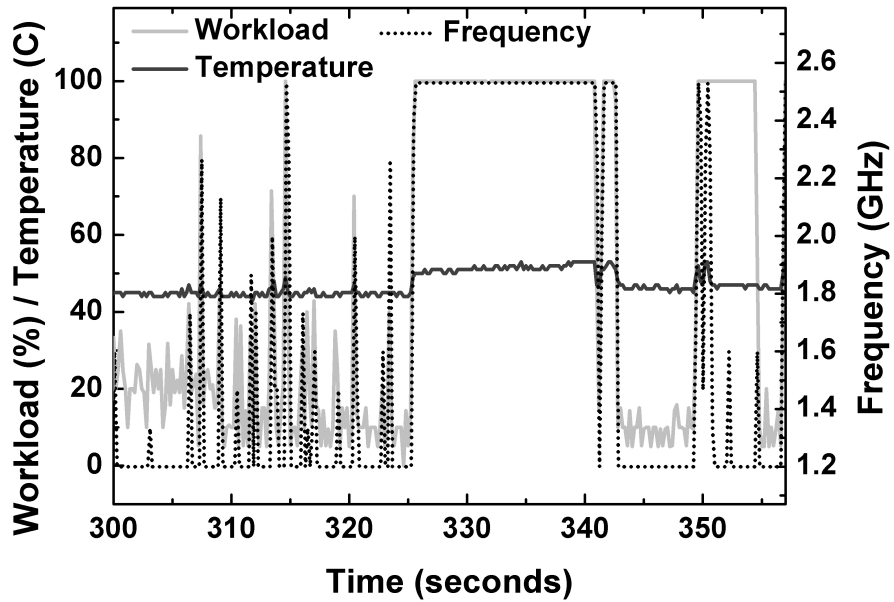


Figure 2.2: Workload and DVFS (Operating Frequency) profile and corresponding temperature of a generic microprocessor.

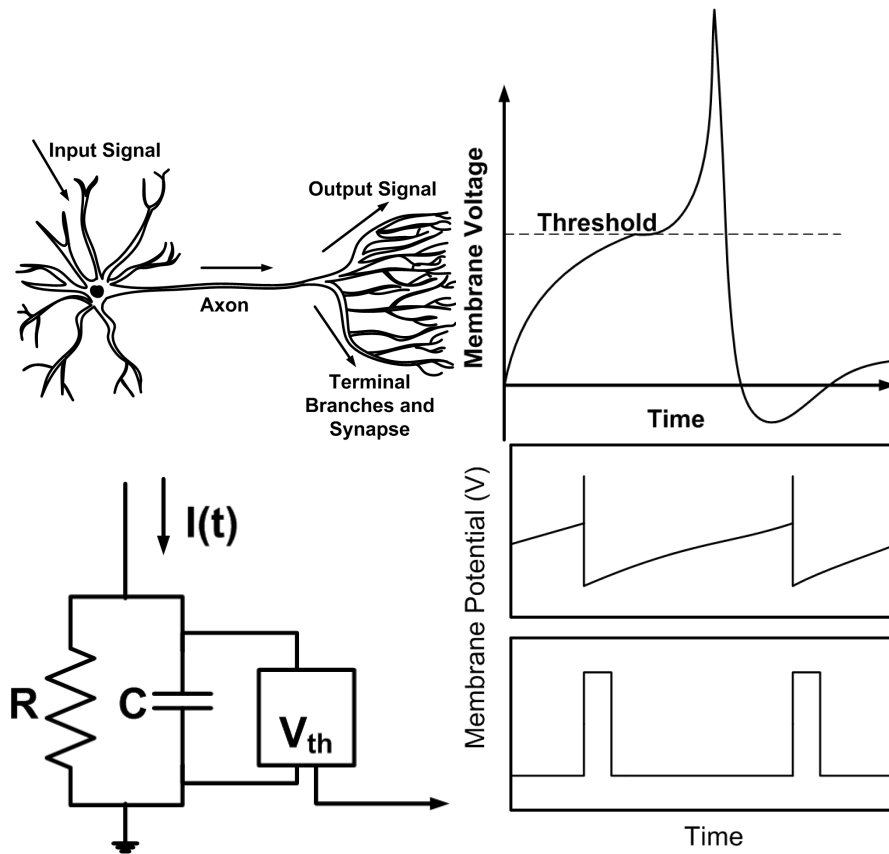


Figure 2.3: (a) Schematic of a biological neuron. (b) Membrane potential of neuron with time. (c) Circuit schematic of a linear integrate and fire neuron. (d) Analog and digital output spikes from the circuit.

## 2.1 Background and Previous Work

### *Spiking Neuron Architecture*

A neural architecture consists of a large number of interconnected neurons with adaptive weights that change with internal or external information flow during the learning phase. First and second generation neurons [22] required weight multiplication with the inputs during learning, making hardware implementation difficult.

Third generation neurons increase the level of biological realism by using individual spikes which allows incorporating spatial-temporal information in communication and computation. These neurons use pulse-coding mechanisms allowing multiplexing of information. Spiking neurons, thus, have been demonstrated to possess higher computational power as compared to the first or second generation neurons. Additionally, weight update

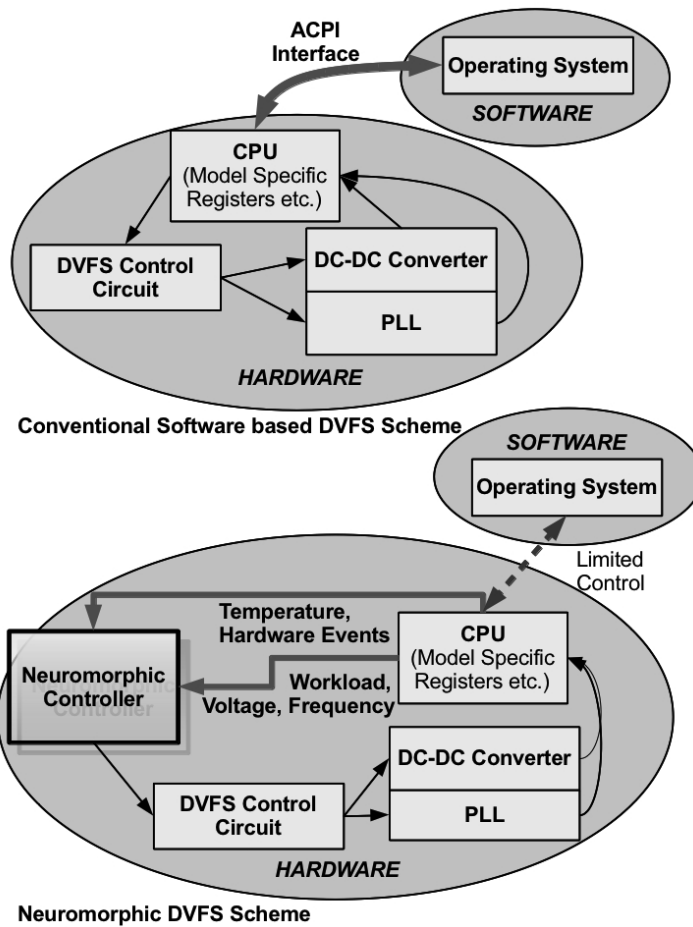


Figure 2.4: Block diagram showing the a conventional DVFS scheme implemented in commercial microprocessors and a comparable Neuromorphic scheme.

involves changing the firing rate of the spiking neurons, making it ideal for hardware implementation. Fig.2.6 shows the schematic of a biological neuron, its equivalent electrical circuit and their membrane potential.

### *Dynamic Power Management*

Microprocessors are typically designed to operate at a target frequency ( $f$ ) and supply voltage ( $V_{DD}$ ) for peak performance. Active power dissipation in a CMOS circuit is given by

$$P = \alpha C V_{DD}^2 f \tag{2.1}$$

where  $C$  is the total switched capacitance and  $\alpha$  is the switching activity factor.

However, peak performance is not always required and energy savings can be achieved by

dynamically varying the supply voltage and frequency of the processor (DVFS) to meet the instantaneous workload demand.

Dynamic voltage and frequency scaling is supported by nearly all commercially available general purpose microprocessors under trademarks such as Intel's Enhanced Speedstep Technology and AMD's Cool'n'Quiet and PowerNow! technologies. These technologies work in using an open standard specification known as Advanced Configuration and Power Interface (ACPI) [23] developed for Operating System-directed configuration and Power Management (OSPM).

Nearly all DVFS schemes proposed are implemented in software, adopting offline profiling to learn the average-case execution time or worst case execution time of tasks and scale voltage/frequency for minimum power/energy consumption [24],[25]. Statistical techniques using CDF and PDF of offline workload profile have been proposed [26]. Power management schemes to predict the future workload based on past workload profile using adaptive and non-adaptive filters have been proposed in [27], [28]. Other schemes include a neural network based power management controller [29], control theory based predictive scheme [30] and one using online learning algorithm [31]. However, all these schemes are implemented in software and they incur an overhead at run-time. Some of these schemes require elaborate workload profiling or complex calculations to find the appropriate  $V_{DD}$  and  $f$  combination and are not appropriate for online implementation. Since the time-varying workload profile is highly nonlinear and aperiodic in nature, it serves as a good testbench to explore the potential of a neuromorphic controller.

#### *Dynamic Thermal Management*

Elevated temperatures severely impact processor performance, increase leakage power dissipation and under extreme conditions, can cause the microprocessor circuit to break-down. All state-of-the-art microprocessor circuits have temperature sensors to estimate the on-chip temperature and can take reactive measures when the temperature reaches a specific threshold.

Numerous reactive and predictive schemes for thermal management have been implemented. Predictive schemes involve workload phase detection [32], neural networks[33] and ARMA (autoregressive moving average) filter based prediction[34]. These schemes

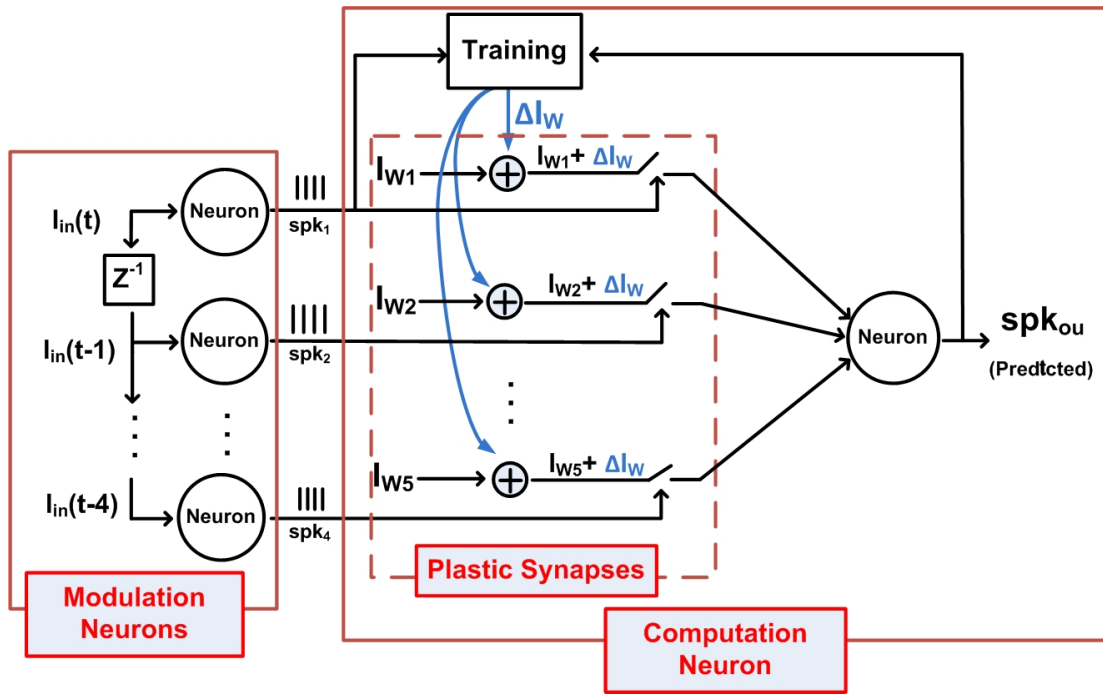


Figure 2.5: Block diagram of the Spiking Neuromorphic Controller.

predict the temperature from workload decomposition and readings from performance counters. Proactive steps involve  $V_{DD}$  and  $f$  scaling, temperature aware task scheduling [35]etc. As in the case of DVFS schemes, all these methods require elaborate calculations, are implemented in software and incur computational overhead on the processor.

The Intel Core i5 processor, used in our experiments, employs a reactive thermal management scheme. When the core temperature, monitored by a digital thermal sensor (DTS), exceeds its maximum junction temperature ( $T_{j,max}$ ), the thermal control circuit (TCC) is activated. The two main reactive steps are reducing the operating frequency and core supply voltage (DVFS) and modulating (or gating) the internal processor core clock. Clock modulation is done by alternately turning the clocks off and on at a duty cycle specific to the processor. Processor performance is decreased by the same amount as the duty cycle. Details of the thermal management steps are provided in the processor data-sheets [36]. A simplified block-diagram of a commercially implemented or a software/OS controlled DVFS scheme is given in Fig. 2.4 along with the block-diagram of a complete hardware based neuromorphic controller. The hardware based controller can work independently or in conjunction with a software based scheme as demonstrated in [37].



## 2.2 CMOS Implementation of Neuromorphic DVS Controller

A block level diagram of a recurrent neuromorphic circuit that predicts the future workload/temperature values based on the past statistics is given in Fig.2.5. The input neurons store the past workload/temperature values which are transmitted to the computation neuron as input. Based on the synapse weights and bias, the computation neuron calculates the next workload/temperature value. The error generated by subtracting the predicted value and the actual value is used to modify the synaptic weights before the next prediction takes place.

A detailed implementation of this controller in analog and digital domains has been discussed in[21]. The details of the digital design and modifications for temperature prediction are described in the next subsection.

### *Analog Design*

An analog spiking neuromorphic controller is implemented in a 45nm CMOS technology using PTM model files<sup>1</sup>[38]. The circuit occupies an area of about  $15 \times 15 \mu\text{m}^2$  and the power consumption is less than  $150\mu\text{W}$ . The circuit consists of analog neurons for modulation and computation, plastic synapses and training.

### Neurons

The building blocks of the controller are simple integrate-and-fire electronic neurons, functionally equivalent to those described in[39]. Network of integrate and fire neurons have been shown to exhibit a wide range of useful computational properties, including feature binding, segmentation, pattern recognition, onset detection, input prediction,[40, 41, 42] etc. A schematic diagram of the circuit implementing the neuronal dynamics is shown in Fig. 2.2. These neurons integrate linearly the total afferent current and when a threshold is crossed they emit a spike. The sub-threshold dynamics can be described by the equation governing the voltage across a capacitor (which represents the membrane potential of the cell) as follows

$$V_m(t + \Delta t) = V_m(t) + \frac{I_w \Delta t}{C}$$

where  $V_w(t)$  is the voltage on the membrane,  $C$  is the soma capacitance( $C_{m1} + C_{m2}$ ),  $\Delta t$  is the time step of  $spk_{in}$ , and  $I_w(t)$  is the postsynaptic current injected. A positive feedback loop,

---

<sup>1</sup>Analog design and simulation by Jounghyuk Suh.

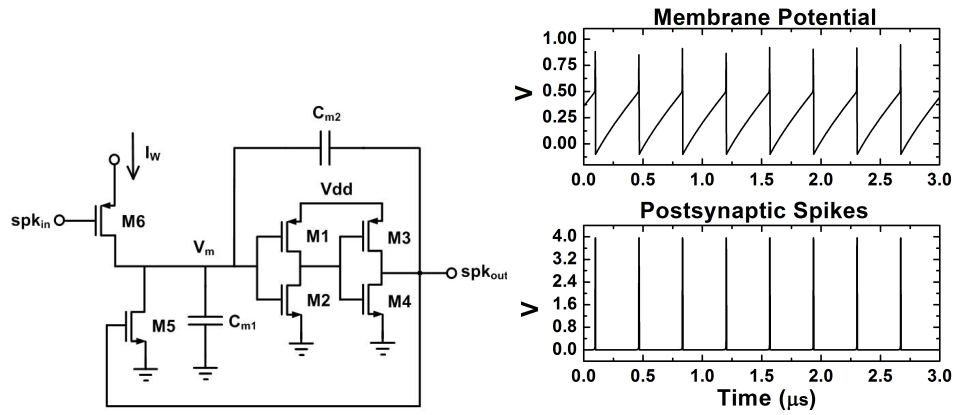


Figure 2.6: (a) Schematic diagram of the integrate and fire neuron circuit (b) Membrane potential of neuron with time.

implemented by the capacitive divider  $C_{m1}-C_{m2}$ , increases  $V_m$  by  $V_{dd}C_{m2}/(C_{m1}+C_{m2})$  [39]. As long as  $V_{spk_{out}}$  is equal to  $V_{dd}$ , the digital switch M5 is kept closed, and the current can discharge the two capacitors causing the membrane potential to decay linearly. As  $V_m$  crosses again, the switching voltage ( $V_{dd}/2$ ) of the inverter M1-M2, the output voltage  $V_{spk_{out}}$  goes back to the ground level (spike inactivation), the integration of the input current can then start again (Fig. 2.2).

**Modulation** In pulse-based (spike) techniques signals take the form of a train of pulses, usually with the signal in an inactive (zero) state most of the time. Such signals have low power (assuming that power consumption is minimal during the zero periods), reasonably noise immune, and easily regenerated if the pulse edge is flattened. The three basic pulse modulation techniques are: pulse height modulation, pulse width modulation, and pulse frequency modulation [43]. Fig. 2.2 shows an integrate and fire neuron performing pulse width modulation. By using the clock signal instead of  $spk_{in}$ , the desired duration of pulses are obtained. As shown in Fig. 2.7, we primed the workload into currents. As the workload increases, the current will increase, and vice versa as the workload decreases the current will decrease proportionally. The simulation result shows an appropriate range to set; this set scale is from 500n to  $15\mu A$ . In this scope, the current travels through the neuron to create a 3MHz to 36MHz pulse encode, shown in Fig. 2.7. The interval between two spikes depends on the input currents ( $I_w$ ). The spike rate can be calculated in the simple case of constant. It is the

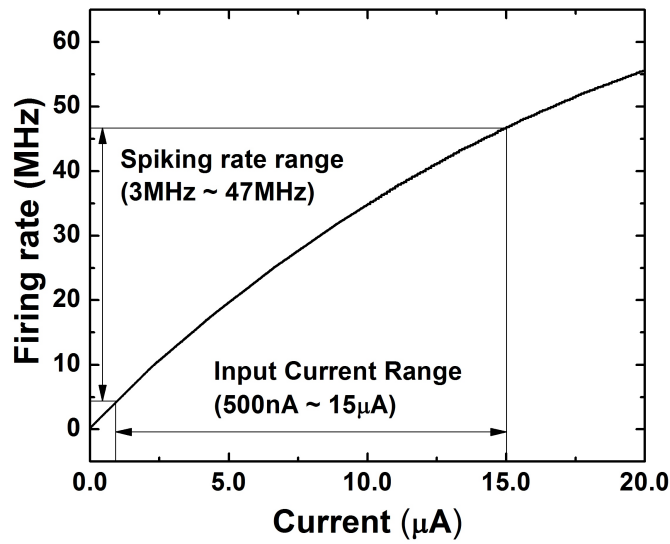


Figure 2.7: Current to spike rate modulation range

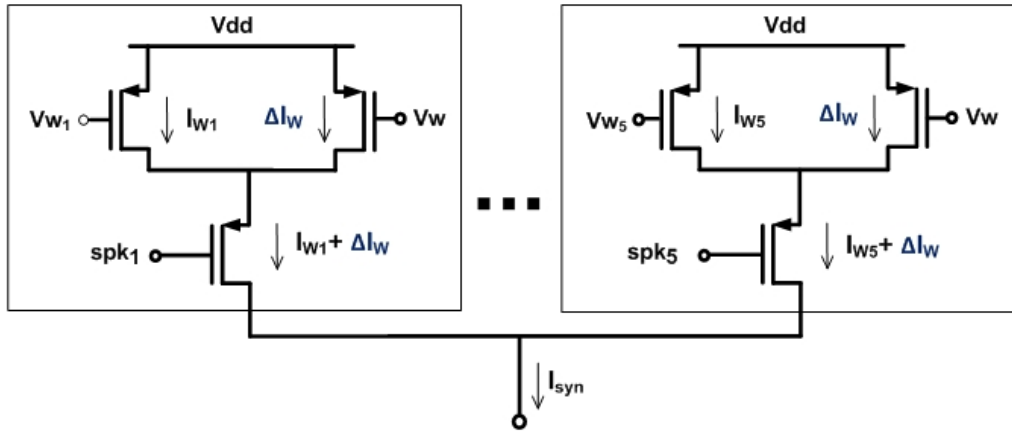


Figure 2.8: Schematic diagram of the plastic synapses

time needed to the membrane potential to reach the switching voltage of inverter M1-2( $V_{dd}/2$ ) starting from the reset potential ( $V_{dd}/2 - V_{dd}C_{m2}/(C_{m1} + C_{m2})$ )

$$firerate = V_{dd} \frac{C_{m2}}{I_w}$$

The component variables  $C_{m1} = C_{m2} = 200\text{fF}$  and  $V_{dd} = 1\text{V}$  are used for this simulation

**Computation** A computation neuron is an excitatory neuron connected by plastic synapses.

The role of synaptic circuits is to convert pre-synaptic voltage into post-synaptic currents injected in the membrane of the computation neuron, with the gain typically termed as the

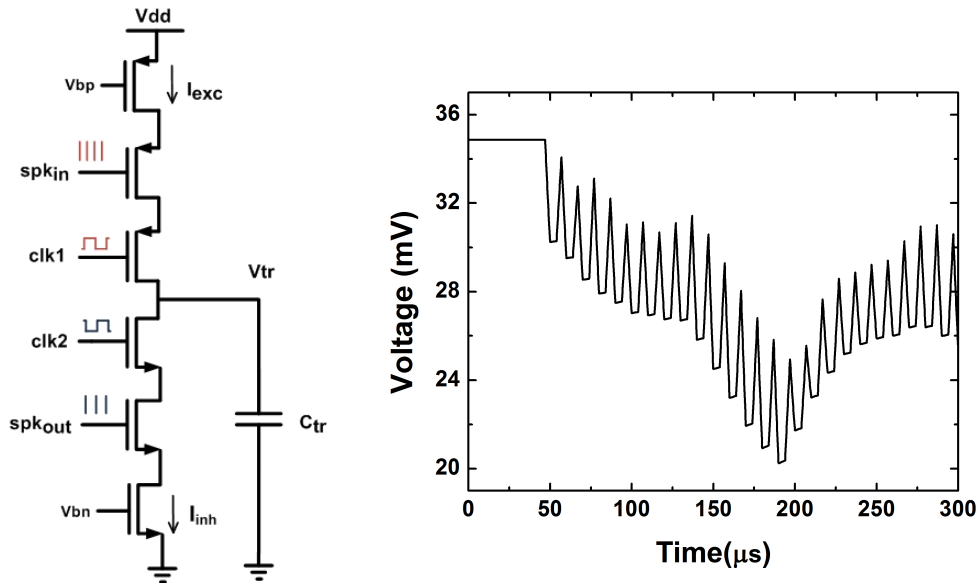


Figure 2.9: (a) Schematic diagram of the training circuit (b) The updates in synaptic weight

synaptic weight. The intensity of the current injected on the neuron is determined by the excitatory postsynaptic current (EPSC). The schematic diagram of the EPSC is shown in Fig. 2.8. The synaptic current is injected into the post-synaptic soma capacitor only upon the occurrence of a pre-synaptic spike ( $spk_{1-5}$  and transistors act as digital switch). The total excitatory afferent current to the neuron ( $I_{syn}$ ) is the sum of all the  $\Delta I_w$  contributions, plus assigned currents.

$$I_{syn} = I_{w1} + I_{w2} + I_{w3} + I_{w4} + I_{w5} + 5\Delta I_w$$

where  $I_{w1}$  is the constant current for the most recent input synapse and  $I_{w2-5}$  is the constant current for the delayed input. Since the latest input has the indication for the next input, the strength of  $I_{w1}$  which starts out as the strongest becomes increasingly weaker as the length of delay becomes longer. This can increase the accuracy of the prediction. The current  $\Delta I_w$  is set by the bias voltage  $V_w$  which is determined by the training procedure.

