

LIMITS OF HOTSPOT DETECTION AND PREDICTION IN MICROPROCESSORS

A DISSERTATION

SUBMITTED TO THE DEPARTMENT OF MECHANICAL ENGINEERING

AND THE COMMITTEE ON GRADUATE STUDIES

OF STANFORD UNIVERSITY

IN PARTIAL FULFILLMENT OF THE REQUIREMENTS FOR THE DEGREE OF

DOCTOR OF PHILOSOPHY

Josef L. Miler

June 2012

© 2012 by Josef Miler. All Rights Reserved.

Re-distributed by Stanford University under license with the author.

This dissertation is online at: <http://purl.stanford.edu/jc805vt4819>

I certify that I have read this dissertation and that, in my opinion, it is fully adequate in scope and quality as a dissertation for the degree of Doctor of Philosophy.

Kenneth Goodson, Primary Adviser

I certify that I have read this dissertation and that, in my opinion, it is fully adequate in scope and quality as a dissertation for the degree of Doctor of Philosophy.

Thomas Kenny

I certify that I have read this dissertation and that, in my opinion, it is fully adequate in scope and quality as a dissertation for the degree of Doctor of Philosophy.

Mehdi Asheghi

Approved for the Stanford University Committee on Graduate Studies.

Patricia J. Gumpert, Vice Provost Graduate Education

This signature page was generated electronically upon submission of this dissertation in electronic format. An original signed hard copy of the signature page is on file in University Archives.

ABSTRACT

Microprocessor hotspots are a major reliability concern with heat fluxes as much as 20 times greater than those found elsewhere on the chip. Chip hotspots also augment thermo-mechanical stress at chip-package interfaces which can lead to failure during cycling. Because highly localized, transient chip cooling is both technically challenging and costly, chip manufacturers are using dynamic thermal management (DTM) techniques that reduce hotspots by throttling chip power. Uncertainty in heat flux profiles and chip thermal response leads to either excessively conservative DTM schemes and underutilized computational potential or device overheating and associated system failure risks. Improved techniques for quantifying uncertainty and accurately predicting transient thermal response are needed for maximizing reliable chip performance.

A review is conducted of recent advancements in sensor design, laboratory thermometry, sensor allocation, and thermal signal processing for dynamic thermal management. Representative examples of DTM implementation are provided. Quantitative error estimates are compared for semiconductor thermal sensors and thermometry techniques, and improvements in thermal sensor placement and signal processing are presented.

A simulation method is developed to determine the accuracy and resolution at which hotspot heat fluxes can be measured using distributed temperature sensors. The model is based on a novel, computationally-efficient, inverse heat transfer solution. The uncertainties in the hotspot location and intensity are computed for randomized chip

heat flux profiles for varying sensor spacing, sensor vertical proximity, sensor error, and chip thermal properties. For certain cases the inverse solution method decreases mean absolute error in the heat flux profile by more than 30%. These results and simulation methods can be used to determine the optimal spacing of distributed temperature sensor arrays for hotspot management in chips.

To enable on-chip modeling of transient temperature response in a semiconductor device subjected to arbitrarily varying power excitations, an original model compression technique is employed. A network identification deconvolution (NID) method is used to characterize device thermal response from either numerical or experimental results. To compute the transient response to an arbitrary power input, a highly-efficient technique based on digital signal processing is employed. An Infinite Impulse Response (IIR) filter dramatically reduces the required computations to achieve accurate response. The technique provides the best possible scaling of overall computation time and significantly reduces memory constraints. This improvement enables implementation of sophisticated runtime dynamic thermal management algorithms for high-power integrated circuit architectures.

In sum, the present doctoral research offers a multi-faceted approach to managing measurement uncertainty in dynamic thermal management schemes and predicting hotspot response to facilitate optimal chip performance within reliable operating conditions.

ACKNOWLEDGMENTS

I thank my advisor, Professor Ken Goodson, for his consistent support during my graduate research. He encouraged me to formulate my research with both technical precision and broad vision. In that context, he gave me the liberty to experiment with new ideas and collaborate across a variety of projects. I thank Professor Mehdi Asheghi for his technical guidance. His advice dramatically shaped the direction of my dissertation work. I thank Professor Juan Santiago for his enthusiastic support during the first few years of graduate school. I recall many fun, thought-provoking conversations with him. I also thank Professors John Eaton, Tom Kenny, and Godfrey Mungal for the advice they gave me at various points throughout my time at Stanford.

I thank Dr. Maxat Touzelbaev for the creative spark he added to our research collaborations, and I thank Dr. Gamal Refai-Ahmed for many stimulating conversations. I am grateful to research colleagues Milnes David, Roger Flynn, Saniya LeBlanc, and Julie Steinbrenner for their kind mentorship, and I thank the many members of our research group for their help. I thank my friends in the San Francisco bay area, especially Eric Leroux, Saahil Mehra, Kevin Rice, and Brennan Sherry, who helped me grow both professionally and personally during graduate school.

I am deeply thankful to my four amazing sisters, Alison, Katie, Michelle, and Kris, for their encouragement, advice, and patience. I dedicate this thesis to my father and mother for their unequivocal love and support. From my mother, I have learned determination; from my father, resourcefulness.

FUNDING ACKNOWLEDGMENTS

I gratefully acknowledge funding support from the Stanford Department of Mechanical Engineering Graduate Teaching and Research Fellowship and from Advanced Micro Devices (AMD) Inc. as part of the Semiconductor Research Consortium (SRC).

TABLE OF CONTENTS

Abstract.....	iv
Acknowledgments	vi
Funding Acknowledgments.....	vii
Table of Contents	viii
List of Tables	ix
List of Illustrations	x
Nomenclature	xiv
Chapter 1: Introduction.....	1
1.1. Thermal Management Challenges for Next-Generation Integrated Circuits... 1	
1.2. Introduction to Dynamic Thermal Management	5
1.3. Outline of Doctoral Research	6
Chapter 2: Thermometry and Error Reduction for Dynamic Thermal Management	9
2.1. Review Studies Related to Dynamic Thermal Management	9
2.2. Representative Realizations of Dynamic Thermal Management	11
2.3. Thermal Sensor Design.....	16
2.4. Laboratory Thermography Techniques	24
2.5. Thermal Sensor Arrays and Signal Processing.....	31
Chapter 3: Uncertainty in Hotspot Detection	38
3.1. Introduction to Hotspot Detection	38
3.2. Simulation Methodology	42
3.3. Spatial Frequency Domain Inverse Heat Transfer Solution	51
3.4. Simulation Results	59
Chapter 4: Fast Calculation of Temperature Evolution in Electronic Devices	70
4.1. Introduction to Transient Hotspot Modeling	70
4.2. Thermal Modeling Approach.....	76
4.3. Model Verification and Applications	84
4.4. Summary	95
Chapter 5: Conclusions.....	97
Bibliography	101

LIST OF TABLES

This thesis does not contain tables.

LIST OF ILLUSTRATIONS

Figure 1: Variability of chip power for Intel Itanium microprocessor family. Figure adapted from [5].3
Figure 2: Variability in power by application for Intel Itanium microprocessor family. Figure adapted from [5].4
Figure 3: Layout of Foxton Technology Controller used for Intel Itanium microprocessor family [6].	...12
Figure 4: Thermal system block diagram for Intel Foxton Technology Controller [12]. 13
Figure 5: Floorplan for AMD Quadcore Opteron processor showing thermal diodes and thermal processing centers [19]. 14
Figure 6: Schematic of digital thermal sensor based on a thermal diode [33]. 17
Figure 7: Representative results of linearity and sensitivity of sensors designed proposed by Aldrete-Vidrio <i>et al.</i> Solid lines indicate least-squares fit to linear approximation and corresponding parameters are shown on the right. [34]. 19
Figure 8: Comparison of proposed thermal measurement technique with simulation tools [35].20
Figure 9: Improvement in thermal measurement due to sensor group calibration [29].21
Figure 10: Temperature estimation error as a function of power estimation error [29].22
Figure 11: Schematic of simulated ring oscillator digital thermal sensor [38].23
Figure 12: Calibrated sensor response for three corner cases [38].24
Figure 13: Error in IR thermography for alternative chip coatings [39].26
Figure 14: Schematic of novel IR thermography technique [41].27
Figure 15: Schematic of fast transient IR thermography technique [41].27
Figure 16: Temperature profile of HFET as measured by Raman spectroscopy and IR thermography. The superior spatial resolution of Raman spectroscopy enables improved peak temperature measurement. [42]28