### Training

The prediction process performed by neural networks is the result of a training procedure, during which the synaptic strengths between neurons are modified according to a learning prescription [42]. A training circuit which is efficiently able to modify and preserve the weights is presented in Fig. 2.9. Transistor M1-M6 and capacitor  $C_{tr}$  implement the training

block. Each transistor acts as digital switches and the current can flow only when a pre-synaptic spike ( $spk_{in}$  and  $spk_{out}$ ) is active.  $Spk_{in}$  is a pre-synaptic spike which is determined by the next input (target).  $Spk_{out}$  is the output of the computation neuron. The sign of the current is determined by the firing rate of those two pre-synaptic inputs. When  $clk1$  is high, the current  $I_{biasp}$  charges the memory capacitor ( $C_{tr}$ ). When  $clk2$  is high, the current  $I_{biasn}$  discharges the memory capacitor. The voltage difference ( $\Delta V_{tr}$ ) can be written as

$$\Delta V_{tr} = \frac{1}{C_{tr}} \int_{t_0}^{t_1} I_{exc}(t) dt + \frac{1}{C_{tr}} \int_{t_0}^{t_2} I_{inh}(t) dt$$

where  $t_1$  is the time duration of the clock signal  $clk1$  and  $t_2$  is the time duration of the clock signal  $clk2$ . The time variation of this voltage difference is shown in Fig. ???. This is the result of the training and it determines the bias voltage of synapse circuits.

### *Digital Design*

For digital design, a discrete-time version of the analog architecture is used. Workload at time  $t$  is assumed to depend on the past workload profile. The following equation describes the prediction function

$$y(t+1) = w[y(t) + 2^{-1}y(t-1) + \dots + 2^{-n}y(t-n)] + \eta \quad (2.2)$$

where  $w$  is the weight of the input,  $\eta$  is the bias and  $n$  is the number of history samples. The workload values are scaled by a factor of  $2^{-n}$  to ensure that the most recent history has the maximum impact on the prediction. Implementing the scaling simply involves binary shift operation. We assume that the sum of the scaled workload values are available as input. We have used 6 bits to denote workload values as integers between 0-100 %.

A block level schematic of the digital design is shown in Fig. 2.10. The design is implemented using two separate clocks,  $sclk$  running at 1 MHz and  $fclk$  running at 200 MHz.  $fclk$  is used for spike generation while  $sclk$  synchronizes the controller. The spike generator consisting of a counter (cntr1) and digital comparator that generates number of spikes is equal to the scaled sum of the workload history.

The second set of counter (cntr2) and comparator in the design form the adaptive spiking neuron where plastic membrane threshold is used for learning. When the number of

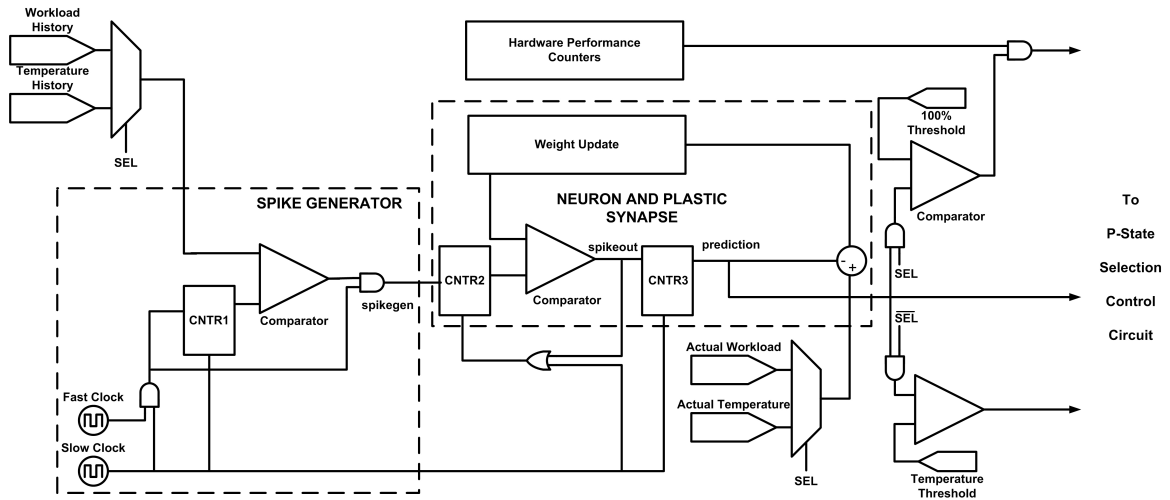


Figure 2.10: Block level schematic of the digital neuromorphic controller circuit. The design uses counters and digital comparators and the processor clock for generating the spike.

output spikes reaches the threshold value, the comparator asserts, which is used to generate back the RESET signal for the counter. The final counter (cntr3) counts the output spikes and gives out a 6-bit prediction value after adding the bias value  $\eta$ .

### Error and Training

Error is calculated by subtracting the predicted value from the actual workload value enabling training. There are two methods for adaptive learning: (1) Updating the threshold of the second counter/comparator pair (cntr2) so that the firing rate of ‘outspike’ changes and (2) modifying the bias value,  $\eta$ . Variation in threshold leads to large changes in output spike rate. This ‘coarse learning’ is only required when the error is very large. For most cases, varying the bias value or ‘fine learning’ gives good prediction. If the error is positive, the bias value is increased by an amount  $\eta$  and if it is negative, the bias is reduced by a value  $\eta/2$ . The circuit is biased towards an increasing workload since a history based predictor typically lags behind the actual values by a unit delay.

### Temperature Prediction

The digital neuromorphic controller for workload prediction can be time-multiplexed and used for temperature prediction as well. As shown in Fig. 2.10, the SEL signal to multiplexer is used to select workload and temperature prediction. If the workload value is above a certain threshold (100%), readings from hardware performance counters are sent to the  $V_{DD}$  and  $f$

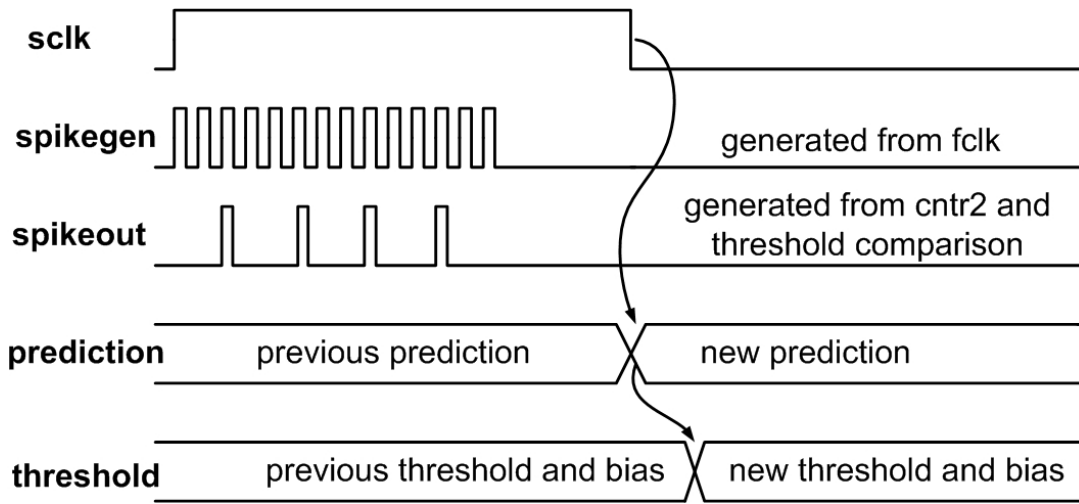


Figure 2.11: Timing diagram showing spike generation, prediction and adaptation phase with respect to the slow clock (*sclk*) at 1MHz.

selection control circuit to determine if the program is memory bound and DVFS can be implemented. During temperature prediction phase, if the temperature value reaches a certain threshold (set much lower than the junction Temperature,  $T_{j,max}$ ), a lower operating voltage and frequency is selected to control power dissipation and reduce the core temperature. A constant positive bias  $\eta$  is applied during temperature prediction to ensure that the predicted temperature is either equal or higher than the actual temperature. Hence, using the neuromorphic controller, the temperature can be lowered proactively and severe steps such as clock modulation that degrade performance, can be prevented.

### Implementation

The complete digital design was synthesized using OpenCell Standard Library [44] based on 45nm PTM model files in Synopsys Design Compiler. The circuit has an active power dissipation of  $9\mu\text{W}$ ,  $3.1\mu\text{W}$  leakage power and occupies  $360\mu\text{m}^2$  area. The timing diagram in Fig. 2.11 shows that spike generation and workload prediction takes place during the HIGH phase of *sclk*. The predicted value is stored in a flip-flop during the falling edge of *sclk* and is used to generate error, new threshold/bias in the LOW phase. The circuit is ready for a new prediction in the next rising edge of *sclk*.

Table 2.1: Processor Specifications

Name	Intel Core i5 460m
Number of Cores	2
Hyperthreading	Yes
Turbo Mode	Yes
Performance Counters	7
Maximum Frequency	2.53 GHz
Minimum Frequency	1.2 GHz
Supply Voltage	1.4V
Number of P-states	11
Max Temperature	105 °C

Table 2.2: Processor Performance Events

Event Name	Event Mask Mnemonic
CLOCK_CYCLES	CPU_CLK_UNHALTED.CORE
INST_RETIRED	INSTRUCTIONS_RETIRED.ANY_P
LLC_MISS	MEM_LOAD_RETIRED.LLC_MISSES
DTLB_MISS	DTLB_LOAD_MISSES.WALK_COMPLETED
ITLB_MISS	ITLB_LOAD_MISSES.WALK_COMPLETED

### 2.3 Workload Characterization and Power Modeling

Each voltage and frequency combination supported by a microprocessor is known as a P-state with P0 corresponding to the highest  $V_{DD}$  and  $f$  value. Detailed specifications of frequency and allowable supply voltage and frequency of the processor is given in Table 2.1. Workload readings were obtained the machine running Linux operating system (kernel 2.6.30) using a custom program that monitors the CPU idle time at regular intervals. Accurate readings up to a resolution of 200ms using the program is possible with negligible load on the processor. We used the ‘ondemand’ frequency scaling governor[45] in Linux that scales the frequency based on workload demand for our experiments. Each supported frequency point has a corresponding operating voltage level, which are used for estimating power savings.

For the Core i5 processor, the Digital Thermal Sensor can be accessed by a Machine Specific Register (MSR), IA32\_THERM\_STATUS. This register stores the DTS value as the difference of junction temperature and the core temperature. The core temperature can be calculated as



$$T_{core} = T_{j,max} - DTS_{value} \quad (2.3)$$

The workload readings denote the percentage of time CPU is busy. However, a memory-bandwidth limited application can show the CPU as busy even when it is not doing any useful work waiting for memory access to complete. During such instances, the CPU P-state can be lowered without any impact on performance [46]. In order to monitor memory access, CPU performance counter data is collected using the *pfmon* command-line utility. The detailed setup is as follows:

- The *pfmon* utility[47] is enabled using a custom-compiled Linux kernel 2.6.30. The program can be setup to access the CPU PMU (Performance Monitoring Unit) to measure performance events such as Unhalted Core Cycles, Number of Instructions Retired, Last Level Cache Miss etc.
- The specific events monitored are:
  1. CPI (Cycles Per Instruction)
  2. LLC\_MISS (Last Level Cache Misses)
  3. DTLB\_MISS (Data TLB Misses)
  4. ITLB\_MISS (Instruction TLB Misses)
- To calculate the percentage of time the CPU spends in memory access or the CPU intensiveness ( $\mu$ ), we use the method described in [31]. The relevant equations[48] are listed below.

$$CPI = \frac{CLOCK\_CYCLES}{INST\_RETIRED}$$

$$CPI_{cache} = \frac{LLC\_MISS \times 130}{INST\_RETIRED}$$

$$CPI_{DTLB} = \frac{DTLB\_MISS \times 30}{INST\_RETIRED}$$

$$CPI_{ITLB} = \frac{ITLB\_MISS \times 30}{INST\_RETIRED}$$

The performance events used in the above equations and their event mask mnemonics for the Core i5 processor are given in Table 2.2. Using the CPI of different events calculated above we can now define  $CPI_{mem}$  as

$$CPI_{mem} = CPU_{cache} + CPI_{DTLB} + CPI_{ITLB}$$

Finally we get the percentage of cycles the CPU spends in useful computation as

$$\mu = \frac{CPI - CPI_{mem}}{CPI} \times 100$$

If  $\mu$  is much lower than 100%, implies that the CPU spends a large percentage of its time in memory access (Cache Miss, TLB Miss etc.) and there is an additional scope to reduce the P-state for power savings without performance hit. Since performance counters are designed into all modern microprocessors, it is reasonable to assume that a dedicated on-chip power and thermal management system would be able to access and utilize them as described in this subsection.

### *Power Modeling*

DVFS schemes are typically explored in software using cycle accurate processor simulators [49] with attached power model simulators [50] which can give estimation of power dissipation and savings with reasonable accuracy. In contrast, numerous schemes are implemented on commercially available machines and power is directly measured using power meters or through readings of current drawn to the machine. This requires some approximation as well since the calculated power includes power dissipation by all the components of the computer and an accurate breakdown of the power consumed by each component.

Since direct power measurement is not possible for our approach and we use real workload and DVFS readings, it is imperative to develop a model to estimate power savings. We have adapted the model from [31] which is reproduced below for clarity. The total power dissipation ( $P_t$ ) of the processor can be divided into two parts: dynamic power ( $P_d$ ) and leakage power ( $P_l$ ) as follows

$$P_t = P_d + P_l \quad (2.4)$$

where

$$P_t = CV_{DD}^2 f + IV_{DD} \quad (2.5)$$

Considering  $P_{max}$  as the power dissipation when the processor is running at maximum  $V_{DD}$  and  $f$ , we can normalize the power dissipation as

$$P_n = \frac{CV_{DD}^2 f}{P_{max}} + \frac{IV_{DD}}{P_{max}} \quad (2.6)$$

Additionally, defining  $\lambda$  and  $\rho$  as the percentage contribution of dynamic and leakage power to the total power dissipation at the highest voltage and frequency we have

$$P_n = \lambda V_n^2 f_n + \rho V_n \quad (2.7)$$

where  $V_n = V/V_{max}$  and  $f_n = f/f_{max}$ . For a state-of-the-art microprocessor, we assume  $\rho$  to be 50% during active operation [51].

Since the number of available voltage levels are typically less than frequency transition points, dynamic power dissipation reduces cubically when both  $V_{DD}$  and  $f$  are scaled and reduces linearly when the processor transitions to higher P-states (only  $f$  scaling) giving linear reduction in power [52]. The normalized power dissipation using the above model for different number of P-states is plotted in Fig. 2.12. In Fig. 2.12, we have assumed number of possible  $V_{DD}$  levels are 3 and frequency levels ranging from 3 to 10 giving a total of 10 P-states.

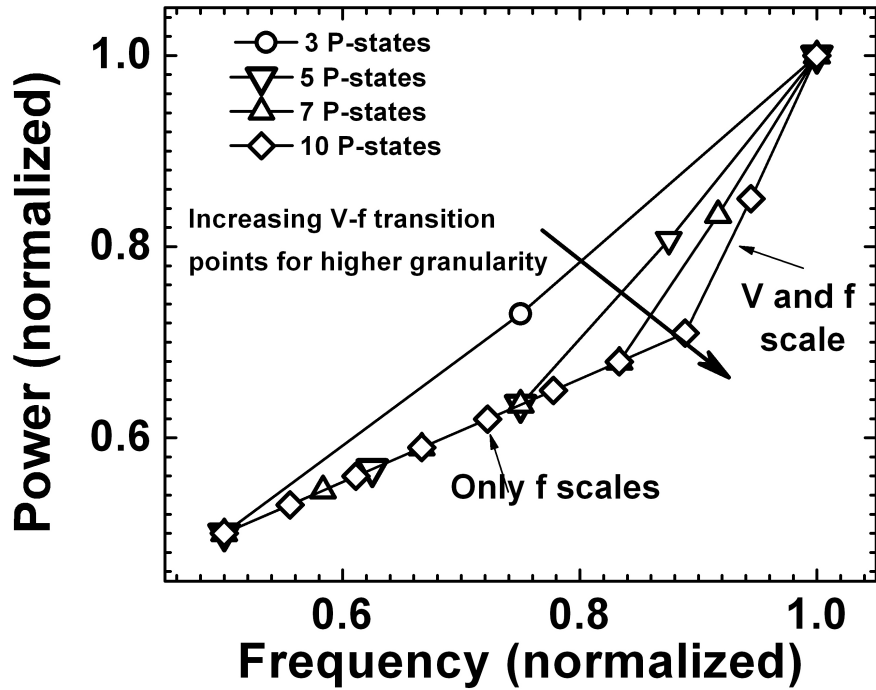


Figure 2.12: Normalized power dissipation vs. normalized frequency for increasing V-f (P-state) transition points. It is assumed that there are three voltage transition points.

Increasing intermediate P-states helps in getting better granularity and power savings, however, there is a transition time interval when the processor is unavailable for any operation [53]. Thus P-state transition incurs a power penalty. Additionally, as seen in Fig. 2.12, with increasing number of available P-states, improvement in power savings is marginal.

Based on the available number of P-states in the processors in our experiments, we estimate the power savings by calculating how much time the processor stays at a given state under the operating systems ‘on-demand’ governor compared to that of our designed neuromorphic controller.

## 2.4 Results and Discussion

To evaluate the performance of the neuromorphic controller, the workload values obtained from the microprocessor are fed as input to the simulated neuromorphic controller. The predicted workload is used to calculate the next supply voltage ( $V_{DD}$ ) and frequency ( $f$ ) or P-state for the microprocessor core. Hence, instead of using the conventional approach of scaling the voltage after the workload value is known, workload prediction allows preemptive

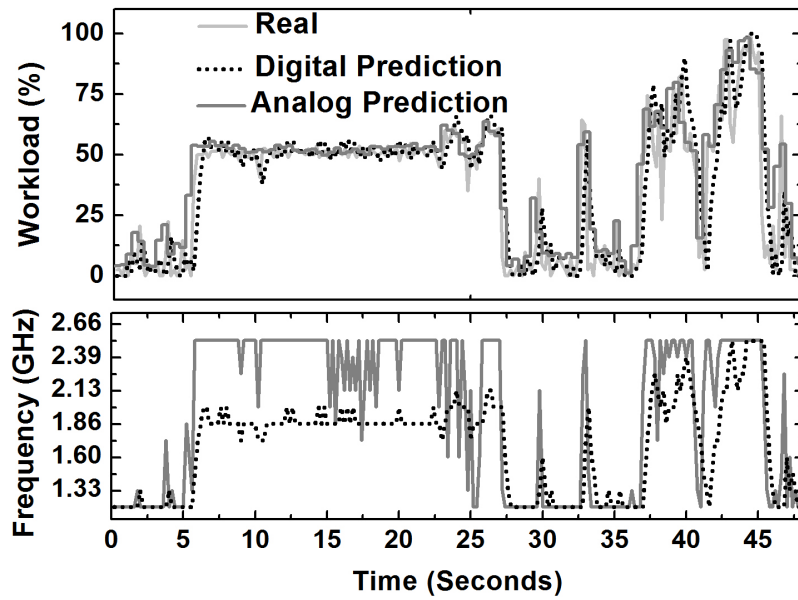


Figure 2.13: Workload prediction and associated P-state frequencies for Intel i5 processor. The processor has more P-states allowing finer granularity in DVFS. The controller predicts workload with high accuracy and results in lower P-states giving power savings.

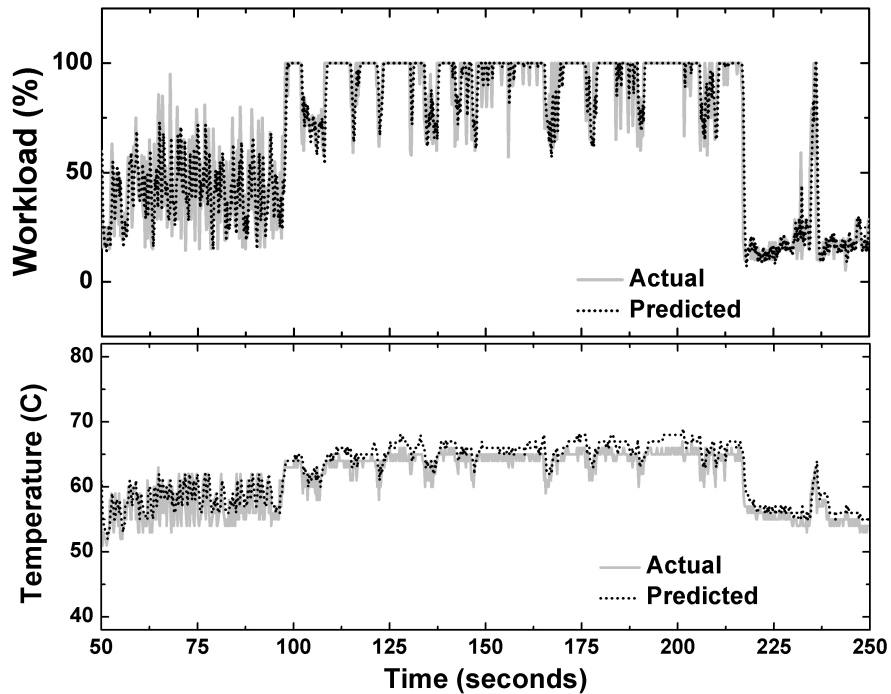


Figure 2.14: Workload and temperature prediction for the processor using the neuromorphic controller. The history based predictor works with an average RMS error of about  $1^{\circ}\text{C}$ .

DVFS. For testing the performance and prediction efficiency of our designed neuromorphic circuit, the predicted workload values are mapped to 3 P-states for the T2400 processor and 11 P-states for the Core i5 processor. When the predicted workload value crosses a threshold, the voltage and frequency are appropriately scaled. Fig. 2.13 shows processor workload from the Core i5 processor running OpenSSL encryption program, predicted values from the controller and the mapped P-states. The predicted mapped P-states are lower than the actual values, thus giving power savings as compared to the OS based frequency scaling governor. Fig. 2.3 shows the workload and corresponding temperature values. The neuromorphic controller accurately predicts the workload and temperature values.

#### *Prediction Accuracy*

Table 2.3 summarizes and compares the simulation results of the analog and digital design. The analog circuit consumes higher power during operation, but its main advantage is non-linear asynchronous operation. This makes it ideal for a general purpose, multiple input and non-linear controller for microprocessors. The digital circuit occupies a larger area since a large number of gates are required to construct the 6-bit counters and comparators.

If the predicted supply voltage is lower than the required value, the error is considered to be a performance hit. However, if the predicted supply voltage is higher than the required value, it is considered as over-estimation. The bias value in the neuromorphic circuit  $\eta$  can be varied to reduce the performance hit and over-estimation. Increasing the bias value consistently predicts a higher workload value than the actual, resulting in lower performance hit at the cost of higher over-estimation and vice versa. This trend is shown in Fig. 2.15. Over-estimation of supply voltage results in higher power dissipation. However, the error in predicting P-states significantly reduces since each workload value with a range of 0 to 100 is mapped to the available number of P-states (3 or 11).

Both designs result in similar performance hit, however over-estimation is higher for the analog design. Worst case prediction error is 17% while the average prediction error is as low as 1.2%. The neuromorphic controller reduces the total error by 50% on an average and over-estimation reduces by about 32% giving power savings. The average error in temperature prediction is 1.6°C with worst case error up to 2.8°C.

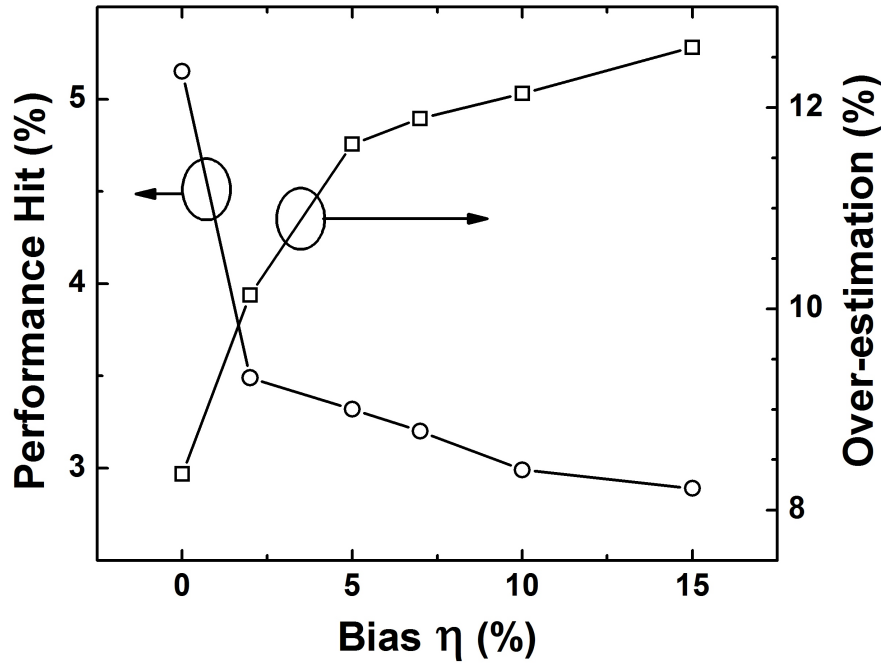


Figure 2.15: Effect of the bias value  $\eta$  on performance hit and over-estimation of the predictor.

Table 2.3: Summary of Results

Metric	Digital
Operation	Linear
Training	Error Subtraction
Frequency	1 MHz
Power	12.1 $\mu$ W
Area	360 $\mu$ m <sup>2</sup>
RMS Workload Prediction Error	1.2%
RMS Temperature Prediction Error	1.6 $^{\circ}$ C

In this work, the neuromorphic controller was verified with commercial microprocessor workloads with a resolution of 100-200ms. However, the controller is functional in the kHz-MHz range. The accuracy and efficiency of the controller depends on workload fluctuation. Based on previous work [2], the workload profile in nanosecond resolution has much lower time-domain fluctuations since only a single thread runs per core and there is much lower context switching. Hence, the controller prediction and efficiency should be equal to, if not better, at higher frequencies making it suitable for per-core DVFS using on-chip DC-DC converters. Further details regarding nano-scale DVFS and on-chip regulator requirements are discussed in detail in the next chapter.

### *Comparison with Hardware-based Schemes*

Conventional DVFS techniques are typically implemented in software making precise comparison between the neuromorphic controller and software based schemes non-trivial. However, software based DVFS schemes can be implemented in hardware. For first-order comparison, we have assumed a control-theory based DVFS scheme [30] implemented in hardware. The following equation is to be implemented:

$$\Delta w_i = K_p \varepsilon(t) + \frac{1}{I} \sum_{T_i} \varepsilon(t) + D \frac{\varepsilon(t) - \varepsilon(t - T_D)}{T_D} \quad (2.8)$$

$$w_{i+1} = w_i + \Delta w_i \quad (2.9)$$

where  $K_p$ ,  $I$  and  $D$  are the controller coefficients with values of 0.5, 28 and 0.00001. Using a similar methodology as our neuromorphic controller of bit-shifts for implementing division, implementing the above equation would still require a dedicated adder/multiplier unit making it significantly larger in area and complexity. The error of the PID controller based scheme in workload prediction is in the range of 10-11%, slightly higher than the nominal mode prediction error of the neuromorphic controller. Thus, it is clearly advantageous to use a neuromorphic controller implemented in hardware with power dissipation overhead of 10-150  $\mu\text{W}$  and greater prediction accuracy as compared to conventional OS-based and similar hardware-based schemes.

### *Power Savings*

Using the power model developed in Section 2.3, the average power savings as compared to the ‘ondemand’ scheme are calculated. Running the processor at a lower P-state results in additional delay due to lower frequency of operation. The associated delay for each sample is calculated by taking the ratio of operating frequency for the ‘ondemand’ and the neuromorphic scheme.

Most of the earlier work on dynamic power management use benchmark programs to quantify energy savings and the efficiency of their algorithms. To emulate realistic workloads



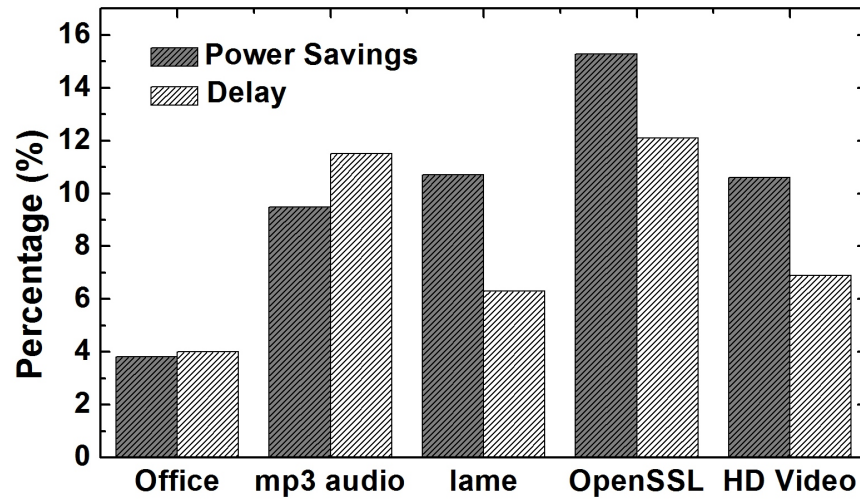


Figure 2.16: Power savings and associated delay for different applications on the Core i5 processor based system.

we run multiple programs simultaneously and take workload/frequency readings. Each reading is for a 10 minute interval, where the dominating process was either Office applications, playing mp3 audio and HD video playback. Two benchmark programs from OpenBench (lame mp3 encoding and OpenSSL encryption) were run as well. For each case, the user was browsing the web and switching between different windows to emulate realistic behavior.

Fig. 2.16 shows the relative power savings and associated delay for each case. In some cases the associated delay is lower compared to the power savings since the processor operates in the quadratic region of the power-frequency curve shown in Fig.2.12. For Office applications and mp3 audio playback, the processor workload is below 50% on average, hence the controller predicts higher P-states (lower  $V_{DD}$  and  $f$ ) and the processor operates in the linear region of the power-frequency curve. This results in higher delay as compared to power savings.

Fig. 2.17 shows the effect of changing the ratio of leakage power to dynamic power on power savings at different P-states. Since dynamic power is quadratically dependent on supply voltage while leakage power depends linearly on  $V_{DD}$ , higher power savings are possible when proportion of leakage power is low. However, our power model is pessimistic since we have not accounted for reduction in leakage power at low temperatures on running the processor at low supply voltage (higher P-state).

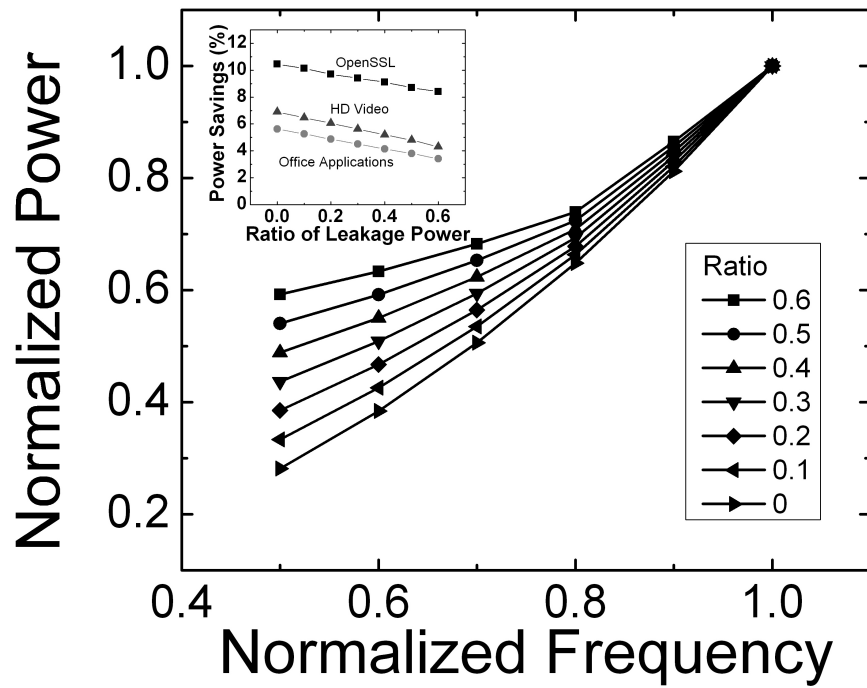


Figure 2.17: Effect of changing the ratio of leakage power to dynamic power on power savings at different P-states. The inset graph shows power savings for different programs at different ratios.

Table 2.4: Calculated Power Savings

Metric	Power Savings	Associated Delay
Total	52.7%	-
Best Case (relative)	15.2%	21.2%
Worst Case (relative)	1.5%	4%
DVFS during memory access	2.8%	0

Table 2.4 summarizes the power savings and associated delay using the neuromorphic controller. The power savings values are in comparison to a processor without any DVFS scheme while the best and worst case values are as compared to the ‘ondemand’ governor. Worst case power savings are observed when the processors are running at full capacity(100% workload), leaving minimal scope for voltage and frequency scaling. However, implementing DVFS when the processor is waiting to complete memory access as described in Section ??, an average of 2.8 % additional power savings is achievable. Fig.2.18 shows one such case for the bzip2 program. Hence, in the worst case, the neuromorphic controller performs about 3% better than an OS based scheme and about 7-9% better on average.

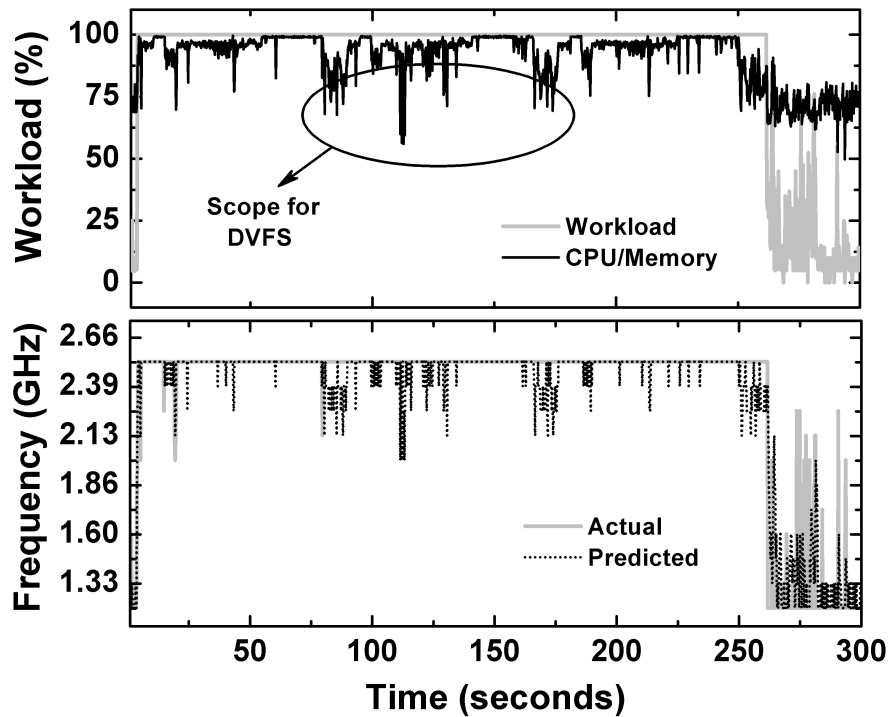


Figure 2.18: Plot showing percentage of CPU utilization/memory access when processor utilization is seen as 100% by the Operating System. The processor can be switched to a lower P-state (shown in lower plot) when ratio is low, giving additional power savings.

Conventional DVFS techniques are typically implemented in software making precise comparison between the neuromorphic controller and software based schemes non-trivial. Commercially available 45nm microprocessors typically have power ratings in the range of 10-45 W. Thus, it is clearly advantageous to use a neuromorphic controller implemented in hardware with power dissipation overhead of 10-150  $\mu$ W and greater prediction accuracy as compared to software based approach.

## 2.5 Fine-grained Power Management

Commercial microprocessors typically employ chip-wide dynamic voltage and frequency scaling. OS-controlled DVFS algorithms through the ACPI interface take up hundreds of CPU cycles to transition from one P-state to another. DVFS is implemented in microsecond or even millisecond resolution since software control has inherent delay. Additionally power supply is provided by off-chip dc-dc converter module with a transient response in the microsecond scale. One of the major hurdles preventing per-core DVFS is the difficulty of integrating high efficiency dc-dc converters on chip.

Recent work on multi-core systems [2] with per-core voltage and frequency scaling regulators [6] has shown that the P-state transition overhead can be in the range of 20ns allowing opportunistic DVFS resulting in additional power savings. Utilizing per-core fine-grained DVFS can result in power savings of up to 18% depending on the workload characteristics. The neuromorphic controller with workload prediction in the MHz range would be ideal for such a scheme since it would allow fast adaptation to changing workload which is one of the limitations of OS-based schemes.

One of the concerns of fine-grained DVFS using a neuromorphic controller is its ability to react to fast changing workloads. To understand this effect, processor activity factor for the benchmark program *ocean* was used from [2] as input to the digital neuromorphic controller. Fig. 2.19 shows the actual and predicted workloads. Two scenarios, where prediction takes place every 10 processor cycles and another every 50 cycles are simulated. As can be seen, the history based neuromorphic controller has good accuracy even when workload prediction is done every 10 cycles.

For better understanding, we plot the neuromorphic controller prediction error vs. update interval, i.e., number of cycles the workload values are updated in Fig. 2.20. Prediction error results in performance hits or wasted energy as described in Fig. 2.15. The update interval is changed from 100 cycles (coarse DVFS) to 1 cycle or fine-grained DVFS. As expected, RMS error reduces significantly when the workload values are updated at intervals of 10 cycles or less. The inset plot shows the effect of prediction error on chip-wide DVFS. For chip-wide DVFS, the average workload of the four workloads is fed as input to the neuromorphic controller. The RMS error is calculated between the predicted workload and actual workload of each core. Even with prediction update every cycle, the RMS error is in the range of 25-58%. This implies that chip-wide dynamic power management is an inefficient method to reduce power consumption.

## 2.6 Conclusion and Future Work

A neuromorphic power and thermal management for on-line prediction and training is implemented in 45nm CMOS technology using analog and digital design approaches. Various performance metrics for both designs were compared and it is found that a hardware based

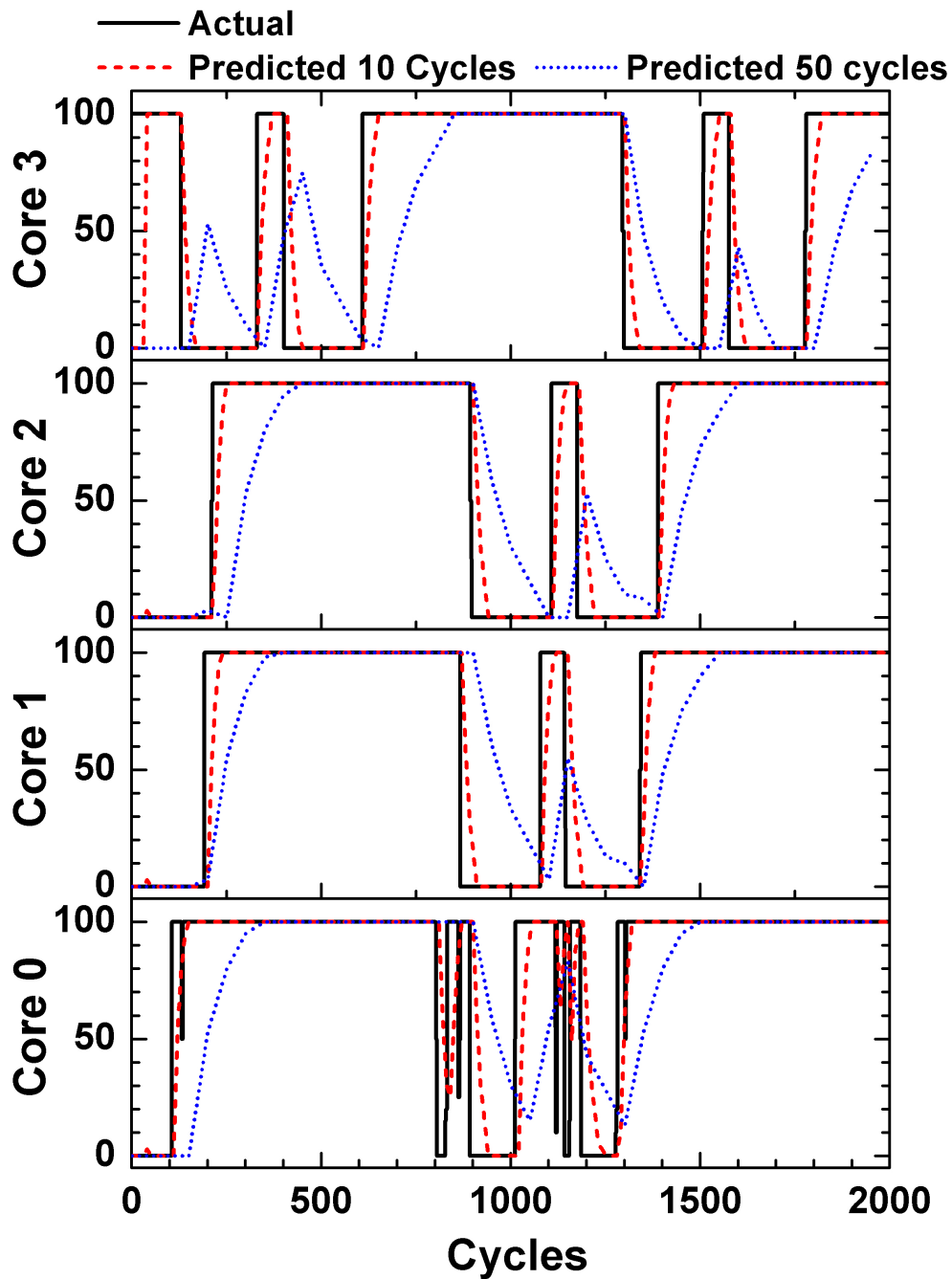


Figure 2.19: Actual vs. predicted workload from neuromorphic controller for the benchmark program *ocean* from [2]. The prediction is updated every 10 cycles and 50 cycles. The amount of error in prediction due to fine-grained operation is much lower and can result in power savings.

neuromorphic controller has better prediction accuracy than currently implemented DVS schemes in commercial microprocessors with minimal area and power overhead. Power savings is achieved by implementing DVFS based on the predicted workload and during

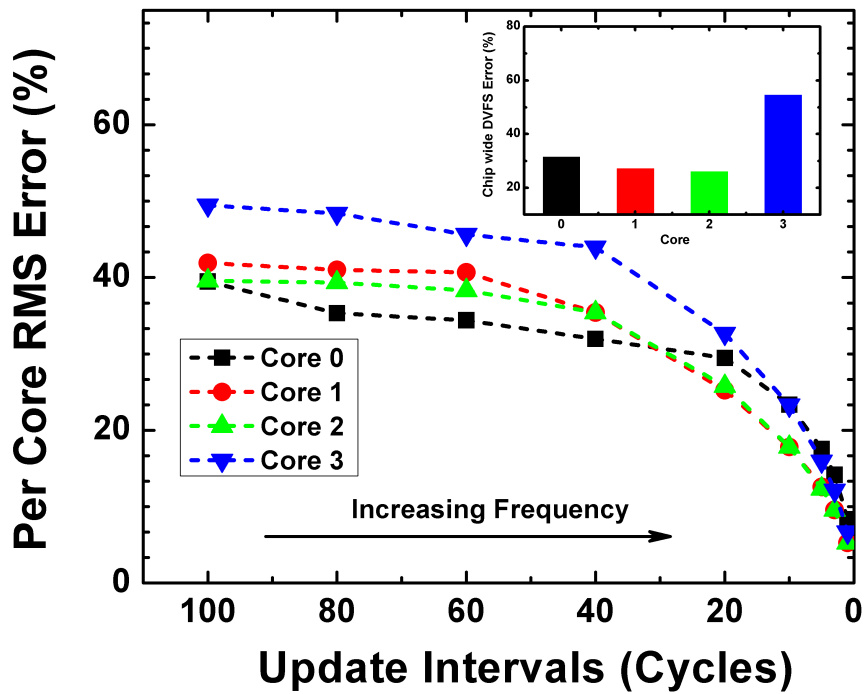


Figure 2.20: Prediction error vs. update times in cycles for the benchmark program *ocean* from [2]. Fine-grained DVFS where workload values are updated at higher frequencies result in low prediction errors. The inset plot shows prediction error when chip-wide DVFS is implemented.

memory stall cycles. Temperature prediction is useful to proactively implement DVFS and prevent performance degradation due to clock gating.

A detailed study of the efficiency of a neuromorphic power controller for fine-grained per-core dynamic power management is presented. It is found that the RMS prediction error is low when workload prediction takes place at higher frequencies (less processor cycles). For actual implementation, it is important to realize high-efficiency on-chip DC-DC converters that can react to such fast-changing workload values. It is also crucial to implement dedicated DC-DC converters per core on the chip. In the next chapter, the requirements for per-core fast transitioning voltage regulators are presented.

### DYNAMIC POWER MANAGEMENT AND REQUIREMENTS FOR ON-CHIP VOLTAGE REGULATORS

Dynamic power management aims to reduce power consumption by adjusting operating voltage and clock frequency adaptively according to changing workloads. Current microprocessors operate at clock frequencies in the range 600MHz to 3GHz depending on the target application such as servers, desktops, laptops or mobile devices. In the last four years, there has been an exponential growth in the use of smart-phones and tablet computers. These battery-operated devices have severe constraints on power dissipation. With improvement in process technology, complex System-on-a-Chip (SoC) with multi-core processors, higher operating frequencies and elaborate techniques for power management have been introduced in the market. Figure 3.1 shows the trend of SoC clock frequency with time. Dual core SoCs operating at around 1GHz were introduced in 2010 and quad-core SoCs with clock frequency of about 2GHz are predicted for 2012/13.

Clock frequency in the GHz range result in workload transients that are on the order of hundreds of nanoseconds. However, the power delivery network typically consists of off-chip voltage regulators operating at frequencies under 5MHz. Additionally the power delivery network consists of finite parasitic resistance and inductance. Hence, these regulators have a transient response time in the range of microseconds or milliseconds and are unable to respond to fast changing workloads and DVFS.

The power delivery network using off-chip regulator also places stringent requirements for decoupling capacitors [54]. A power distribution network consists of the power supply (regulator), a current load (processor) and the interconnect lines connecting the supply to the load. The power supply is modeled as an ideal voltage source providing the power and ground voltage levels ( $V_{DD}$  and  $V_{GND}$ ). The current load represents the actual processor (transistors/gates) and is modeled as a variable current source. The power and ground lines have finite parasitic resistance and inductance. Hence there are resistive voltage drops  $\Delta V=IR$  and inductive voltage drops  $\Delta V=L di/dt$ . The voltage levels at the load terminals deviate from the nominal power supply levels dropping to  $V_{DD}-IR-L di/dt$ . These fluctuations

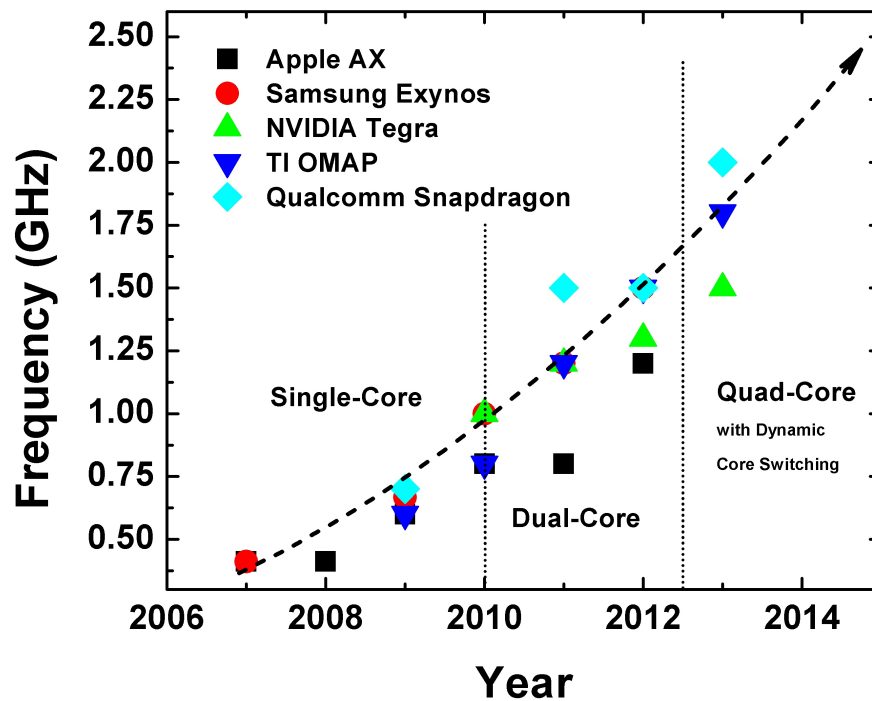


Figure 3.1: Trend of SoC clock frequency and number of cores with time from major manufacturers.

in the supply voltage levels are referred to as power supply noise. To reduce power supply noise, instantaneous charge is supplied to the load through on-chip and off-chip decoupling capacitors. The power delivery network and decoupling capacitance placement has to be carefully designed for the processor to operate without failure.

Based on the above discussion, the main issues with off-chip regulators are as follows:

1. Off-chip regulators are bulky and occupy a significant area of the PCB.
2. Parasitic resistance and inductance between the regulator and processor prevents the regulator from reacting to fast changing workloads.
3. Parasitic resistance and inductance on the power delivery network poses stringent requirements on off-chip and on-chip decoupling capacitance.
4. With the advent of multi-core processor systems, having a common power plane across cores prevents opportunistic DVFS and results in wasted energy as shown in Fig. 2.20 from Chapter 2.