Figure 17: Temperature response measured by built-in sensor and corresponding phase shift response for transient interferometric mapping [57].	30
Figure 18: Schematic of transient interferometric mapping set-up. Abbreviations are: SLD: superluminescent diode; AOM: acousto-optical modulator; PBSC: polarizing beam splitter cube; BS: beam splitter; L: lens; DUT: device under test; FPA: focal plane array camera [57].	31
Figure 19: Schematic of two-stage hotspot interpolation technique proposed by Long <i>et al.</i> [59]. A low-resolution sensor array (a) is used to estimate the hotspot location and then a local high-resolution array (b) is activated in the vicinity of the hotspot.	33
Figure 20: Error results for full thermal characterization using various thermal sensor array interpolation schemes [18].	34
Figure 21:: Error results for hotspot estimation using various thermal sensor array interpolation schemes [18].	34
Figure 22: Schematic of two techniques proposed for (a) off-line sensor set-up and (b) runtime temperature estimation [64].	35
Figure 23: Schematic of interconnect lines required for monitoring a single core [60].	36
Figure 24: Layout of interconnect lines for single core monitoring [60].	37
Figure 25: Schematic of model geometry. An arbitrary heat flux profile is applied on the bottom boundary. The boundary condition on all sidewalls is adiabatic; the boundary condition on the top surface is uniform heat transfer.	43
Figure 26: Representative images of each of the four main steps in the simulation methodology. The inputted heat flux profile (a) is used as a reference for determining the error in (d) the calculated heat flux profile.	44
Figure 27: Block diagram of numerical approach used for determining hotspot detection accuracy. FFT and IFFT refer to the Fast Fourier Transform and the Inverse Fast Fourier Transform, respectively.	47
Figure 28: Heat flux profiles used for resolution study, referred to as Case I and II respectively. Both heat flux profiles have equivalent average heat flux and produce similar temperature response profiles. The solution methods are tested for their ability to correctly resolve these heat flux profiles.	50
Figure 29: Schematic of two-port terminal network [67].	52

Figure 30: Representative plots of inverse solution transfer function. Plots show two-dimensional shape of the transfer function (a) without filtering and (b) with filtering. (c) Values of the transfer function for varying x-direction spatial frequency and for y-direction frequency of zero (labeled “on-axis”) as well as for maximum y-direction frequency (labeled “off-axis”). The filter roll-off occurs at approximately 4000 [m ⁻¹].	56
Figure 31: Average mean absolute error (MAE) for varying numbers of randomized heat flux profiles for (a) variable heat flux and (b) binary heat flux. Results for both cases are independent of the number of heat flux profiles for more than 50 heat flux profiles.	58
Figure 32: Demonstration of the averaging technique for (a) variable heat flux and (b) binary heat flux. Results for 50 heat flux profiles are shown. The bold black line indicates the average value.	60
Figure 33: Effects on uncertainty of variable versus binary inputted heat flux profile for varying vertical proximity between sensor and circuit level. The binary heat flux profile results in substantially lower MAE.	61
Figure 34: Uncertainty in calculated heat flux profile for varying convective heat transfer coefficient. The inverse solution method is much less sensitive to heat transfer coefficient than the direct interpretation method.	62
Figure 35: Uncertainty in calculated heat flux profile for varying sensor error at a vertical proximity of (a) 2.575 μm and (b) 7.53 μm . The inverse solution method is susceptible to sensor error at high spatial frequency. The MAE for the direct interpretation method is not affected by varying sensor error.	64
Figure 36: Uncertainty in calculated heat flux profile for varying vertical proximity between the sensor and circuit levels for zero sensor error. For most cases, large changes in vertical proximity yield modest improvements in heat flux uncertainty.	65
Figure 37: Plot of minimum accurate sampling frequency as a function of vertical proximity between chip and sensor level for heat transfer coefficient of $10^4 \text{ W/m}^2\text{-K}$. The inverse solution method is accurate in the shaded region. The direct interpretation technique is inaccurate across the entire domain.	67
Figure 38: Plot of minimum accurate sampling frequency as a function of vertical proximity between chip and sensor level for heat transfer coefficient of $10^5 \text{ W/m}^2\text{-K}$. The inverse solution method is accurate in the shaded region. The direct interpretation technique is inaccurate across the entire domain.	67