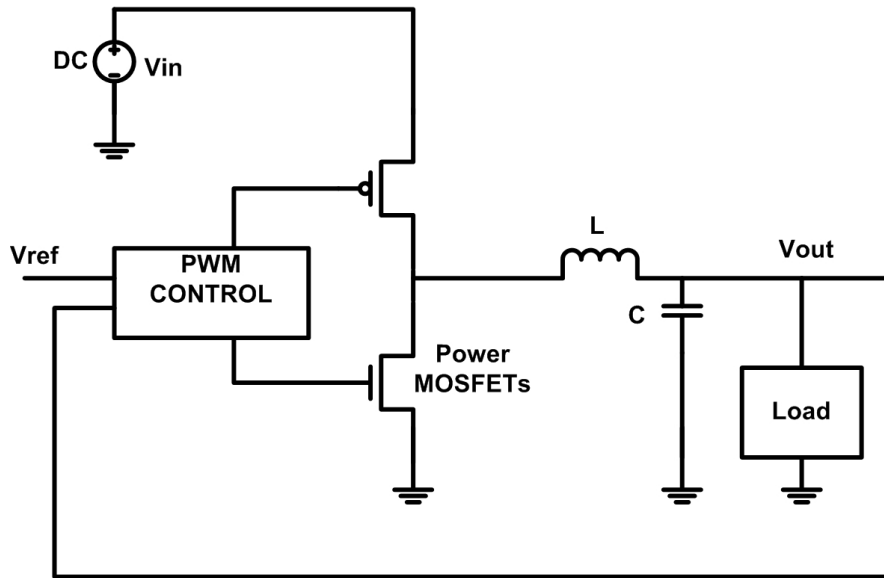


Figure 3.2: Simple schematic of a switch-mode regulator.

In order to mitigate these issues, design and implementation of monolithic voltage-regulator adjacent to load has been an area of active research [13, 55, 56, 57, 6]. The main hurdle preventing the implementation of on-chip regulators is the size of the passive capacitor and inductor of a switch mode dc-dc converter. In the next section we discuss the basics of voltage regulator design.

### 3.1 Switch-mode converter basics

A switch mode converter consists of power switches between the input and output and an output low-pass LC filter as shown in Figure 3.2. The gate of the power switches on and off with a duty cycle  $D$  and it can be shown that in an average sense, the output DC voltage is equal to

$$V_{OUT} = D \times V_{IN} \quad (3.1)$$

The voltage at the output of the switches consists of high frequency harmonics due to switching in addition to this DC average value. The low-pass LC output filter removes the unwanted high frequency harmonics and giving a DC output voltage with a magnitude set by the duty cycle of the switching signal. The values of the passive  $L$  and  $C$  is determined by the allowed current ripple across the inductor and voltage ripple on the capacitor. We concentrate

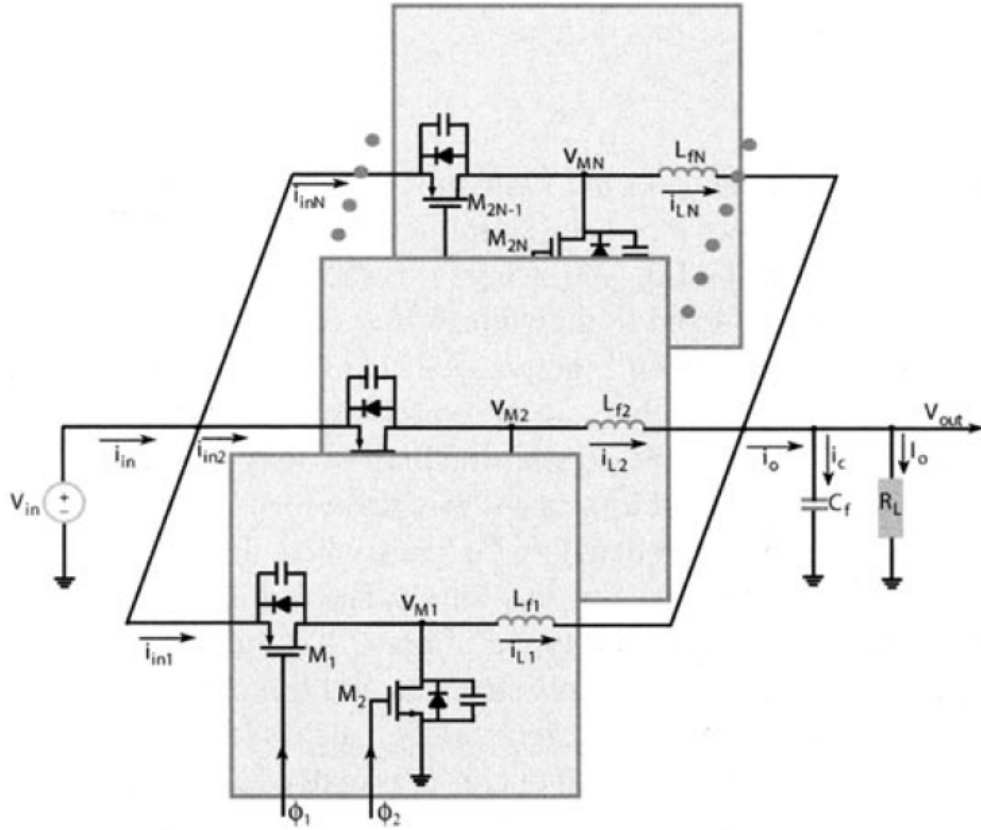


Figure 3.3: Interleaved DC-DC converter architecture from [3].

on the inductor requirement in this work. By utilizing an interleaved topology where multiple switches and L combinations are used in parallel with the triggering pulses applied with  $2\pi/n$  phase difference (where  $n$  is the number of stages), the output current ripple can be cancelled relaxing the inductor and output capacitor requirement as shown in Figure 3.3 [3] and 3.4 [4]. Multi-phase designs are best-suited for on-chip implementation since they allow smaller capacitance, have better transient response and the current density of interconnects is limited. Increasing the number of interleaved stages increases the number of power switches resulting in additional switching losses.

Based on the volt-second balance theory, the inductor value can be chosen using the following equation:

$$L = \frac{(V_{in} - V_{out})D}{2f_s\Delta I_o} \quad (3.2)$$

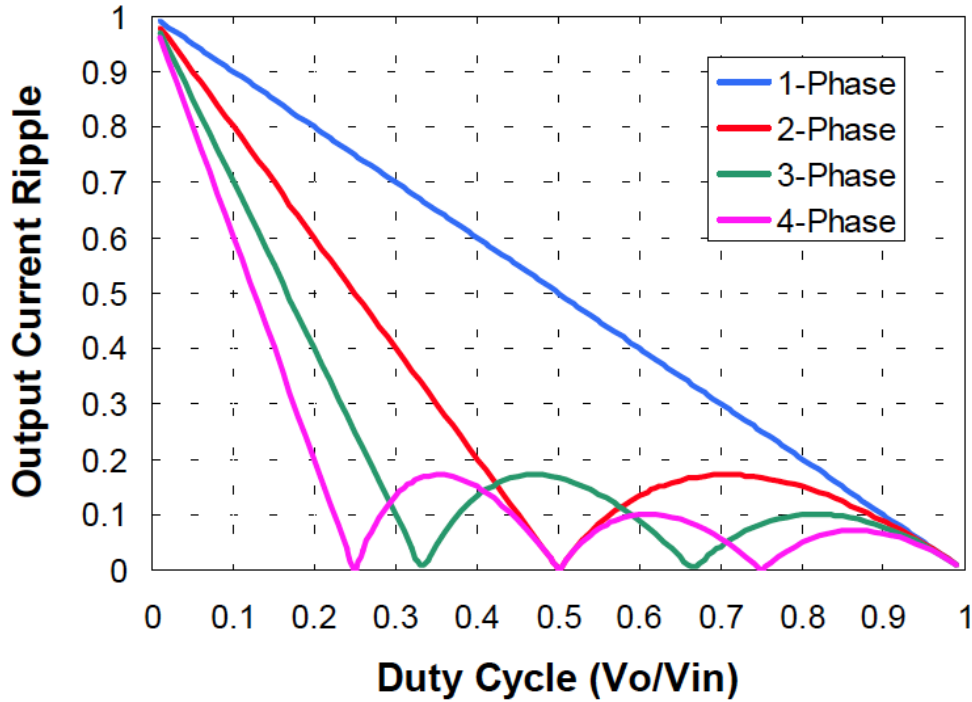


Figure 3.4: Current ripple cancellation using interleaved topology. Figure from [4].

where  $V_{in}$  is the input voltage,  $V_{out}$  is the output voltage,  $D$  is the switching duty cycle,  $f_s$  is the switching frequency and  $\Delta I_o$  is the maximum allowable current ripple.  $\Delta I_o$  is typically kept at 10-20% of the maximum output/load current specification. Using Equation 3.2 and assuming nominal values of  $V_{out}=1.2$ ,  $V_{in}=2.4$  and  $I_o=0.5A$ , the inductor value required vs. switching frequency is plotted in Fig. 3.5. For a 4-phase topology and high-frequency operation, inductor values are in the sub 10nH range, suitable for on-chip implementation especially using integrated magnetic materials. However, switching power loss increases at higher frequencies, reducing the regulator efficiency as shown in Fig. 3.5.

On-chip spiral inductors have relatively higher DC resistance compared to off-chip solenoidal inductors. The resistive loss from the inductor is the major contributor to the total losses in an integrated dc-dc converter [3]. Thick and wide metal wires in the order of  $10\mu m \times 20\mu m$  are required to keep resistive losses to a minimum. Hence it is very important to quantify the effect of resistive loss on regulator efficiency. We develop a very simple efficiency model of the regulator that takes into account losses in the switches and losses due to the series DC resistance of the on-chip inductor.

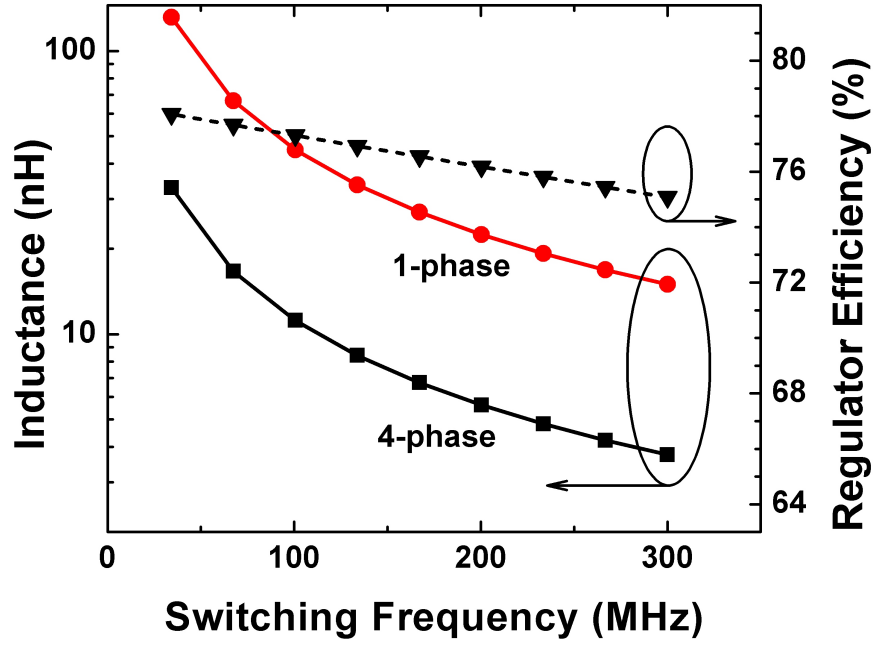


Figure 3.5: Value of filter inductor for a voltage regulator as a function of frequency. As can be clearly seen, to use integrated on-chip inductors, high switching frequency operation is required. However, regulator efficiency reduces at higher frequencies due to switching losses.

#### Efficiency Model

The efficiency of a dc-dc converter is given by

$$\eta = \frac{P_{out}}{P_{out} + P_{losses}} \times 100\% \quad (3.3)$$

where  $P_{out} = V_{out} \times I_{out}$ . In our simplified first-order model, we take into account power loss in the inductor due to low Q-factor and switching losses.

$$P_{inductor} = (I_{outmax} + \Delta I_{out} \times \sqrt{2})^2 \times R_L \quad (3.4)$$

where  $R_L$  is the inductor series resistance. Losses due to the power MOSFET switching is given by

$$P_{switching} = CV_{in}^2 f_s \quad (3.5)$$

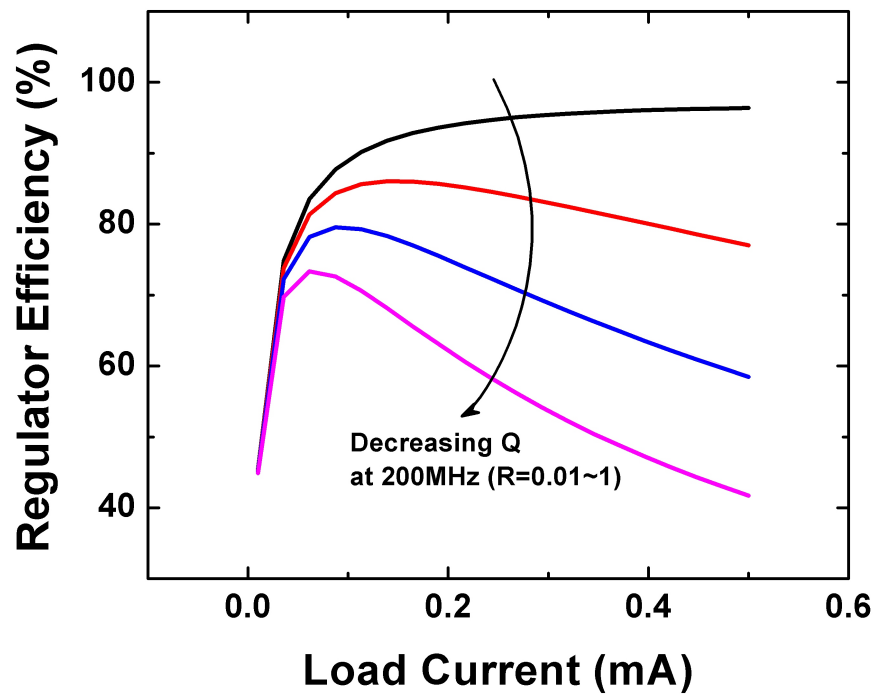


Figure 3.6: Efficiency of voltage regulator for different Q inductors.

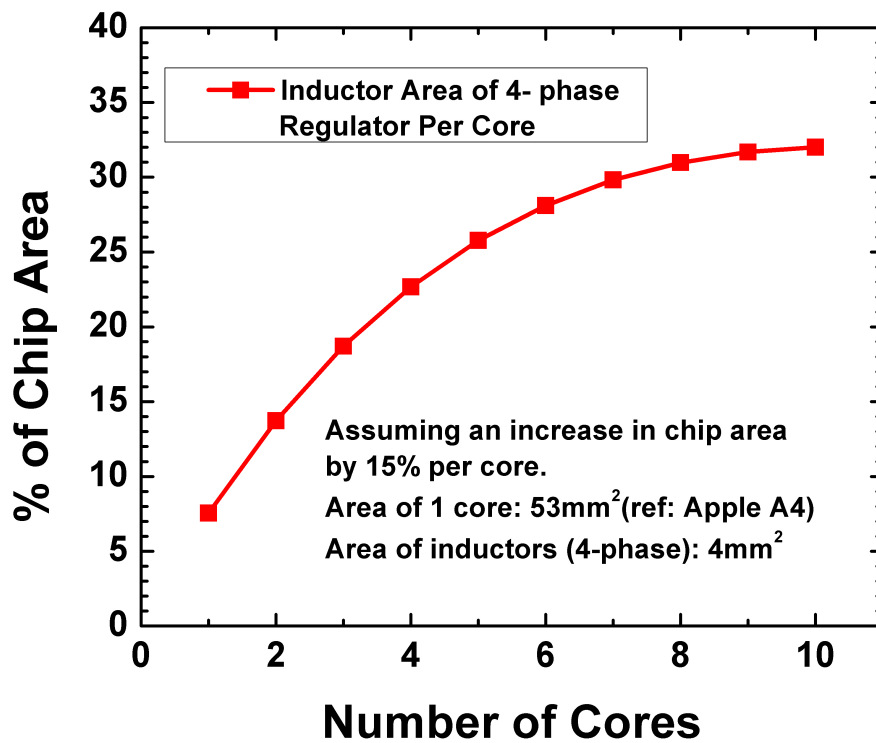


Figure 3.7: Plot of projection of area occupied by on-chip inductors in a 4-phase dc-dc converter with the number of cores. Area of a single core is assumed to be 53 mm<sup>2</sup> from [5].

Plotting the efficiency of the converter using the above equations for different DC resistance of the inductor shown in Fig. 3.6, we find that regulator efficiency is severely affected by high series resistance. It is important fabricate high-Q integrated inductors with minimal area for them to be useful in on-chip regulators.

Additionally, for multi-core systems, per-core dynamic voltage and frequency scaling is important to reduce wasted energy. As shown in Fig. 2.20 from Chapter 2, the prediction error of the neuromorphic controller reduces when workload prediction is updated at high frequencies and per-core DVFS is implemented. Multiple on-chip dc-dc converters are required to implement DVS per core of the chip. Since on-chip inductors occupy significant area, integrating multiple inductors can be prohibitively large. Fig. 3.7 plots the projection of percentage of area occupied by on-chip spiral inductors in a 4-phase buck converter vs. the number of cores. Each inductor area is assumed to be  $500 \times 500 \mu\text{m}^2$  from [6] and the area of a single core is assumed to be  $53 \text{ mm}^2$  (Apple A4 SoC area) from [5]. For scaling the area of the entire chip with the number of cores an increase in 10% is assumed per core. From these first-order assumption, it can be seen that merely the inductors can take up as much as 30% of the chip area for 8 cores or more. Hence, to realize multiple on-chip dc-dc converters it is imperative to reduce the area of the inductors.

The next section discusses previous work to develop on-chip voltage regulators with integrated inductors. Based on results from previous publications, specifications are developed for the design of inductors with magnetic materials.

### 3.2 Review of On-Chip Voltage Regulators

In order to reduce the requirements of the passive LC filter and allow on-chip integration, recent work has concentrated on high frequency buck converters with fast transient response. Voltage regulators using integrated inductors with magnetic materials have also been proposed [11, 12], however the inductors still occupy large area with values in the range of about  $1 \mu\text{H}$  operating at a few MHz. Interleaved topologies with inductor values below  $10\text{nH}$  operating above  $200\text{MHz}$  have been proposed in [55, 8]. The feasibility of on-chip inductors with integrated magnetic materials in voltage regulators was explored in detail in [13] and upto 6.8X reduction in off-chip decoupling capacitance was reported. Very high frequency

Table 3.1: Related Work

Publication	Operating Frequency MHz	Core Material	Number of phases	Inductance per phase nH
[11]	2	Ti/FeTaN	1	1000
[12]	5	FeBN	1	1000
[59]	25	air-core	4	220
[60]	80	air-core	4	26
[9]	45	air-core	2	11
[8]	233	air-core	4	6.8
[57]	100	air-core	8	1.8
[61]	170	air-core	4	2
[62]	300	air-core	4	2
[6]	250	air-core	4	1
[58]	3000	air-core	4	0.32

switching regulator at 3GHz was reported in [58] using pH inductors but with efficiency as low as 48% making them as good as linear regulators. Recently a 3-level dc-dc converter using an interleaved a mixed-topology consisting of input flying capacitor and output LC filter was proposed to further relax the inductor requirements [6]. Successful implementation using air-core inductors with Q-factor of 4 at 220MHz was demonstrated in this work.

Table 3.1 summarizes the inductor requirements of previously published works on integrated on-chip buck converters. For reference, the first two rows show results from inductors with magnetic cores. Most designs using on-chip air-core inductors operate at several MHz and utilize interleaved topology. The quality factor of the inductors is in the range of 4 to 20 at their operating frequency.

An active area of research to reduce the area of on-chip inductors is the use of integrated magnetic materials. Scaled magnetic inductors can allow multiple dc-dc converters dedicated to each core in multi-core processors. Details of design, simulation, fabrication and measurement of CMOS process compatible planar inductors with integrated magnetic materials and their potential applications are presented in the next chapter. Based on the system-level study and exploration of previous works conducted in this section, the requirements for on-chip planar inductors with integrated magnetic materials as part of the output LC filter of buck converters can be summarized as follows:

1. Inductance in the range of 1nH to 20nH.
2. High quality factor or series resistance less than  $0.5 \Omega$  to keep resistive losses to a minimum. The operating or switching frequency should be above 100 MHz to allow voltage to transition with fast changing workloads.
3. Area significantly less than air-core counterparts (at least 2X-4X lower) so that multiple dc-dc converters can be realized on chip for per-core dynamic voltage scaling.



### DESIGN OF ON-CHIP SPIRAL INDUCTORS WITH INTEGRATED MAGNETIC MATERIALS

Incorporating magnetic material with on-chip inductor is one of the most researched approaches to increase inductance (L), quality factor (Q) and the silicon chip-area efficiency. A number of coil and winding geometries such as stripline, stripe, toroidal, solenoid and meander structures with different configurations of magnetic materials have been investigated [63, 64, 65, 66, 67, 68, 69]. Toroidal and solenoid structures are difficult to integrate on-chip and the processing steps are not CMOS compatible. Air-core inductors using planar spiral geometry is preferred due to ease in processing. However, inductance values in the range of 1-10nH with quality factors lower than 10 and inductance density lower than 100nH/mm<sup>2</sup>. This translates to occupied chip area of a few hundred micrometers by a few hundred micrometers.

Extending the design of discrete inductors on PCBs using solenoid structures with high permeability magnetic core to amplify the linked flux, magnetic materials can be added around planar inductors along the flux-path. Adding magnetic materials results in increase in inductance and quality factor along with reduced capacitance, resistance and area. Spiral inductors with single or double layers of magnetic material show very limited enhancement of inductance (up to 2X) [64, 65, 66]. The highest L density has been reported by using flux-enclosed structure of Co-Zr-Ta lamination around spiral inductors [63].

The primary goal of this work is to design on-chip spiral inductors with magnetic materials for use in integrated voltage regulators based on the specifications discussed in Chapter 3. In this chapter we begin with an analytical study of stripline inductors with magnetic material to get a conceptual understanding of the affecting parameters. Next, measurement results from different fabricated structures with NiFe and CoZrTa magnetic materials are presented. A short discussion about the potential of inductors with integrated magnetic materials for resonant clock distribution for microprocessors is presented. Finally, we explore the feasibility of these inductors for on-chip voltage regulators.

#### 4.1 Simulation and Analytical Study

It is essential to gain a qualitative and quantitative understanding about the effect of each constitutive parameter as well as the structural variations on the performance of on-chip spiral inductors with magnetic material prior to actual fabrication. Using finite-element 3D electromagnetic field simulators such as Ansys HFSS is a straightforward method which allows us to modify material parameters such as permeability and conductivity and create novel structures such as magnetic dots, filaments and films around the inductor trace metals. However, 3D electromagnetic solvers are computationally intensive and require a long time to converge to a solution. Hence, developing analytical equations provide useful insight about the effect of each constitutive parameter and allows us to design simulations efficiently.

##### *Analytical Model of Stripline Inductors with Magnetic Materials*

Previous work on deriving analytical equations for inductor structures with magnetic materials assume a single stripline inductor with a single magnetic film without any insulating layer covering the inductor metal from all sides.

As shown in the figure, region ‘0’ represents the inductor metal trace and region ‘1’ represents the magnetic film surrounding it. The figure denotes half of the stripe in the x-direction. The thickness of the metal wire is  $t_0$  and that of the magnetic film is  $t_1$ . The width of the metal wire is  $w$  and the length is  $l$ . According to [70], solving Maxwell’s equation gives the inductance of such a structure as

$$L = \frac{\mu_0 l}{2\alpha_0 \omega} \sqrt{\frac{\mu_1 \sigma_0}{\mu_0 \sigma_1}} \tanh(\alpha_1 t_1) \quad (4.1)$$

where  $\mu_i$  and  $\sigma_i$  are the permeabilities and conductivities of the respective mediums.

Additionally

$$\alpha_i = \frac{1+j}{\delta_i} \quad (4.2)$$

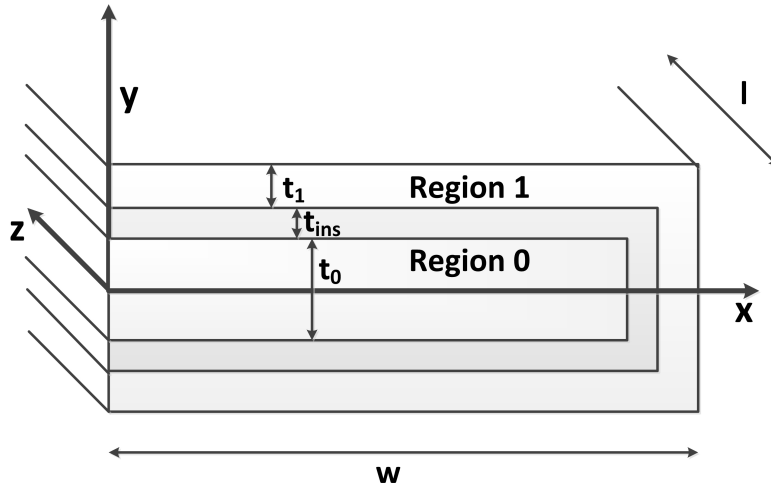


Figure 4.1: Structure of stripline inductor with magnetic material used to derive analytical equations. Region 0 is inductor metal, Region 1 is magnetic material and  $t_{ins}$  represents thickness of insulating layer.

$$\delta_i = \sqrt{\frac{2}{\omega\mu_i\sigma_i}} \quad (4.3)$$

Since the expression of L itself contains real and imaginary components, the imaginary component denotes inductance while the real part represents eddy current losses in the magnetic material. These equations are valid when the film thickness  $t_1 \ll t_0$ , since the effect of  $t_1$  in the y-direction is not considered. A similar expression has been derived in[71] by considering the reluctance of the magnetic structure. Deriving analytical expressions by this method is comparatively simpler since it does not involve solving complex differential equations from Maxwells equations and can be extended for more complex structures. Additionally, the reluctance method accurately models the effect of insulation layer between the inductor wire and the magnetic film.

#### *Comparison between Analytical Model and 3D EM solver*

On comparing the results of both approaches with that of 3D electromagnetic field solvers such as Ansoft HFSS, it is found that an expression that borrows features from both approaches accurately models the device behavior when magnetic film thickness is comparable to that of the inductor metal as well as on including insulation layer. Since eddy current loss is maximum in the magnetic vias[63], the high frequency response of the inductor

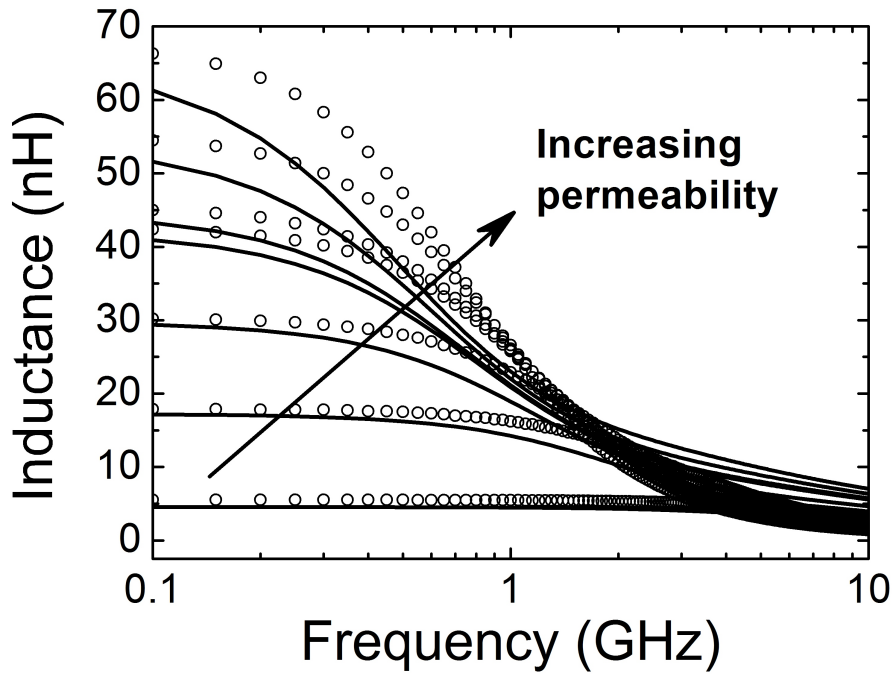


Figure 4.2: Comparison of analytical model with 3D EM simulations in HFSS for varying permeability ( $\mu$ ) of magnetic material

structure is accurately captured on replacing the thickness of the magnetic film  $t_1$  by the effective thickness  $\sqrt{2t_1}$ , which is the thickness at the corners of the film. The width of the magnetic film  $w$  is replaced by the total width of the magnetic film,  $2w$ . If there is an insulation layer, the total width is  $2(w + t_{ins})$ . Using these modified equations as given below, we find that the analytical expressions accurately model the results from HFSS simulations.

$$L = \frac{\mu_0 l}{4\alpha_0(w + t_{ins})} \sqrt{\frac{\mu_1 \sigma_0}{\mu_0 \sigma_1}} \tanh(\alpha_1 t_1) \quad (4.4)$$

Fig.4.1 shows the validation of the analytical model by comparing with 3D electromagnetic simulations in Ansys HFSS. The modified model is scalable to all significant materials and thickness of the insulating layer. Some important insights gained from the analytical model are:

1. Inductance increases on using high permeability ( $\mu$ ) magnetic material. However, since the skin depth  $\delta$  is also dependent on  $\mu$ , the frequency response suffers slightly on increasing  $\mu$  as seen Fig.4.2.

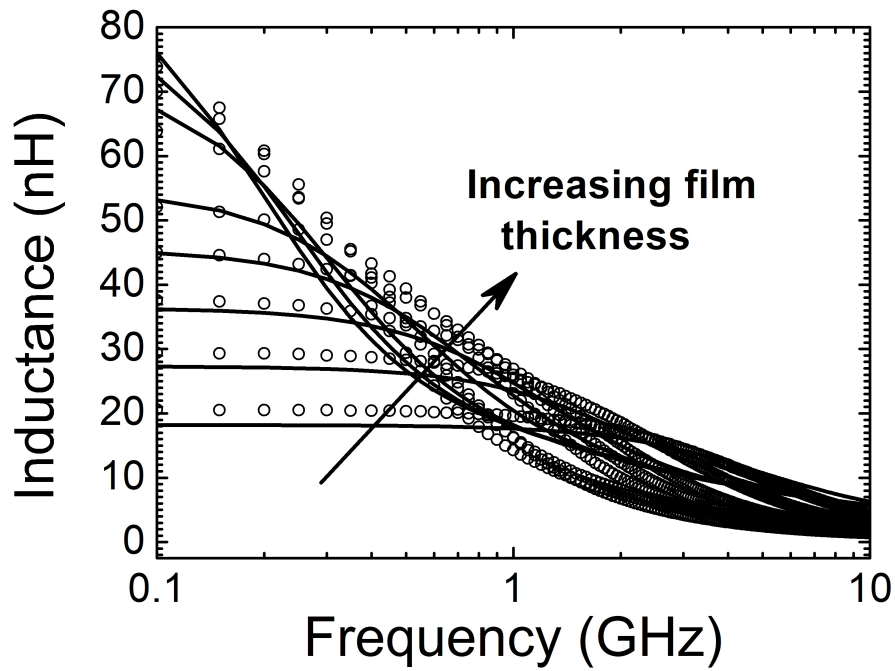


Figure 4.3: Comparison of analytical model with 3D EM simulations in HFSS for varying magnetic film thickness ( $t_1$ ).

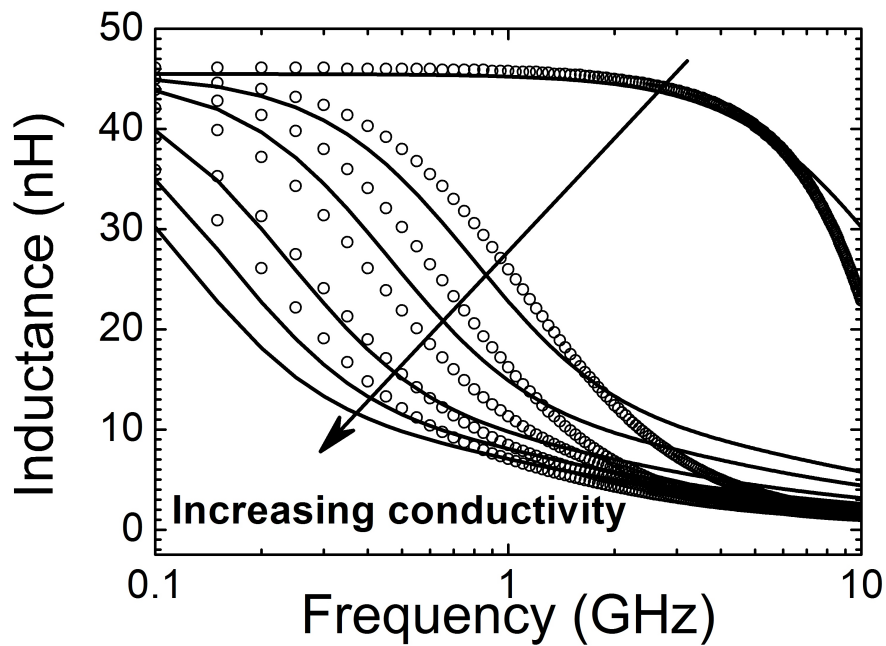


Figure 4.4: Comparison of analytical model with 3D EM simulations in HFSS for varying conductivity ( $\sigma$ ) of magnetic material.

2. There is high inductance enhancement on using thicker magnetic film since there is more high permeability magnetic material in the flux path. However, thicker magnetic

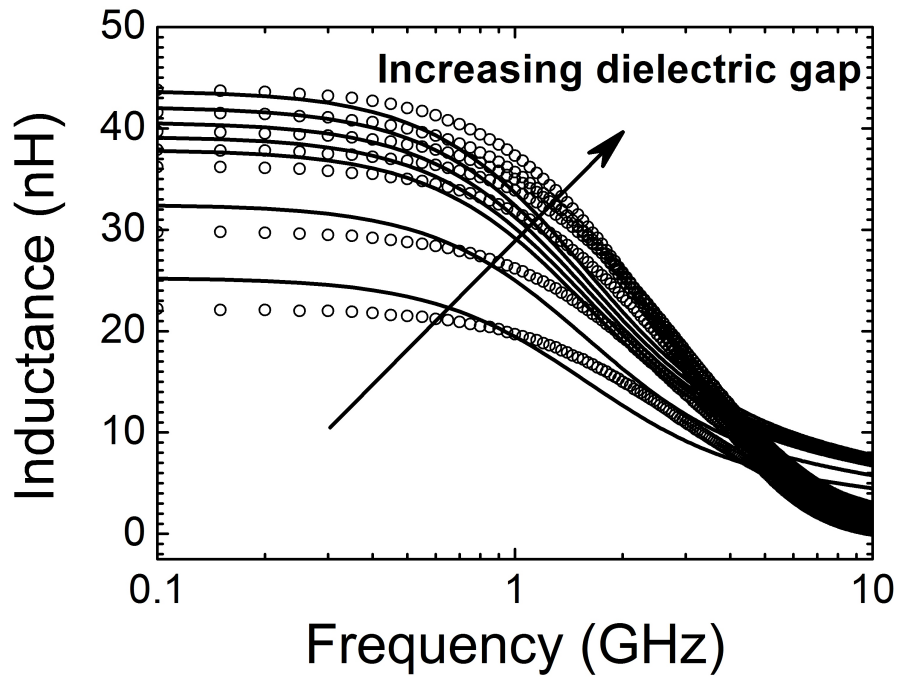


Figure 4.5: Comparison of analytical model with 3D EM simulations in HFSS for increasing dielectric gap ( $t_{ins}$ ) between inductor metal and magnetic material.

material results have more eddy currents resulting in a significant reduction in frequency response as seen in Fig.4.1.

3. Eddy currents scale with the conductivity of the magnetic material significantly affecting frequency response of inductance as seen in Fig.4.1. Hence, a material with high permeability and low conductivity would be ideal for on-chip inductors with integrated magnetic materials.
4. Inductance decreases on using an insulating layer between the inductor metal and magnetic film. This scales inversely with the insulation thickness since magnetic flux is strongest near the metal wires (Fig.4.1).

#### *Optimization of Inductor Structure*

Most previous approaches to integrate magnetic materials with on-chip inductors involve using a single or double magnetic thin film above and below the inductor metal layer. Additionally, to reduce eddy currents the magnetic material can be cut into stripes or dots. Simulations in

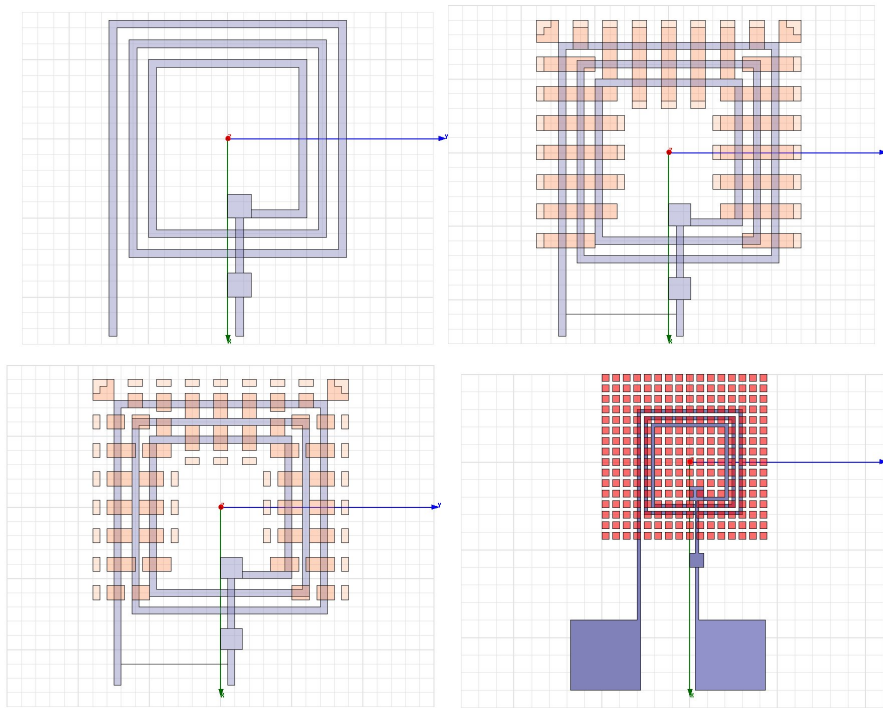


Figure 4.6: Simulation structures of bare inductor, inductor with continuous magnetic rings, broken stripes and magnetic dots in HFSS.

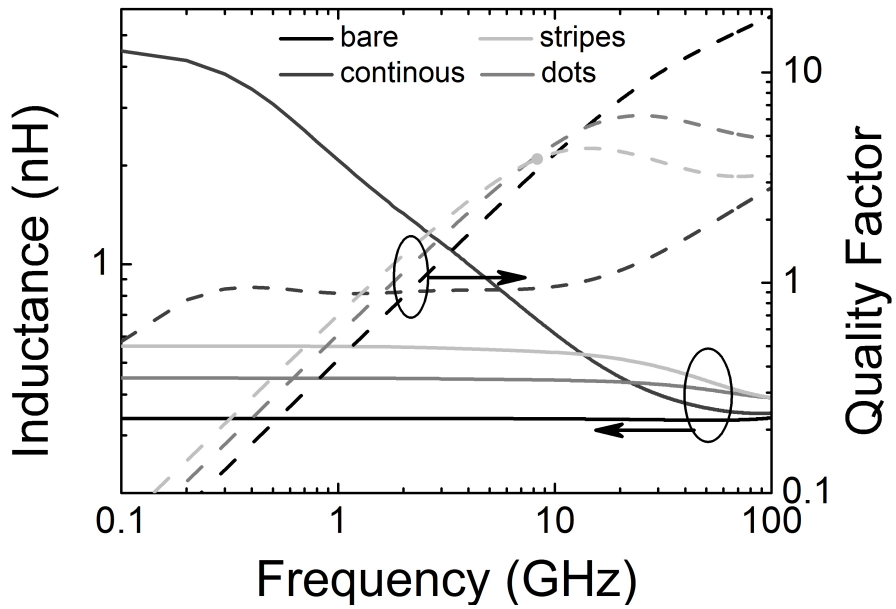


Figure 4.7: The complex permeability spectra of different patterned NiFe laminations.

HFSS show that any gap in the flux path significantly reduces the inductance value. However, since the eddy currents reduce as well, very high frequency response is observed.

Extensive simulations of different structures in HFSS show that the maximum inductance enhancement is obtained using a continuous magnetic material in the direction parallel to the flux path. Hence, the magnetic material should surround the inductor metal wires completely. Any gap in the form of dots or stripes reduces the inductance enhancement that is up to 2X the value of the bare inductor. However, a continuous structure also suffers from maximum eddy current loss and hence, low frequency response, typically in the Mhz range.

Hence, there is a clear tradeoff between high inductance value (gain) and frequency response (bandwidth) and achieving a good balance between the two based on the target application is one of the goals of this work. Based on the insights developed from the analytical models and simulation results from 3D electromagnetic field solvers, inductors with different structural modifications were fabricated and characterized. The fabrication and measurement results are presented in the next section.

#### 4.2 Device Fabrication and Measurement *Inductors with magnetic dots for high frequency applications*

In this work, through experiments and simulations, we have demonstrated that the high-frequency response can be effectively improved by using patterned magnetic dots, instead of a continuous magnetic thin-film. In addition, the frequency response is further enhanced by reducing the size of magnetic dots. For our experiments, we have fabricated single-turn spiral inductors with top and bottom magnetic layer. Using large and complex structures for on-chip inductor results in a low self-resonant frequency [72]. Since we have a single-turn and reasonably small sized inductor, self-resonance occurs at a very high frequency, allowing GHz range operation. The primary source of poor frequency response is identified as eddy current losses in the magnetic thin film, which is mainly a result of mutual coupling between the magnetic and inductor metal layer.

#### Fabrication

Inductors with magnetic material were fabricated on quartz substrate using photolithography and electron beam lithography<sup>1</sup>. Devices fabricated include inductors without any magnetic

---

<sup>1</sup>Fabrication of inductors were done by Wei Xu, Tawab Dastagir and Hao Wu.



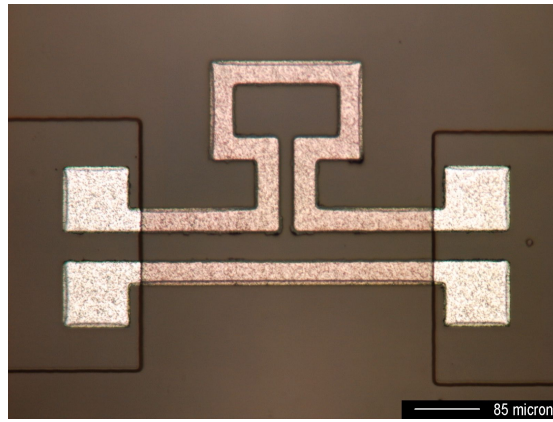


Figure 4.8: Bare Inductor

material, with a continuous magnetic film, square dots of sizes  $16 \times 16 \mu\text{m}^2$ ,  $8 \times 8 \mu\text{m}^2$  and  $4 \times 4 \mu\text{m}^2$ . The top view of different configurations and a cross-section view of the inductors with magnetic dots is shown in Fig.4.10. All dot patterns have the same area ratio of 64%, defined as the area of the magnetic dots to the complete film. HP8720D network analyzer and Cascade GS and SG probes were utilized for two-port measurements.

#### Results and Discussions

Measured and simulated inductors ( $L$ ) of different patterned devices are shown in Fig 4.16. Inductor without magnetic material (bare) shows a small but consistent inductance throughout the entire frequency range, while with double layer continuous NiFe film the inductance drops off at around 300 MHz. For the structure with the patterned magnetic dots, we observe a significant improvement in the high frequency response in spite of slightly lower values of inductance as compared to the continuous film case.

On maintaining the same area ratio of 64% and reducing the area of the dots, the frequency response improves from  $16 \times 16 \mu\text{m}^2$  case to the  $4 \times 4 \mu\text{m}^2$  case. The inductance of the  $4 \times 4 \mu\text{m}^2$  patterned dots inductor remains constant up to nearly 10 GHz, which is limited by the test equipment. Quality factors for these inductors shown in Fig 4.16 indicate that  $Q$  peaks close to 8.87 at 9 GHz for patterned inductor. Scaling down the size of the NiFe dots improves the quality factor. The device with  $4 \mu\text{m}$  dots achieves the highest  $Q$ . These observations are reproduced by 3D EM simulations in Ansoft HFSS by the broken lines in Fig 4.16.

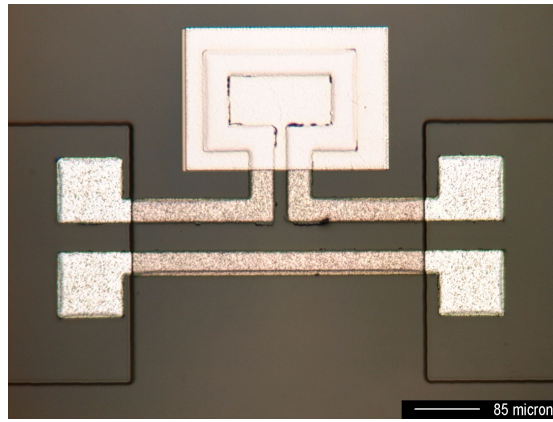


Figure 4.9: Inductor with magnetic film

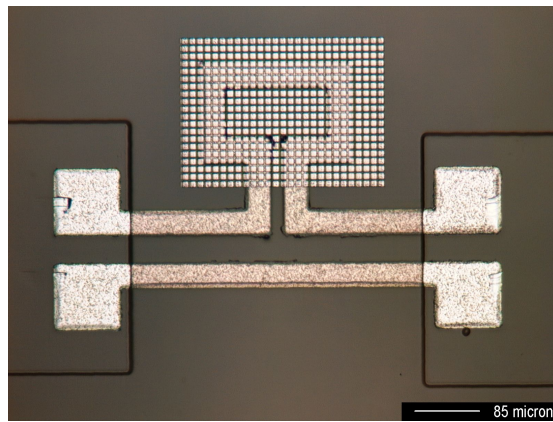


Figure 4.10: Inductor with patterned NiFe dots.

#### Equivalent circuit model and parameter extraction

The decrease of  $L$  at high frequency can be attributed to eddy current loss in the magnetic layer which mutually couples with the inductor metal layer, reducing  $L_s$  and quality factor.

Based on the equivalent circuit diagram given in Fig. 4.12 and assuming that all capacitances are negligible, for the primary coil we have

$$R_s I_1 + j\omega L_s I_1 - j\omega M I_2 = 0 \quad (4.5)$$

In Fig 4.12,  $L_m$  is the inductance of the metallic magnetic material. Thus, in the secondary coil we have

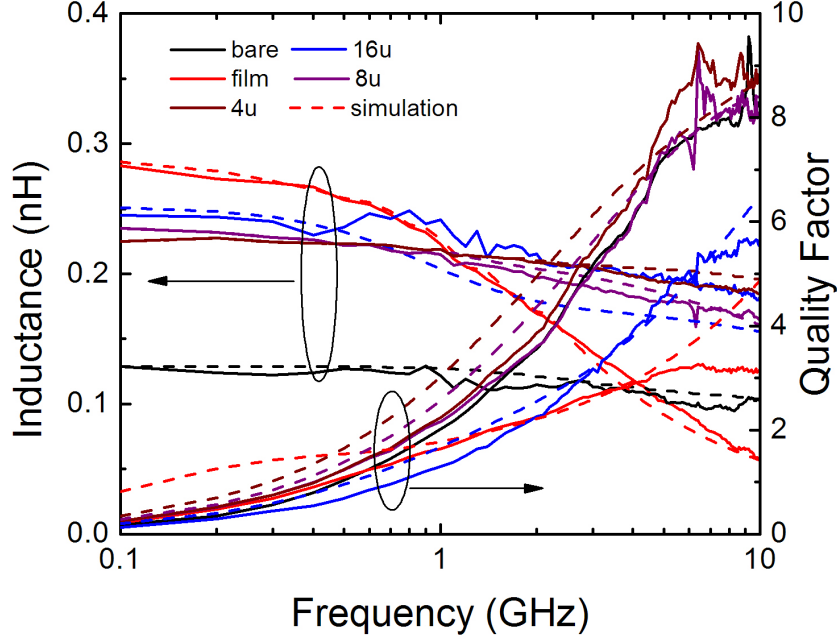


Figure 4.11: (a) Measured and simulated inductance ( $L$ ) vs. frequency for different patterns (bare, film and dots). (b) Measured and simulated quality factor ( $Q$ ) vs. frequency for different patterns. Broken lines show simulation results.

$$R_m I_2 + j\omega L_m I_2 - j\omega M I_1 = 0 \quad (4.6)$$

Since permalloy is highly conducting, we assume  $R_m \ll \omega L_m$  at high frequencies, which gives  $I_2 = (M/L_m)I_1$ . Using this in (4.5), we have

$$R_s I_1 + j\omega I_1 \left( L_s - \frac{M^2}{L_m} \right) = 0 \quad (4.7)$$

Thus, with increasing frequency the mutual inductance  $M$  increases due to stronger inductive coupling and reduces the measured inductance by a factor of  $M^2/L_m$ .

Breaking the magnetic layer into discrete dots effectively reduces the eddy current loops. Square magnetic dots of  $16 \times 16 \mu\text{m}^2$ ,  $8 \times 8 \mu\text{m}^2$  and  $4 \times 4 \mu\text{m}^2$  are fabricated maintaining the same volume ratio. The eddy current loops are now confined within the dots, reducing  $L_m$ , the mutual coupling  $M$  and the eddy current loss at high frequency.