Figure 39: Example of a network circuit model of a chip die on a heat spreader attached to a heat sink. This type of model is too computationally intensive for runtime implementation [66].	72
Figure 40: Foster and Cauer RC ladder network representation of thermal system.	75
Figure 41: Network deconvolution using responses in time and frequency.	87
Figure 42: Identification of system poles and semi-infinite limit.	88
Figure 43: Transfer function identification.	89
Figure 44: IIR Response for varying inputted time discretization.	90
Figure 45: Schematic of chip system used for numerical simulation and proposed modeling method.	91
Figure 46: Thermal modeling of chip-spreader geometry shown on Figure 45. Compared are the results of simulations using commercial solver with the output of an IIR filter based on 11 stage network. The power step is at 100 W. The maximum transient errors are less than 0.4% of the steady state response.	92
Figure 47: Comparison between numerical model and IIR digital filter output subject to square-wave input power excitations.	93
Figure 48: Comparison in computational efficiency of different methods for evaluation of convolution integrals. The recursive IIR digital filter is superior to existing convolution methods and achieves best possible scaling due to its constant computation overhead for each time step.	94

NOMENCLATURE

Variables

a	Width of chip, m
b	Length of chip, m
f_r	Sensor spatial frequency in the radial-direction, m^{-1}
f_x	Sensor spatial frequency in the x-direction, m^{-1}
f_y	Sensor spatial frequency in the y-direction, m^{-1}
G_{inv}	Inverse solution transfer function
h	Convective heat transfer coefficient, $\text{W}/\text{m}^2\text{-K}$
k	Thermal conductivity, $\text{W}/\text{m-K}$
N	Number of random heat flux profiles tested
P	Power, W
Q''	Heat flux, W/m^2
R''_{th}	Chip vertical thermal resistance per unit area, $\text{W}/\text{m}^2\text{-K}$
t_0	Thickness of chip, m
t	Time, s
T	Temperature, $^{\circ}\text{C}$
τ	Time constant, s
w_t	System response function in time domain
w_r	System response function in real frequency domain
w_i	System response function in imaginary frequency domain

Subscripts

c	Circuit level
s	Sensor level
l	Low resolution
f	Full resolution
e	Includes sensor error
i	Index in x-direction spatial domain
j	Index in y-direction spatial domain

CHAPTER 1: INTRODUCTION

1.1. Thermal Management Challenges for Next-Generation Integrated Circuits

Thermal management plays a central role in microprocessor reliability and performance. For example, a 10-15°C increase in chip operating temperature can lead to a 50% reduction in device lifetime [1]. Time-dependent dielectric breakdown (TDDB) has an exponential dependence on temperature [2]. Circuit performance is degraded by overheating owing to reduced electron mobility, and the temperature dependence of leakage power causes positive thermal feedback. Addressing these thermal reliability concerns has become a central challenge in the development of next generation microprocessors for both high-performance and mobile applications.

In high-performance applications, chip manufacturers are pushing towards three-dimensional integrated circuit (3D-IC) architectures, imposing unprecedented heat output per surface area and introducing new interfaces and design constraints.