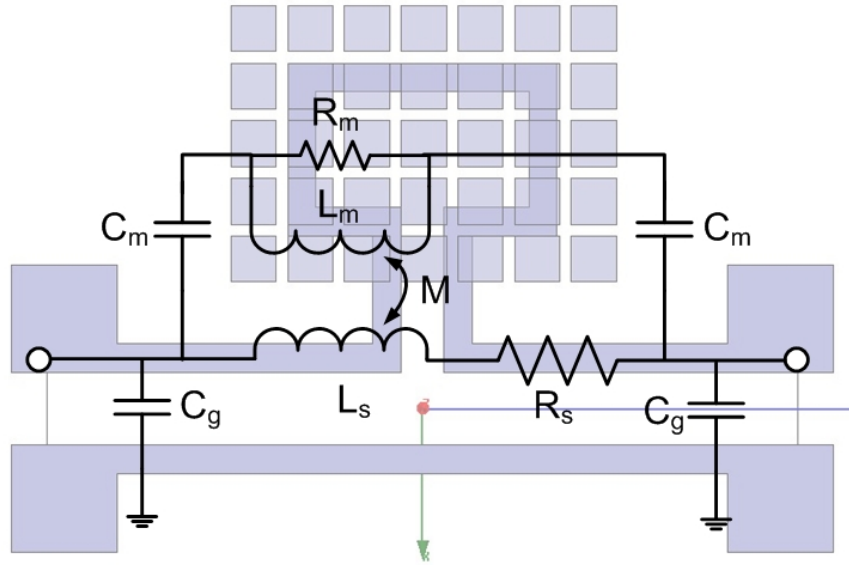


Figure 4.12: (a) Equivalent circuit model used to extract inductance. The eddy current loop inside the metallic magnetic layer forms a secondary inductor  $L_m$  that mutually couples with  $L_s$ . The mutual inductance increases with frequency, reducing the inductance seen through the primary inductor ports at high frequency.

The factor  $M^2/L_m$  is frequency dependent since eddy currents are enhanced at high frequencies. It can be extracted from measurement as follows:

1. From the equivalent circuit in Fig. 4.12,  $M$  is assumed to be negligible at low frequency (0.1 GHz) and the measured inductance is  $L_s$ .
2. At high frequencies the measured inductance is effectively  $(L_s - M^2/L_m)$ .
3. At the desired frequency, the difference of the measured inductance at low frequency (0.1 GHz) and at that frequency, gives  $M^2/L_m$ .

$M^2/L_m$  at 1 GHz and 10 GHz are given in the Table 4.1. As expected, mutual coupling is much higher in the case of a continuous magnetic film; it reduces as the size of the magnetic dots is decreased.

For further confirmation, current density is plotted inside the magnetic film and magnetic dot structure in the 3D EM simulator as shown in Fig. 4.13; there is a large eddy

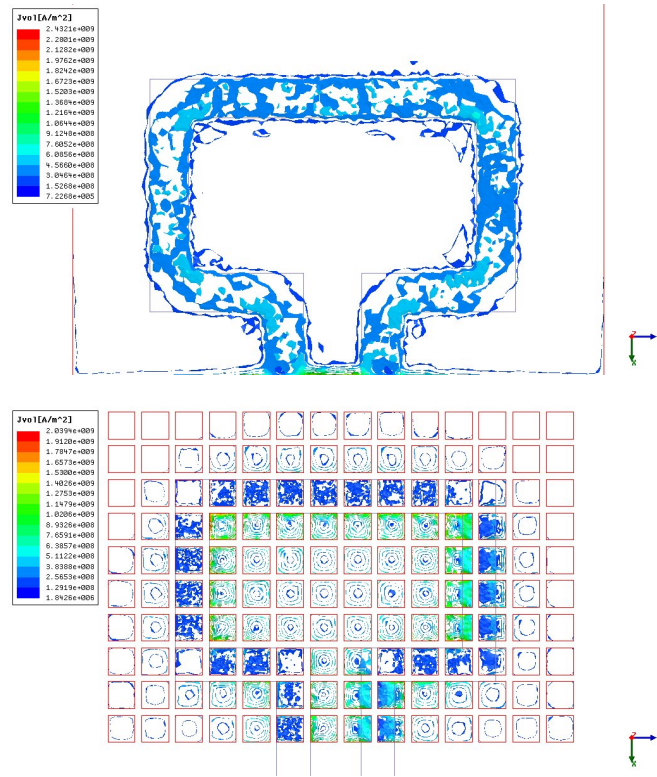


Figure 4.13: 3D EM simulator plot showing magnitude of current density in the magnetic film and magnetic layer with  $16 \mu\text{m}$  dots. Large eddy current loops are formed in the film compared to that of magnetic dots, resulting in poor frequency response.

current loop in the film structure (resulting in larger  $L_m$  and  $M$ ) while the eddy current loops in the magnetic dots are limited by the dot size.

#### *Inductors with Magnetic Rings*

Based on the simulation and analytical study described in Section, devices with magnetic material providing a continuous high permeability flux path were fabricated for high inductance enhancement. In this work, we have demonstrated spiral inductors with patterned

Table 4.1: Extracted values from measurement

Inductor Type	$L_s$	$M^2/L_m$	$M^2/L_m$
	0.1 GHz	1 GHz	10 GHz
Thin film	0.283 nH	0.06 nH	0.225 nH
$16 \times 16 \mu\text{m}^2$	0.269 nH	0.03 nH	0.082 nH
$8 \times 8 \mu\text{m}^2$	0.235 nH	0.02 nH	0.07 nH
$4 \times 4 \mu\text{m}^2$	0.2248 nH	0.005 nH	0.0385 nH

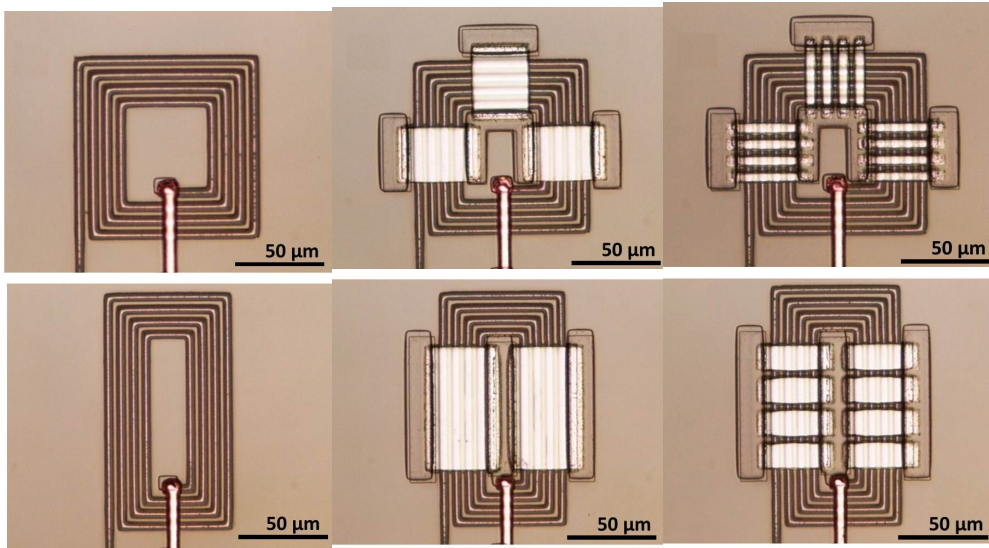


Figure 4.14: Optical images of the fabricated rectangular and square inductors (a) without NiFe (bare), (b) with NiFe thin film enclosed, (c) with patterned NiFe rings.

permalloy rings at 100  $\mu\text{m}$  scale achieving enhancements of 6X in L and 3X in Q-factor at frequencies as high as 200 MHz. An inductance density around 770 nH/mm<sup>2</sup> is achieved, which is higher than other reported values of NiFe based inductors [64, 68].

#### Fabrication

Inductors with laminated NiFe thin films were fabricated on quartz substrate using electron beam lithography (EBL) and magnetron sputtering for pattern definition and metallization<sup>2</sup>. Fig. 4.14 shows the optical images of fabricated inductors and Fig. 4.15 shows the fabrication process flow. In order to control the stripe domains formed inside the thick NiFe film which result in degraded permeability, NiFe laminations were used by depositing layers of NiFe thin film with 5 nm Cr spacer layer in between [73]. By varying the thicknesses of each NiFe layer while keeping the total film thickness the same (1  $\mu\text{m}$ ), various values of permeability can be obtained. In this work, three types of lamination have been integrated into inductors with corresponding permeability of 170 (film A), 300 (film B) and 500 (film C) Gauss/Oe. Different types of spiral inductors were fabricated with magnetic rings, helping shed light on the optimization strategy.

<sup>2</sup>Fabrication of inductors were done by Wei Xu, Tawab Dastagir and Hao Wu



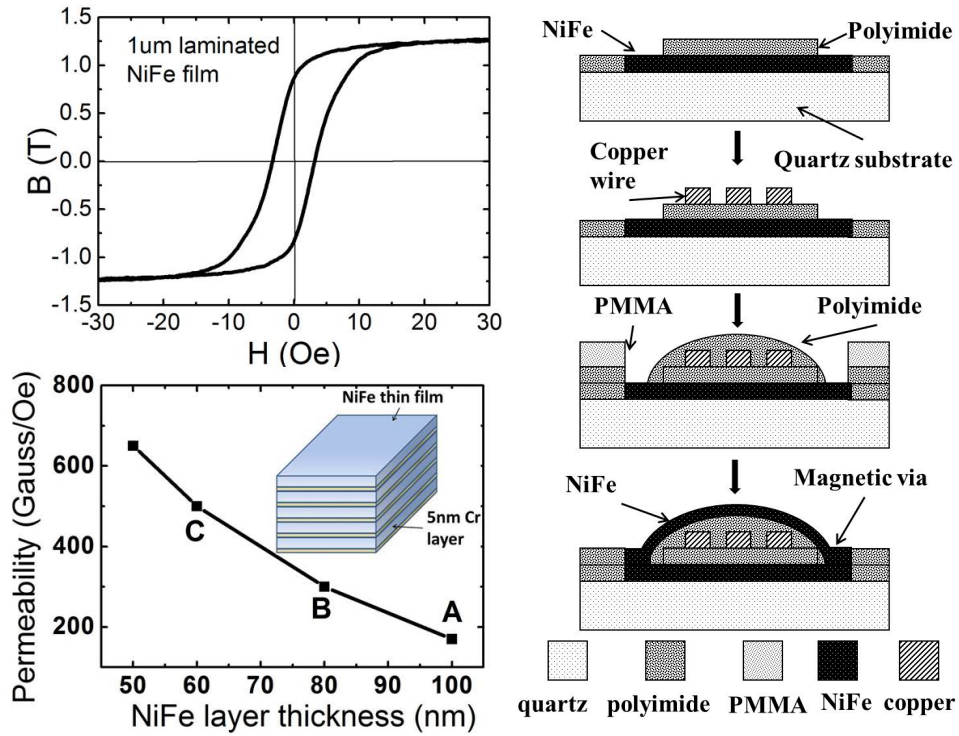


Figure 4.15: (a) B-H magnetic hysteresis loop of 1 μm laminated NiFe film (50nm NiFe/5nm Cr) and layer thickness dependence of permeability measurement. Inset shows schematic of laminated NiFe film. (b) Process flow for on-chip inductor with NiFe.

## Results and Discussions

The DC inductance of a rectangular line is approximately given by

$$L = \mu_0 \mu_r \frac{t_m l}{2w} \quad (4.8)$$

where  $\mu_r$  is the relative permeability,  $l$  is the length,  $w$  is the width and  $t_m$  is the thickness of magnetic film. A detailed theoretical analysis is given in [74].

### Permeability Dependent Inductance and Quality Factor

Fig. 4.16 shows L and Q plots vs. frequency for rectangular inductors with single NiFe ring structure using three types of laminated films. Among these, inductor with film C shows largest enhancement of 6X in L and 3X in Q-factor at 200 MHz. It clearly demonstrates the benefit of using high permeability magnetic material to enhance the performance of inductors

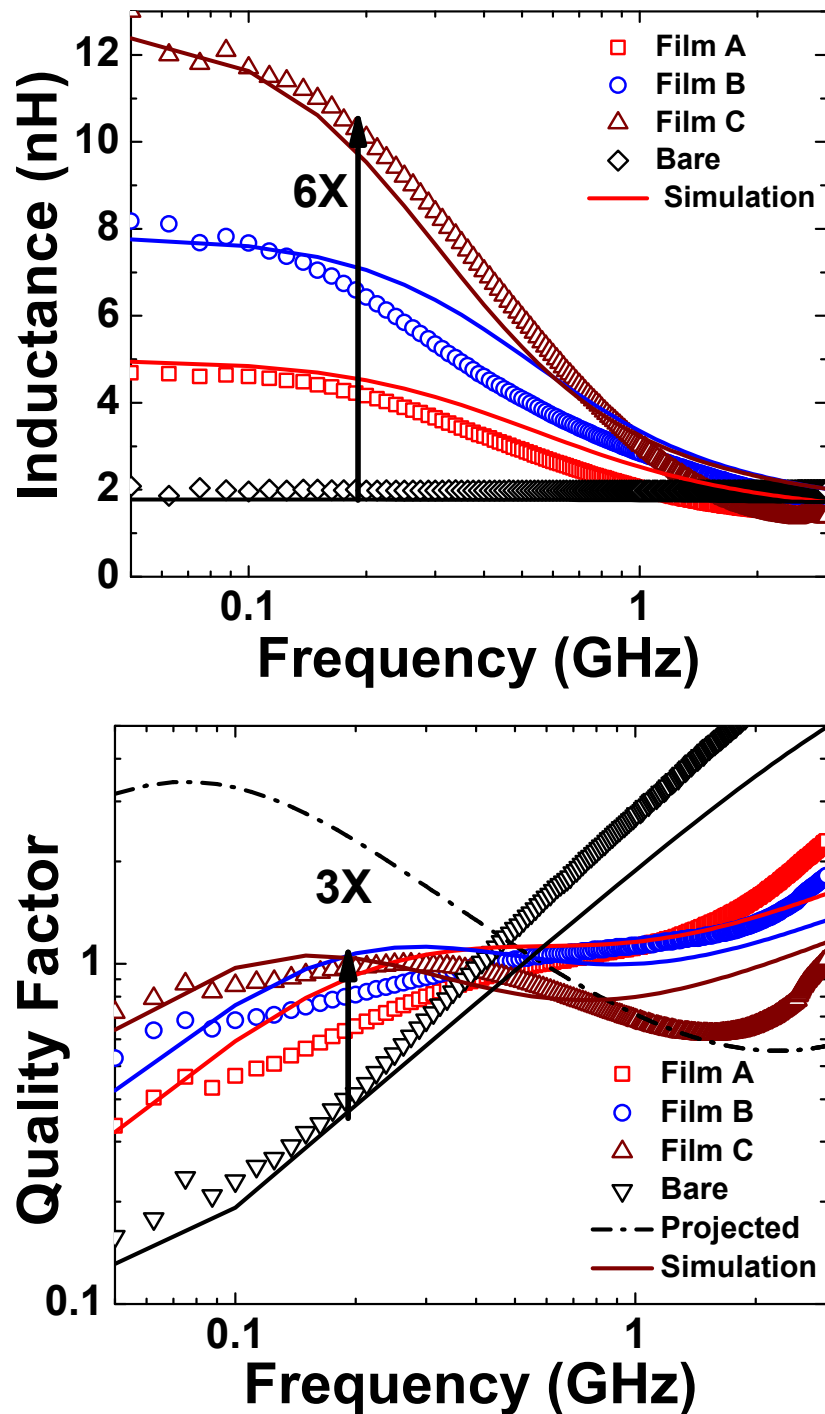


Figure 4.16: Measured and simulated curves of  $L$  and  $Q$  vs. frequency for rectangular inductors with single NiFe ring. Devices with film A ( $\mu=170$ ), B ( $\mu=300$ ) and C ( $\mu=500$ ) were investigated.  $L$  and  $Q$ -factor increase up to 6X and 3X are measured from device C at 200 MHz. Projected  $Q$  was simulated by using 4  $\mu\text{m}$  thick Cu wire with ideal conductivity.

at low frequencies. However, at high frequencies  $L$  and  $Q$  decrease mainly due to the eddy current loss inside NiFe laminations. The total loss can be expressed [75] as follows:



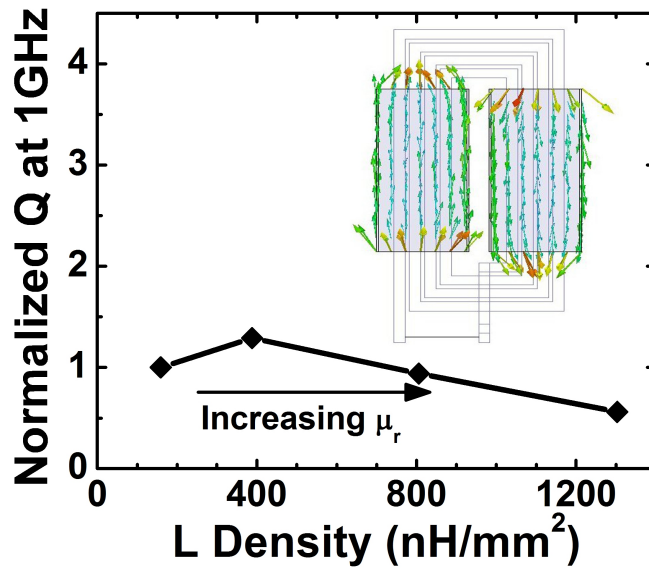


Figure 4.17: Simulation result showing normalized Q at 1GHz plotted against increasing L density. Inset shows eddy current vectors on a simulated structure in HFSS.

$$P = \frac{\omega^2 |\Phi|^2 \sigma l \delta}{8b} \left( \frac{\sinh(\frac{a}{\delta}) - \sin(\frac{a}{\delta})}{\cosh(\frac{a}{\delta}) - \cos(\frac{a}{\delta})} \right) \quad (4.9)$$

where  $a$ ,  $l$ ,  $b$ ,  $\sigma$  are the thickness, length, width and the conductivity of each lamination,  $\Phi$  is the magnetic flux inside the lamination. In this work,  $a \leq 100$  nm, and the skin depth  $\delta$  of NiFe at 1 GHz is approximately  $0.4 \mu\text{m}$ . Since  $a / \delta \ll 1$ , the power loss per unit volume under a uniform magnetic field can be reduced to

$$P = \frac{\pi^2 f^2 B^2 a^2 \sigma^2}{6} \quad (4.10)$$

where  $B$  is the magnetic flux density inside lamination. Since  $B$  is proportional to the permeability of lamination, the eddy current loss is then square dependent of the permeability. Thus, at high frequencies, a higher permeability of NiFe film results in a larger eddy current loss and lower Q-factor as consistently observed in Fig. 4.16. Simulation results show that by changing the thickness of the copper wire from  $2 \mu\text{m}$  to  $4 \mu\text{m}$  while keeping the rest of the parameters same, quality factor greater than 3 can be achieved. ('projected' trace of Fig. 4.16).

To further investigate the permeability dependent L and Q, extensive simulations and comparison of extracted experimental data at both low and high frequencies are performed.

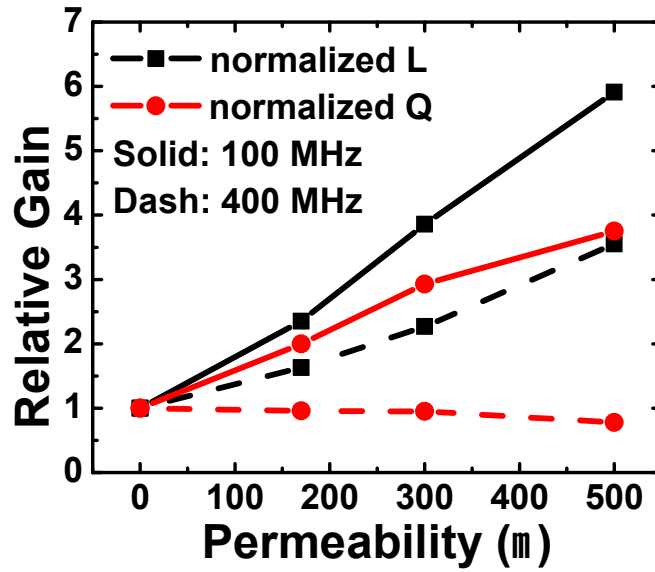


Figure 4.18: Measurement result showing relative gain in L and Q at low and high frequencies.

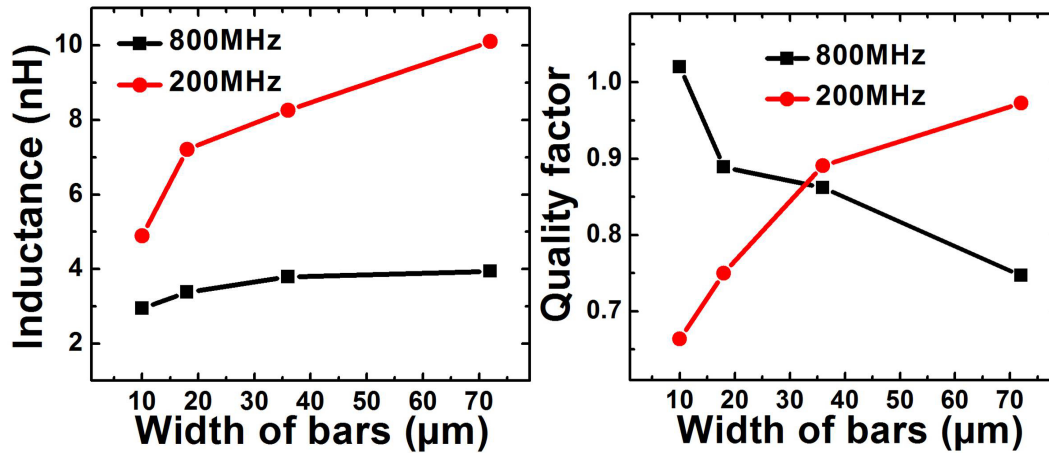


Figure 4.19: The dependence of inductance and quality factor with respect to the width of the magnetic bar at 200 MHz and 800 MHz.

Simulation results in Fig. 4.17 show that at high frequencies (1 GHz), quality factor drops for inductor with higher permeability values of NiFe lamination. The values are normalized to the Q factor of bare inductor. A similar trend is also observed from our measurements as summarized in Fig. 4.18. As expected, at 100 MHz, normalized L and Q increase owing to increased  $\mu$ . However, at 400 MHz, even though L is much higher than the bare inductor, Q factor drops. Additionally, the Q for  $\mu=500$  is lower as compared to  $\mu=170$ , which confirms that the higher permeability magnetic material incurs higher loss due to eddy current, as described in equation 4.10.

## Effect of Magnetic Material Patterning

In order to study the width dependence of eddy current loss, the complete NiFe film is then split into bars with the total volume kept the same. Since the aspect ratio of the bar increases as its width reduces, the local demagnetizing effect increases resulting in a smaller value of permeability as shown in Fig. 4.18. Fig. 4.19 show plots of L and Q vs. change in width of the magnetic bars at 200 MHz and 800 MHz. As expected, with narrower magnetic bars, permeability decreases and L improvement drops from 5.3X (1 ring) to 2.5X (6 rings) at 200 MHz. Similar trend is observed in Q-factor. However, at high frequency (800 MHz), the trend reverses. The Q-factor decreases as the width of magnetic bars increases even though the L value is still higher. It can be concluded that as the width increases, eddy current loss increases which prevents Q-factor from improving at the same rate as L.

### *Inductors using CoZrTa*

The measurement results from Sections 4.2 and 4.2 show that using high conductivity magnetic material such as NiFe can give limited performance enhancement in frequency response and quality factor owing to very high eddy current losses in the magnetic material. Based on the discussion in Section 4.1 it can be concluded that a magnetic material with higher resistivity would help in increasing the bandwidth and quality factor of the inductors. Following this lead, inductors with integrated amorphous CoZrTa (CZT) layers were fabricated. CoZrTA has a resistivity of 100  $\Omega$  cm, which is about five times higher than that of Permalloy (20  $\Omega$  cm). Fig.4.20 shows optical images of the fabricated devices.

Fig.4.21 shows inductance and quality factor vs. frequency for inductors with CZT magnetic rings. For each configuration, the frequency response and quality factor at high frequencies is significantly higher than the fabricated inductors with Permalloy. This is due to the higher resistivity of the magnetic material. The frequency response of the inductors reaches up to about 3-4 GHz and quality factors of 3 at 3 GHz. Using thicker copper wires, the quality factor at low frequency can be further improved.

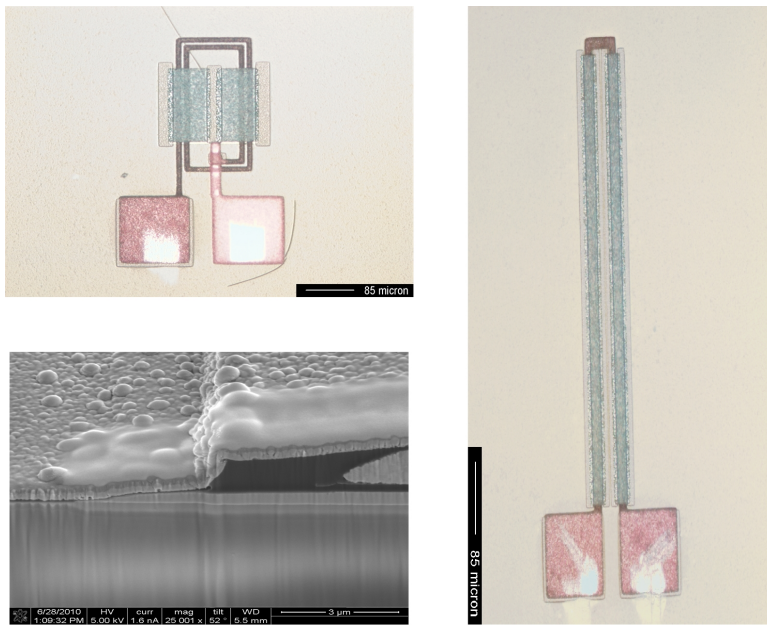
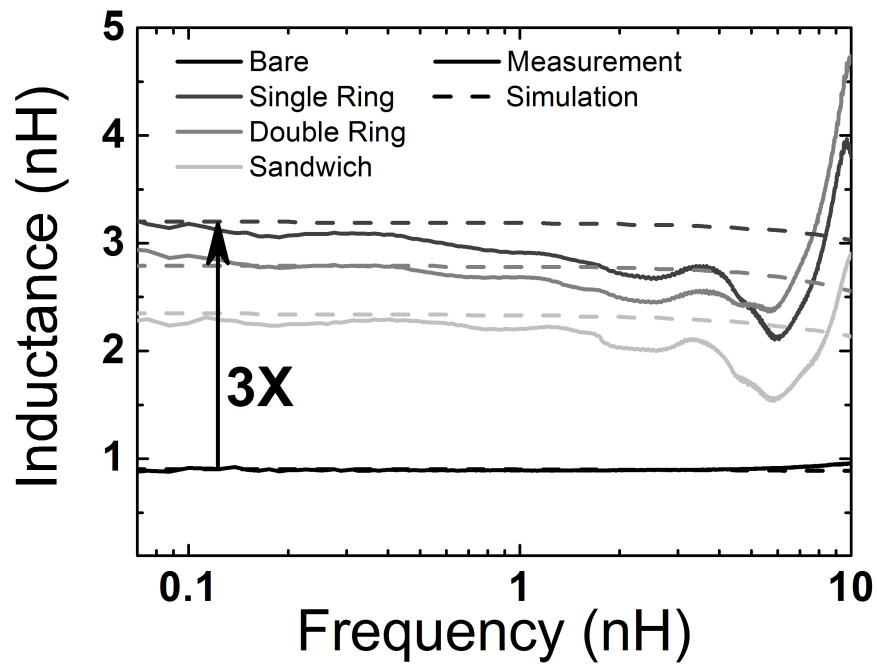


Figure 4.20: Optical images of rectangular inductor (2 turn), stripeline inductor with CoZrTa (CZT) magnetic material and cross section showing inductor metal line and magnetic via connecting the top and bottom films.



### Conclusion

The effect of using high permeability magnetic materials with on-chip spiral inductors to enhance inductance density and quality factor are analyzed by extensive experiment and

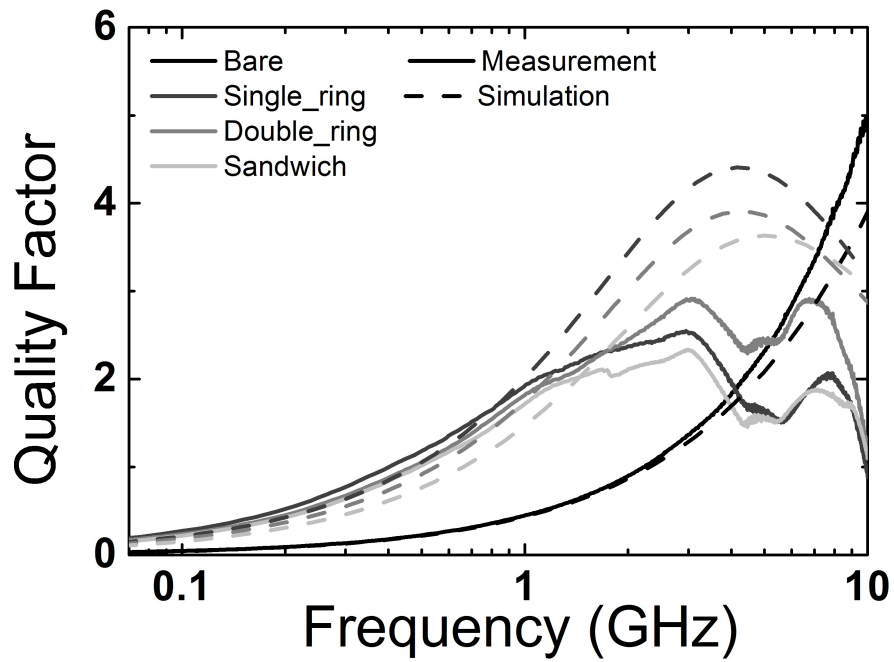
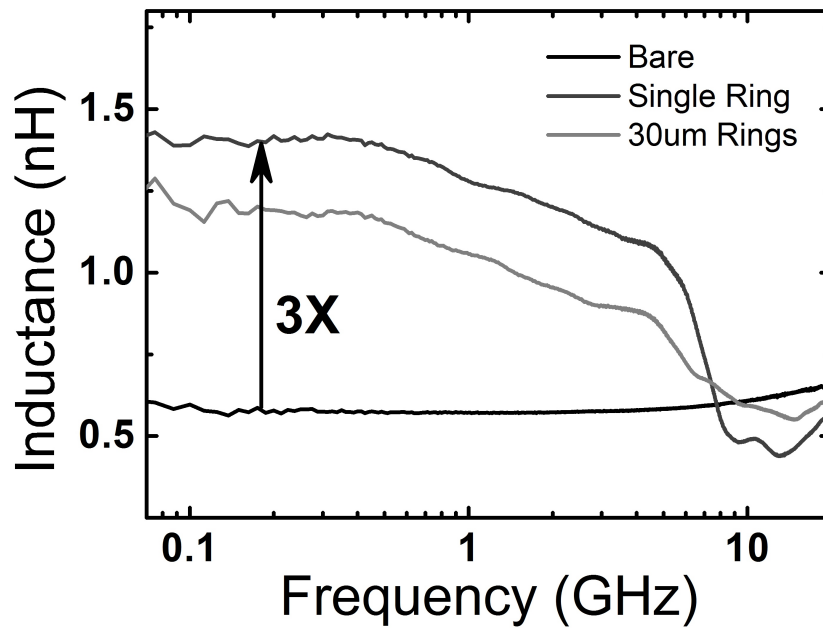


Figure 4.21: Measurement results showing L and Q from rectangular 2 turn inductor with CZT magnetic material.



simulation. Enhancements of 6X in L and 3X in Q-factor at 200 MHz are reported for spiral inductors that also achieve high inductance density of up to  $770 \text{ nH/mm}^2$  using permalloy as the magnetic material. The permeability dependence of inductance, eddy current loss in the magnetic material and quality factor is investigated at both low and high frequencies. The

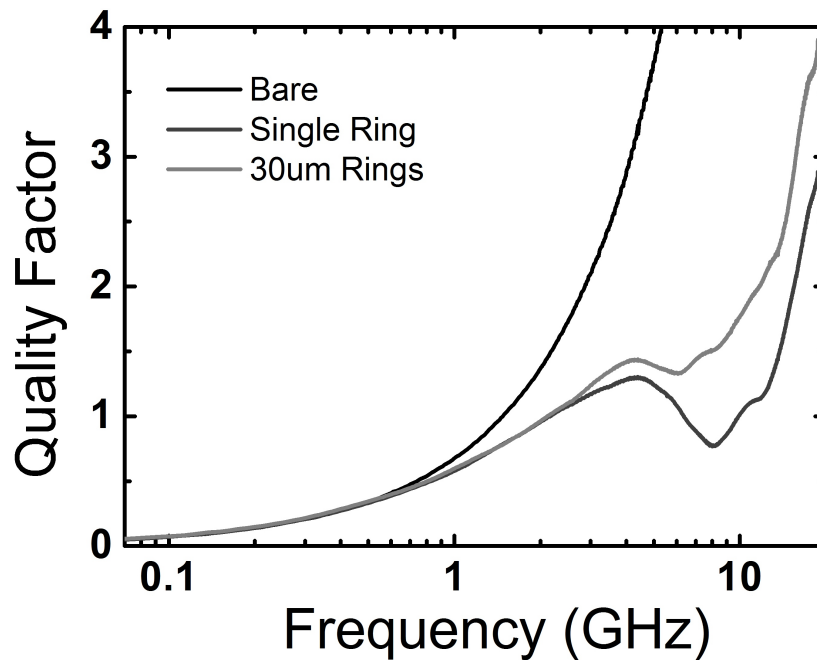


Figure 4.22: Measurement results showing L and Q from stripline inductor with CZT magnetic material.

results show that on-chip inductors with high permeability magnetic material can have significant performance enhancement below 200 MHz. However, at higher frequencies, the high permeability results in severe eddy current loss inside the magnetic material, degrading the performances of inductors. Using high resistivity CoZrTa as the magnetic material results in 3X improvement in inductance, frequency response of up to 4 GHz and peak quality factor of 3 at 3GHz. Hence, patterned magnetic rings with low conductivity are suitable for high frequency applications and continuous magnetic films with thick copper wires for low frequency applications.

#### 4.3 Potential Application: Resonant Clock Distribution

Using inductors with magnetic materials allows savings in chip area, lower inductive cross-talk and possibly higher inductance value and quality factor. Based on the magnetic material properties and the structure chosen, these inductors can be used for different applications. A simulation study exploring the usage of scaled magnetic inductors for resonant clock distribution in microprocessors is presented in [76].

#### 4.4 Potential Application: Magnetic Inductors for Integrated Voltage Regulators

The inductors designed, fabricated and measured in the previous sections were aimed for high frequency operation. The goal was to achieve high inductance and quality factors with minimal area consumption to enable applications in the 1GHz and higher frequency range such as resonant clocking. In this section, we explore the potential of magnetic inductors for on-chip integrated voltage regulators (buck converters) to provide efficient power supply to microprocessors and SoC circuits and allow fine-grained dynamic voltage scaling. The inductor designs for high frequency applications can be scaled for sub-300MHz operation in DC-DC converters by fabricating thicker and wider inductor metal wires.

Based on the inductor specifications derived in Chapter 3, a simulation study is undertaken to identify the optimum design of magnetic inductors for on-chip buck converters. For monolithic integrated buck converters, it has been observed that the resistive loss in the inductor has the maximum contribution towards losses reducing its efficiency. For a given CMOS technology the  $L/R_L$  is constant [56], i.e., if higher inductance is required, the area, required number of turns and length of the metal trace is higher increasing the series resistance. Thicker and wider metals are used for large inductors to keep the resistance low. Smaller inductors (below 10nH) have lower series resistance but the inductance is low requiring higher switching frequency to keep current ripple below the specifications. This reduces efficiency due to switching losses.

Fig. 4.23 shows the required inductance, efficiency and outer diameter of the inductors vs. switching frequency. Simple closed form equations are used to estimate the inductance and area [84]. The metal thickness and width are appropriately scaled for larger inductors so that the equivalent DC series resistance is the same for all cases. Hence, the efficiency curve is not affected for different inductors. It is found that for very large inductors with diameter of 600  $\mu\text{m}$  and higher, very thick and wide ( $10\mu\text{m} \times 20\mu\text{m}$  or higher) wires are required.

Assuming that up to 4X reduction in outer diameter is achieved by incorporating inductors with integrated magnetic rings, from Fig. 4.23 we can observe that for a given inductance value, the area reduces significantly. The same inductance can be achieved with

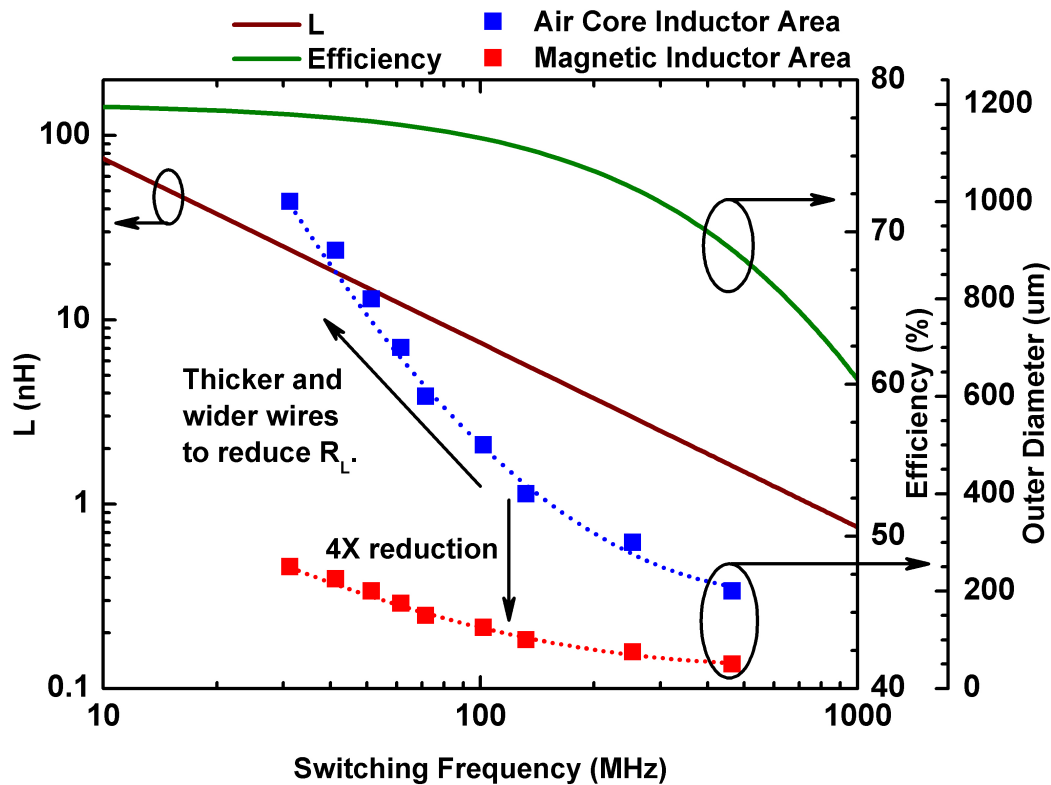


Figure 4.23: Plot showing inductance requirement for on-chip dc-dc converter and efficiency with switching frequency. Small inductors are sufficient if switching frequency is high, however efficiency suffers due to higher switching losses. The area of corresponding air-core on-chip spiral inductors are also plotted. Assuming up to 4X reduction in area due to incorporating magnetic materials, smaller inductors with similar performance can be fabricated resulting in area savings and better efficiency due to lower  $R_L$ .

lesser number of turns and resulting in lower series resistance. This has double benefits.

Significant area savings of up to 16X is possible and lower series resistance can improve the efficiency of the regulator. Reducing the area allows multiple dc-dc converters to be fabricated on chip allowing per-core dynamic voltage scaling. From Fig. 2.20 in Chapter 2 and Fig. 3.7 in Chapter 3, we can conclude the using dedicated dc-dc converters per core utilizing scaled magnetic inductors can enable fine-grained per-core dynamic voltage scaling resulting in power savings.

The general inductor requirements for a multi-phase buck converter operating in the range of 100 MHz to 300 MHz is inductance between 1nH to 10nH (based on the allowable ripple current) and quality factor between 3-10 in the frequency of interest. For comparison purposes, we target a specific design presented in [6] since its design goal is fine-grained



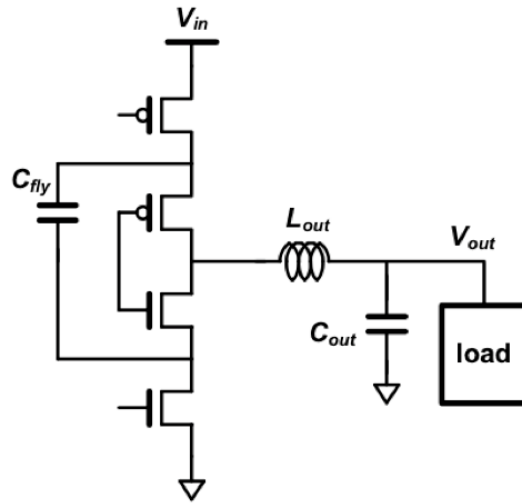


Figure 4.24: Circuit schematic (single phase) of a 3-level on-chip dc-dc converter presented in [6].

dynamic voltage scaling in multi-core SoCs. The 3-level dc-dc converter presented in [6] consists of a 4-phase mixed topology (switch-cap and buck) allowing smaller inductor size. The converter achieves peak efficiency of 77% and can operate at a maximum frequency of 200 MHz. The required inductance per phase is 1nH implemented in 130nm CMOS process with L/R of 2.5nH/ $\Omega$ . This translates to a quality factor of about 4 at 200 MHz. The circuit diagram of a single phase of the design is given in Fig. 4.24 and the fabricated chip micrograph in Fig. 4.25. The on-chip air-core inductors take up 32% of the total circuit area and each inductor is about 500x500 $\mu\text{m}^2$  large. HFSS simulations are run to identify the optimum configuration of inductor (width, diameter, magnetic ring thickness etc.) with magnetic materials to achieve similar L/R values.

#### *Simulation Results*

Since previously presented magnetic inductors possess very low quality factors at frequencies less than 500 MHz, for on-chip regulator applications, the metal width and thickness of the inductor need to be significantly higher. This also allows higher current density in the inductor wires, since currents in the range of 50-200mA per phase needs to be supplied to the load. 3D HFSS simulations are run with the configuration as described in Table 4.2. Fig. 4.26 shows the top view of the simulated structure with 9 magnetic rings per side.

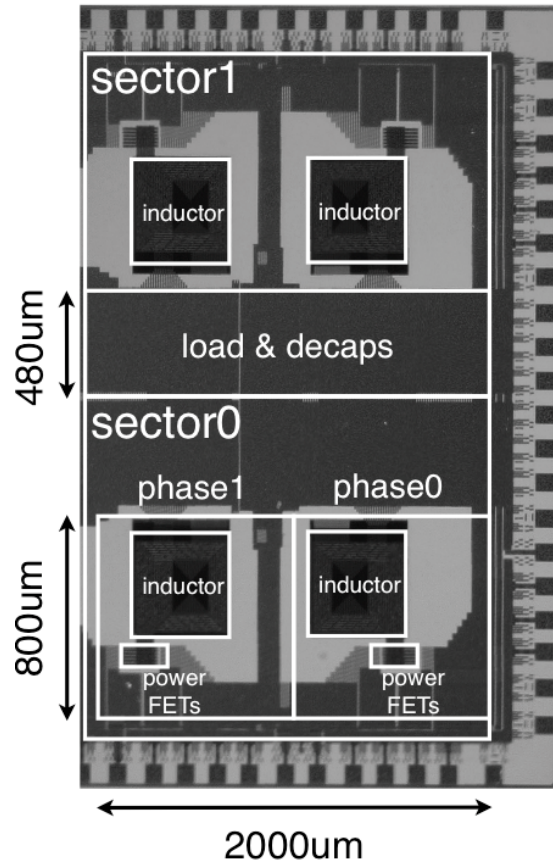


Figure 4.25: Chip micrograph (single phase) of a 3-level on-chip dc-dc converter from [6]. Each inductor occupies an area of about  $500 \times 500 \mu\text{m}^2$  taking up about 32% of the silicon real estate.

Table 4.2: Simulation Parameters

Parameter	Value
Metal	Copper
Number of Turns	1
Metal Width	10 $\mu\text{m}$
Metal Thickness	5 $\mu\text{m}$
Metal Spacing	5 $\mu\text{m}$
Outer Diameter	100 $\mu\text{m}$
Magnetic Ring Thickness	2 $\mu\text{m}$
Magnetic Ring Width	5 $\mu\text{m}$
Magnetic Material	CoZrTa
Number of Rings per side	3, 6, 9
Relative Permeability	800

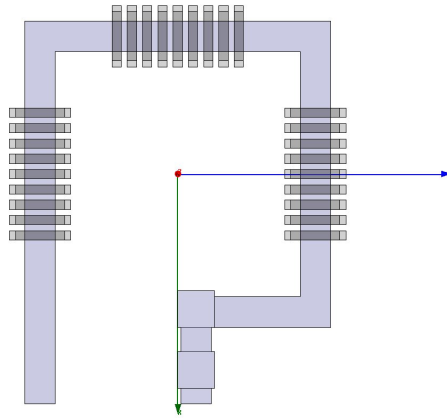


Figure 4.26: Top view of single turn inductor with magnetic rings. Three different configurations with 3, 6 and 9 rings per side are simulated to understand the effect on inductance enhancement, frequency response and Q-factor.

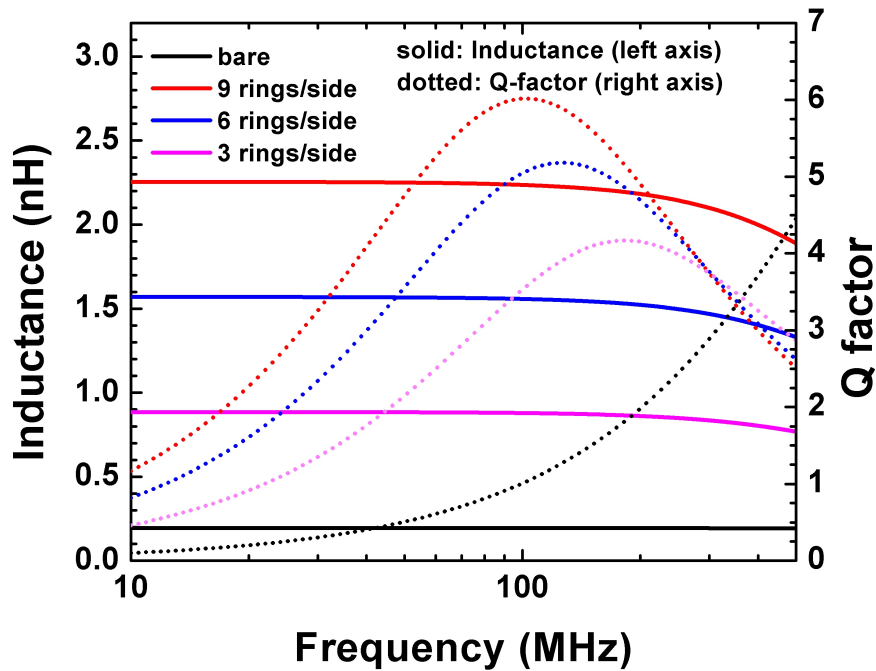


Figure 4.27: Inductance and Quality factor vs frequency for single turn thick metal inductor with 3, 6 and 9 magnetic rings per side for use in on-chip voltage regulator. The frequency region of interest is in the range of 100MHz to 250MHz.

Simulation results showing inductance and quality factor versus frequency are presented in Fig. 4.27. The single turn inductor with 3 magnetic rings per side shows an inductance of about 0.9 nH, close to the required target of 1nH from [6]. The quality factor is exactly 4 at 200MHz as required by the design. Increasing the number of magnetic rings per

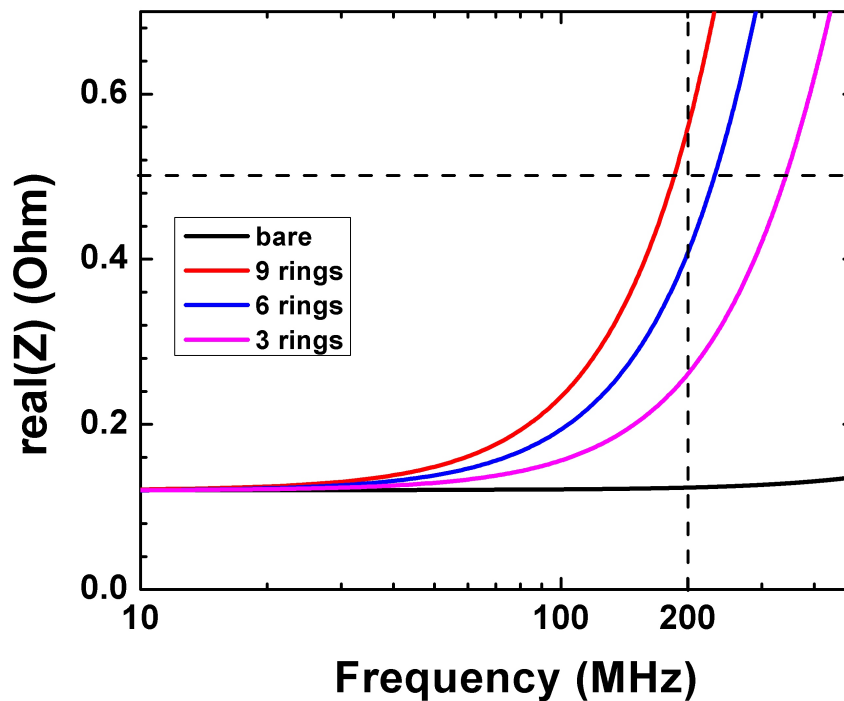


Figure 4.28: Real Impedance of the magnetic inductor with frequency. The dotted lines show allowable resistance of the inductor for high efficiency buck converter. Adding magnetic rings increase eddy current losses limiting their operational frequency to around 200MHz.

side enhances inductance and peak quality factor up to 6 but reduces the frequency response due to additional eddy current loss in the rings. For the 9 ring configuration, peak Q moves to 100 MHz. Adding more number of turns with larger outer diameter can significantly improve the inductance in the range 10 nH. For our target design a single turn inductor with outer diameter of 100  $\mu\text{m}$  is sufficient. Fig. 4.28 plots the real impedance of the magnetic inductors vs. frequency. The real impedance represents the frequency dependent resistance seen at the output of the LC tank of the on-chip buck converter. As shown in Fig. 3.6, Chapter 3, the efficiency of the converter is above 60% for a resistance of about 0.5  $\Omega$ . Considering 0.5  $\Omega$  to be the cut-off, we can see that magnetic inductor with 9 rings can be operated up to 200 MHz beyond which the AC resistance becomes prohibitively high reducing the total efficiency. The inductors with 6 and 3 magnetic rings suffer from relatively lower eddy current losses and can be used at switching frequencies.

### *Area Savings*

The magnetic inductors occupy an area of about  $10^4 \mu\text{m}^2$  area which is 25X smaller compared to the air-core inductors shown in Fig. 4.25. Considering 4 such inductors for each phase, the total chip area savings will be about 30% occupying merely 2% of the total converter area. This reduction in chip area has tremendous potential in implementing multiple on-chip buck converters with integrated inductors for supplying separate power to individual cores and components of an SoC. Assuming a general design where the magnetic inductor has an outer diameter up to 4X smaller than its air-core counterpart, area savings of up to 16X is possible. This allows the fabrication of multiple on-chip dc-dc converters dedicated to individual cores performing per-core fine-grained dynamic voltage scaling.

### 4.5 Conclusion

A thorough study of on-chip spiral inductors with integrated magnetic materials has been presented in this chapter. With detailed theoretical background, 3D electromagnetic simulations and measurement of fabricated inductors it is found that there is a trade-off between inductance enhancement, frequency response and quality factor due to additional eddy current losses in conductive magnetic materials. A short study of magnetic inductors for on-chip resonant clock distribution is presented.

Based on the requirements of inductance and series resistance for integrated dc-dc converters for microprocessors, a system-level study of the design trade-offs and advantages of implementing magnetic inductors as a part of the output filter is presented. It is found that using magnetic inductors can reduce chip-area, improve regulator efficiency allowing implementation of per-core dc-dc converter with additional advantages such as reduction in decoupling capacitance.

### CONCLUSION

With ever-increasing market penetration of desktops, servers, laptops and mobile devices such as smart phones and tablets constantly connected to cloud servers, the specifications for low power high performance microprocessors has become more stringent. Design of modern microprocessors is an extremely complex task and optimization at each level of hierarchy is important to achieve the final goal of energy savings without hampering performance. In this dissertation, we have addressed this issue from both, system and device level optimizations, with each level affecting the next.

The important results presented can be summarized as follows:

- A low power neuromorphic controller for dynamic power and thermal management was designed, simulated and shown to perform better compared to commercial solutions.
- The neuromorphic controller is suitable for per-core fast-transitioning workloads and results in minimal error while providing predictive DVFS.
- Actual implementation of per-core fine-grained DVFS requires high-efficiency on-chip dc-dc converters to modulate the supply voltage as required. A system-level study of dc-dc regulators shows that area of on-chip spiral inductors is prohibitively high for integration of multiple converters. Additionally, high DC series resistance of the inductors reduces the converter efficiency.
- In order to reduce the area of integrated inductors, on-chip spiral inductors using integrated magnetic materials were designed, simulated, fabricated and measured to quantify inductance and area gains and trade-offs in cost and efficiency. Their potential applications in on-chip voltage regulators were explored and it was found that high efficiency converters with very low area overhead is possible.
- Scaled magnetic inductors enable multiple on-chip implementations of power converters helping realizing fine-grained per-core DVFS for multi-core systems. The design

trade-offs in selecting inductor sizes, corresponding chip area savings and the effect on efficiency is presented.

### 5.1 Future Work

Some potential areas of future work are identified as follows:

1. It is important to realize the cost-efficiency trade-off while implementing inductors with magnetic materials. Since the inductor fabrication requires multiple additional processing steps, the processing costs can be high and extensive research is required to explore and develop cost-effective solutions for large scaling manufacturing of inductors with integrated magnetic materials.
2. Additional integration of software programs, performance counters, power controller and voltage regulator is required.
3. A detailed design of the feedback control system that changes the voltage regulator duty cycle based on the output signals of the neuromorphic controller is important.

## REFERENCES

- [1] G. Cauwenberghs, “Neuromorphic cognitive engineering: Large-scale silicon neural systems,” *in 4<sup>th</sup> Decade of the Mind Conf*, 2009.
- [2] W. Kim, M. Gupta, G. Y. Wei and D. Brooks, “System level analysis of fast, per-core dvfs using on-chip switching regulators,” feb. 2008, pp. 123-134.
- [3] S. Abedinpour, et. al., “Monolithic distributed power management for systems-on-chip (soc),” *Annals of Telecommunications*, vol. 59, pp. 938–973, 2004.
- [4] P. Xu, et al. , “Investigation of candidate topologies for 12v vrm,” *in APEC 2002*.
- [5] [Online]. Available: <http://www.ifixit.com/Teardown/Apple-A4-Teardown/2204/1>
- [6] W. Kim, et. al., “A fully-integrated 3-level dc/dc converter for nanosecond-scale dvs with fast shunt regulation,” *in ISSCC 2011*, feb. 2011, pp. 268 –270.
- [7] D. Ma and R. Bondade, “Enabling power-efficient dvfs operations on silicon,” *Circuit and Systems Magazine, IEEE*, vol. 10, no. 1, pp. 14 –30, 2010.
- [8] P. Hazucha et al. , “A 233-mhz 80four-phase dc-dc converter utilizing air-core inductors on package,” *JSSC*, vol. 40, no. 4, pp. 838 – 845, april 2005.
- [9] S. Abedinpour , et al, “A multistage interleaved synchronous buck converter with integrated output filter in 0.18um sige process,” *TPE*, vol. 22- 6, pp. 2164 –2175, 07.



- [10] J. Wibben and R. Harjani, "A high efficiency dc-dc converter using 2nh on-chip inductors," in *VLSI Circuits, 2007 IEEE Symposium on*, 2007, pp. 22 –23.
- [11] C.-S. Kim, S. Bae, H.-J. Kim, S.-E. Nam, and H.-J. Kim, "Fabrication of high frequency dc-dc converter using ti/fetan film inductor," *Magnetics, IEEE Transactions on*, vol. 37, no. 4, pp. 2894 –2896, Jul. 2001.
- [12] K. H. Kim, J. Kim, H. J. Kim, S. H. Han, and H. J. Kim, "A megahertz switching dc/dc converter using febn thin film inductor," *Magnetics, IEEE Transactions on*, vol. 38, no. 5, pp. 3162 – 3164, Sep. 2002.
- [13] G. Schrom, P. Hazucha, J.-H. Hahn, V. Kursun, D. Gardner, S. Narendra, T. Karnik, and V. De, "Feasibility of monolithic and 3d-stacked dc-dc converters for microprocessors in 90nm technology generation," in *Low Power Electronics and Design, 2004. ISLPED '04. Proceedings of the 2004 International Symposium on*, 2004, pp. 263 – 268.
- [14] E. C. Giacomo Indiveri and R. J. Douglas, "Artificial cognitive systems: From vlsi networks of spiking neurons to neuromorphic cognition," *Cognitive Computation*, vol. 1, no. 2, pp. 119–127, June 2009.
- [15] B. Bavarian, "Introduction to neural networks for intelligent control," *IEEE Control Systems Magazine*, vol. 8, no. 2, pp. 3 –7, apr 1988.
- [16] J. H. Wijekoon and P. Dudek, "Compact silicon neuron circuit with spiking and bursting behaviour," *Neural Networks*, vol. 21, no. 2-3, pp. 524 – 534, 2008.
- [17] G. Indiveri, E. Chicca, and R. Douglas, "A vlsi array of low-power spiking neurons and bistable synapses with spike-timing dependent plasticity," *IEEE Trans. on Neural Networks*, vol. 17, no. 1, pp. 211 –221, jan. 2006.
- [18] P. Lichtsteiner, C. Posch, and T. Delbruck, "A 128 times 128 120 db 15 us latency asynchronous temporal contrast vision sensor," *JSSC*, vol. 43, no. 2, pp. 566 –576, feb.

2008.

- [19] V. Chan, S.-C. Liu, and A. van Schaik, "Aer ear: A matched silicon cochlea pair with address event representation interface," *IEEE TCAS I*, vol. 54, no. 1, pp. 48–59, jan. 2007.
- [20] S. Rusu and et al., "A 45nm 8-core enterprise xeon processor," *A-SSCC 2009.*, pp. 9–12, nov. 2009.
- [21] S. Sinha, J. Suh, B. Bakkaloglu, and Y. Cao, "Workload-aware neuromorphic design of low-power supply voltage controller," in *Proceedings of ISLPED 2010*. New York, NY, USA: ACM, 2010, pp. 241–246.
- [22] W. S. McCulloch and W. Pitts, "A logical calculus of the ideas immanent in nervous activity," *Bulletin of Mathematical Biology*, vol. 5, no. 4, pp. 114–133, December 1943.
- [23] Advanced configuration and power interface. [Online]. Available: <http://www.acpi.info/>
- [24] X. Zhong and C.-Z. Xu, "Energy-aware modeling and scheduling for dynamic voltage scaling with statistical real-time guarantee," *IEEE Trans. on Computers*, vol. 56, no. 3, pp. 358–372, march 2007.
- [25] H. Kim, H. Hong, H.-S. Kim, J.-H. Ahn, and S. Kang, "Total energy minimization of real-time tasks in an on-chip multiprocessor using dynamic voltage scaling efficiency metric," *Computer-Aided Design of Integrated Circuits and Systems, IEEE Transactions on*, vol. 27, no. 11, pp. 2088–2092, 2008.
- [26] J. Lorch and A. Smith, "Pace: a new approach to dynamic voltage scaling," *IEEE Trans. on Computers*, vol. 53, no. 7, pp. 856–869, july 2004.
- [27] A. Sinha and A. Chandrakasan, "Dynamic voltage scheduling using adaptive filtering of workload traces," *Int. Conf. on VLSI Design, 2001.*, pp. 221–226, 2001.

- [28] S.-Y. Bang, K. Bang, S. Yoon, and E.-Y. Chung, "Run-time adaptive workload estimation for dynamic voltage scaling," *IEEE Trans. on Computer-Aided Design of Integrated Circuits and Systems*, vol. 28, no. 9, pp. 1334–1347, sept. 2009.
- [29] A. Golda and A. Kos, "Neural processor as a dynamic power manager for digital systems," in *MICAI '08*. Berlin, Heidelberg: Springer-Verlag, 2008, pp. 333–342.
- [30] Y. Gu and S. Chakraborty, "Control theory-based dvs for interactive 3d games," june 2008, pp. 740–745.
- [31] G. Dhiman and T. Rosing, "System-level power management using online learning," *IEEE Trans. on Computer-Aided Design of Integrated Circuits and Systems*, vol. 28, no. 5, pp. 676–689, may 2009.
- [32] R. Cochran and S. Reda, "Consistent runtime thermal prediction and control through workload phase detection," in *Design Automation Conference (DAC), 2010 47th ACM/IEEE*, 2010, pp. 62–67.
- [33] T. Mitra and R. Jayaseelan, "Dynamic thermal management via architectural adaptation," in *Design Automation Conference, 2009. DAC '09. 46th ACM/IEEE*, 2009, pp. 484–489.
- [34] A. Coskun, T. Rosing, and K. Gross, "Utilizing predictors for efficient thermal management in multiprocessor socs," *Computer-Aided Design of Integrated Circuits and Systems, IEEE Transactions on*, vol. 28, no. 10, pp. 1503–1516, 2009.
- [35] J. Yang, X. Zhou, M. Chrobak, Y. Zhang, and L. Jin, "Dynamic thermal management through task scheduling," in *Performance Analysis of Systems and software, 2008. ISPASS 2008. IEEE International Symposium on*, 2008, pp. 191–201.
- [36] *Intel® Core™ i7-600, i5-500, i5-400 and i3-300 Mobile Processor Series (Volume One)*, Intel Corporation, Jan. 2010.

- [37] M. Floyd, B. Brock, M. Ware, K. Rajamani, A. Drake, C. Lefurgy, and L. Pesantez. (2010) Harnessing the adaptive energy management features of the power7 chip. [Online]. Available: <http://www.hotchips.org/conference-archives/hot-chips-22/>
- [38] [Online]. Available: [ptm.asu.edu](http://ptm.asu.edu)
- [39] C. Mead, *Analog VLSI and neural systems*. Boston, MA, USA: Addison-Wesley Longman Publishing Co., Inc., 1989.
- [40] W. Maass and C. M. Bishop, Eds., *Pulsed neural networks*. Cambridge, MA, USA: MIT Press, 1999.
- [41] E. Chicca and et al., “A vlsi recurrent network of integrate-and-fire neurons connected by plastic synapses with long-term memory,” *IEEE Trans. on Neural Networks*, vol. 14, no. 5, pp. 1297 – 1307, sept. 2003.
- [42] S. Mitra, S. Fusi, and G. Indiveri, “Real-time classification of complex patterns using spike-based learning in neuromorphic vlsi,” *IEEE Trans. on Biomedical Circuits and Systems*, vol. 3, no. 1, pp. 32 –42, feb. 2009.
- [43] L. S. Smith, “Implementing neural models in silicon,” *Handbook of nature-inspired and innovative computing*, 2004.
- [44] [Online]. Available: [www.nangate.com](http://www.nangate.com)
- [45] V. Pallipadi and A. Starikovskiy, “The ondemand governor,” in *Linux Symposium '06*, 2006, pp. 223–238.
- [46] L. T. Clark, F. Ricci, and W. E. Brown, “Dynamic voltage scaling with the xscale embedded microprocessor,” *Adaptive Techniques for Dynamic Processor Optimization*, 2008.

- [47] perfmon2: the hardware-based performance monitoring interface for linux. [Online]. Available: <http://perfmon2.sourceforge.net/>
- [48] (2009, Sept.) Using intel® vtune™ performance analyzer to optimize software for the intel(r) core(tm) i7 processor family. [Online]. Available: <http://software.intel.com/en-us/articles/using-intel-vtune-performance-analyzer-to-optimize-software-for-the-intelr-core-tm-i7-processor-family/>
- [49] K. Meng, R. Joseph, R. P. Dick, and L. Shang, “Multi-optimization power management for chip multiprocessors,” in *PACT '08*. New York, NY, USA: ACM, 2008, pp. 177–186.
- [50] D. Brooks, V. Tiwari, and M. Martonosi, “Wattch: a framework for architectural-level power analysis and optimizations,” *SIGARCH Comput. Archit. News*, vol. 28, pp. 83–94, May 2000.
- [51] D. J. Frank, R. Puri, and D. Toma, “Design and cad challenges in 45nm cmos and beyond,” in *ICCAD '06*. New York, NY, USA: ACM, 2006, pp. 329–333.
- [52] A. Naveh, E. Rotem, A. Mendelson, S. Gochman, R. Chabukswar, K. Krishnan, and A. Kumar, “Power and thermal management in the intel core duo processor,” *Intel Technology Journal*, May 2006.
- [53] J. Park, D. Shin, N. Chang, and M. Pedram, “Accurate modeling and calculation of delay and energy overheads of dynamic voltage scaling in modern high-performance microprocessors,” in *Proceedings of ISLPED 2010*. New York, NY, USA: ACM, 2010, pp. 419–424.
- [54] M. Popovich, “High performance power distribution networks with on-chip decoupling capacitors for high performance power distribution networks with on-chip decoupling

capacitors for nanoscale integrated circuits,” Ph.D. dissertation, University of Rochester, 2007.

- [55] G. Schrom, P. Hazucha, J. Hahn, D. Gardner, B. Bloechel, G. Dermer, S. Narendra, T. Karnik, and V. De, “A 480-mhz, multi-phase interleaved buck dc-dc converter with hysteretic control,” in *Power Electronics Specialists Conference, 2004. PESC 04. 2004 IEEE 35th Annual*, vol. 6, june 2004, pp. 4702 – 4707 Vol.6.
- [56] G. Schrom, P. Hazucha, F. Paillet, D. Gardner, S. Moon, and T. Karnik, “Optimal design of monolithic integrated dc-dc converters,” in *Integrated Circuit Design and Technology, 2006. ICICDT '06. 2006 IEEE International Conference on*, 0-0 2006, pp. 1 –3.
- [57] G. Schrom, P. Hazucha, F. Paillet, D. J. Rennie, S. T. Moon, D. S. Gardner, T. Kamik, P. Sun, T. T. Nguyen, M. J. Hill, K. Radhakrishnan, and T. Memioglu, “A 100mhz eight-phase buck converter delivering 12a in 25mm<sup>2</sup> using air-core inductors,” in *Applied Power Electronics Conference, APEC 2007 - Twenty Second Annual IEEE*, 25 2007-march 1 2007, pp. 727 –730.
- [58] M. Alimadadi, S. Sheikhaei, G. Lemieux, S. Mirabbasi, and P. Palmer, “A 3ghz switching dc-dc converter using clock-tree charge-recycling in 90nm cmos with integrated output filter,” in *Solid-State Circuits Conference, 2007. ISSCC 2007. Digest of Technical Papers. IEEE International*, feb. 2007, pp. 532 –620.
- [59] P. Li, L. Xue, P. Hazucha, T. Karnik, and R. Bashirullah, “A delay-locked loop synchronization scheme for high-frequency multiphase hysteretic dc-dc converters,” *Solid-State Circuits, IEEE Journal of*, vol. 44, no. 11, pp. 3131 –3145, nov. 2009.
- [60] N. Sturcken, M. Petracca, S. Warren, L. P. Carloni, A. V. Peterchev, and K. L. Shepard, “An integrated four-phase buck converter delivering 1a/mm<sup>2</sup> with 700ps controller delay and network-on-chip load in 45-nm soi,” in *Custom Integrated Circuits Conference (CICC), 2011 IEEE*, sept. 2011, pp. 1 –4.

- [61] J. Wibben and R. Harjani, "A high-efficiency dc-dc converter using 2 nh integrated inductors," *Solid-State Circuits, IEEE Journal of*, vol. 43, no. 4, pp. 844–854, april 2008.
- [62] S. Kudva and R. Harjani, "Fully integrated on-chip dc-dc converter with a 450x output range," in *Custom Integrated Circuits Conference (CICC), 2010 IEEE*, sept. 2010, pp. 1–4.
- [63] D. S. Gardner, G. Schrom, F. Paillet, B. Jamieson, T. Karnik, and S. Borkar, "Review of on-chip inductor structures with magnetic films," *IEEE Trans. on Mag.*, vol. 45, no. 10, Sp. Iss. SI, pp. 4760–4766, Oct 2009.
- [64] W. Xu, S. Sinha, F. Pan, T. Dastagir, Y. Cao, and H. Yu, "Improved frequency response of on-chip inductors with patterned magnetic dots," *IEEE Electron Device Letters*, vol. 31, no. 3, pp. 207–209, Mar 2010.
- [65] C. Yang, F. Liu, T.-L. Ren, L.-T. Liu, G. Chen, X.-K. Guan, A. Wang, and H.-G. Feng, "Ferrite-integrated on-chip inductors for rf ics," *Electron Dev. Let.*, vol. 28, no. 7, pp. 652–655, July 2007.
- [66] R.-F. Jiang, N. Shams, M. Rahman, and C.-H. Lai, "Exchange-coupled irmn/cofe multilayers for rf-integrated inductors," *Magnetics, IEEE Transactions on*, vol. 43, no. 10, pp. 3930–3932, Oct. 2007.
- [67] C. Yang, F. Liu, X. Wang, J. Zhan, A. Wang, T.-L. Ren, L.-T. Liu, H. Long, Z. Wu, and X. Li, "Investigation of on-chip soft-ferrite-integrated inductors for rf ics-part ii: experiments," *IEEE Tran. on Electron Dev.*, vol. 56, no. 12, pp. 3141–3148, dec. 2009.
- [68] Y. Zhuang, M. Vroubel, B. Rejaei, J. N. Burghartz, and K. Attenborough, "Magnetic properties of electroplated nano/microgranular nife thin films for rf application," vol. 97, no. 10. AIP, 2005, p. 10N305.

- [69] M. Yamaguchi, S. Bae, K. H. Kim, K. Tan, T. Kusumi, and K. Yamakawa, "Ferromagnetic rf integrated inductor with closed magnetic circuit structure," in *IEEE MTT-S International*, June 2005, p. 4 pp.
- [70] A. Gromov, V. Korenivski, K. Rao, R. van Dover, and P. Mankiewich, "A model for impedance of planar rf inductors based on magnetic films," *Magnetics, IEEE Transactions on*, vol. 34, no. 4, pp. 1246–1248, Jul. 1998.
- [71] B. Jamieson, T. O'Donnell, P. McCloskey, D. S. Gardner, and S. Roy, "Optimization of magnetic enhancement layers for high-frequency stripline micro-inductors," *Journal of Magnetism and Magnetic Materials*, vol. 322, no. 9-12, pp. 1527–1531, 2010, proceedings of the Joint European Magnetic Symposia. [Online]. Available: <http://www.sciencedirect.com/science/article/B6TJJ-4XH5MPY-5/2/b30ce35aa75f44a06f619172638ed9f2>
- [72] Y. Zhuang, B. Rejaei, E. Boellaard, M. Vroubel, and J. Burghartz, "Integrated solenoid inductors with patterned, sputter-deposited cr/fe10co90/cr ferromagnetic cores," *Electron Device Letters, IEEE*, vol. 24, no. 4, pp. 224–226, April 2003.
- [73] T. Dastagir, W. Xu, S. Sinha, H. Wu, Y. Cao, and H. Yu, "Tuning the Permeability of Permalloy Films for On-chip Inductor Applications," *Applied Physics Letters*, to be published.
- [74] V. Korenivski and R. B. van Dover, "Magnetic film inductors for radio frequency applications," *J. Appl. Phys.*, vol. 82, p. 5247, 1997.
- [75] M. Latour, "Note on losses in sheet iron at radio frequencies," *Proceedings of the Institute of Radio Engineers*, vol. 7, no. 1, pp. 61–71, FEB 1919.
- [76] S. Sinha, W. Xu, J. Velamala, T. Dastagir, B. Bakkaloglu, H. Yu, and Y. Cao, "Enabling resonant clock distribution with scaled on-chip magnetic inductors," in *Computer*



*Design, 2009. ICCD 2009. IEEE International Conference on, 2009, pp. 103–108.*

- [77] C. Anderson et.al., “Physical design of a fourth-generation power ghz microprocessor,” *ISSCC*, pp. 232–233, 451, 2001.
- [78] F. O’Mahony, C. Yue, M. Horowitz, and S. Wong, “A 10-ghz global clock distribution using coupled standing-wave oscillators,” *Solid-State Circuits, IEEE Journal of*, vol. 38, no. 11, pp. 1813–1820, Nov. 2003.
- [79] J. Wood, T. Edwards, and S. Lipa, “Rotary traveling-wave oscillator arrays: a new clock technology,” *Solid-State Circuits, IEEE Journal of*, vol. 36, no. 11, pp. 1654–1665, Nov 2001.
- [80] S. Chan, K. Shepard, and P. Restle, “Uniform-phase uniform-amplitude resonant-load global clock distributions,” *Solid-State Circuits, IEEE Journal of*, vol. 40, no. 1, pp. 102–109, Jan. 2005.
- [81] S. Chan et. al., “Distributed differential oscillators for global clock networks,” *Solid-State Circuits, IEEE Journal of*, vol. 41, no. 9, pp. 2083–2094, Sept. 2006.
- [82] A. Muhtaroglu, G. Taylor, T. Rahal-Arabi, and K. Callahan, “On-die droop detector for analog sensing of power supply noise,” *VLSI Circuits Symposium*, pp. 193–196, 2003.
- [83] B. Mesgarzadeh, M. Hansson, and A. Alvandpour, “Jitter characteristic in charge recovery resonant clock distribution,” *Solid-State Circuits, IEEE Journal of*, vol. 42, no. 7, pp. 1618–1625, July 2007.
- [84] S. Jenei, B. Nauwelaers, and S. Decoutere, “Physics-based closed-form inductance expression for compact modeling of integrated spiral inductors,” *Solid-State Circuits, IEEE Journal of*, vol. 37, no. 1, pp. 77–80, jan 2002.