Researchers are pursuing exotic microfluidic heat exchangers (e.g. [3], [4]) to cool these systems, yet these techniques require substantial integration complexity, pose reliability concerns, and occupy valuable regions of the microprocessor. An alternative approach is being pursued using through silicon vias (TSVs) to dissipate heat between stacked levels in the chip [5]. While this technique offers improvements in chip conduction, excessively dense arrays of TSVs would be required to adequately reduce source-to-sink thermal resistance; even so, the problem of dissipating the corresponding heat fluxes at the sink would remain unsolved. Research activities in this field are increasingly intense as it remains unclear which thermal management

solutions can address the challenges of three-dimensional integrated circuit (3D-IC) architectures.

In mobile applications, which are rapidly becoming a dominant player in the consumer electronics sector, cost constraints require novel approaches for minimizing package thermal resistance and accommodating transient thermal excursions while remaining extremely inexpensive and robust to highly variable ambient conditions. Strong consumer demand for compact, attractive form factors further exacerbates the problem. Attention is focused on the chip-to-package thermal resistance as the design space for package-to-ambient thermal resistance is constrained by end user use. Forced convection air cooling techniques are inappropriate due to cost, size, and reliability. Advancements in closed-loop liquid cooling such as vapor chambers offer compelling performance characteristics but are too costly for most consumer applications. Improvements may be achieved by enhancing the thermal conductivity of packaging materials, possibly via nanoscale inclusions. Alternately, wafer thinning and improved thermal interface materials (TIMs) may provide much needed decreases in junction-to-package thermal resistance. Pending unexpected breakthroughs in these fields, however, it appears mobile applications are approaching their power limitations due to thermal constraints, leaving recent improvements in mobile battery technologies underutilized.

Dynamic thermal management (DTM) schemes offer a complementary technique for managing microprocessor hotspots in both high-performance and mobile applications. By dynamically re-routing power on the chip in response to thermal signals, DTM

schemes improve chip performance while ensuring reliable operating conditions.

DTM schemes resolve the challenge of designing for highly variable power profiles, unknown ambient conditions, and degrading package thermal performance.

Variations in chip power profiles result from increased process variability, unpredictable application loads, and variations in ambient temperature conditions [6].

Joule heating models for well-characterized circuits are accurate, but models for circuit leakage power, which is a strong function of operating temperature and process variation, require large uncertainty intervals. Chip power output is necessarily a function of computational tasks which can vary greatly between applications. Power-aware application development has been discussed but remains beyond the planning scope of thermal management solutions. Figure 1 shows a representative range of power variability within the Intel Itanium microprocessor family [6].

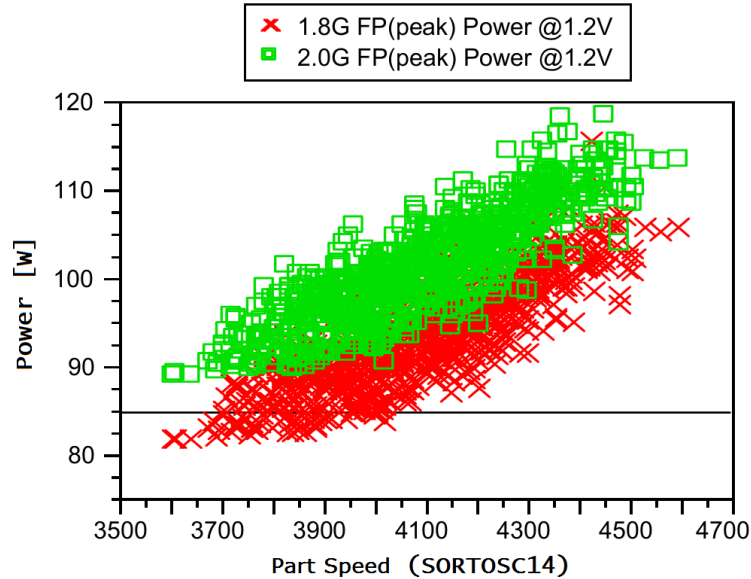
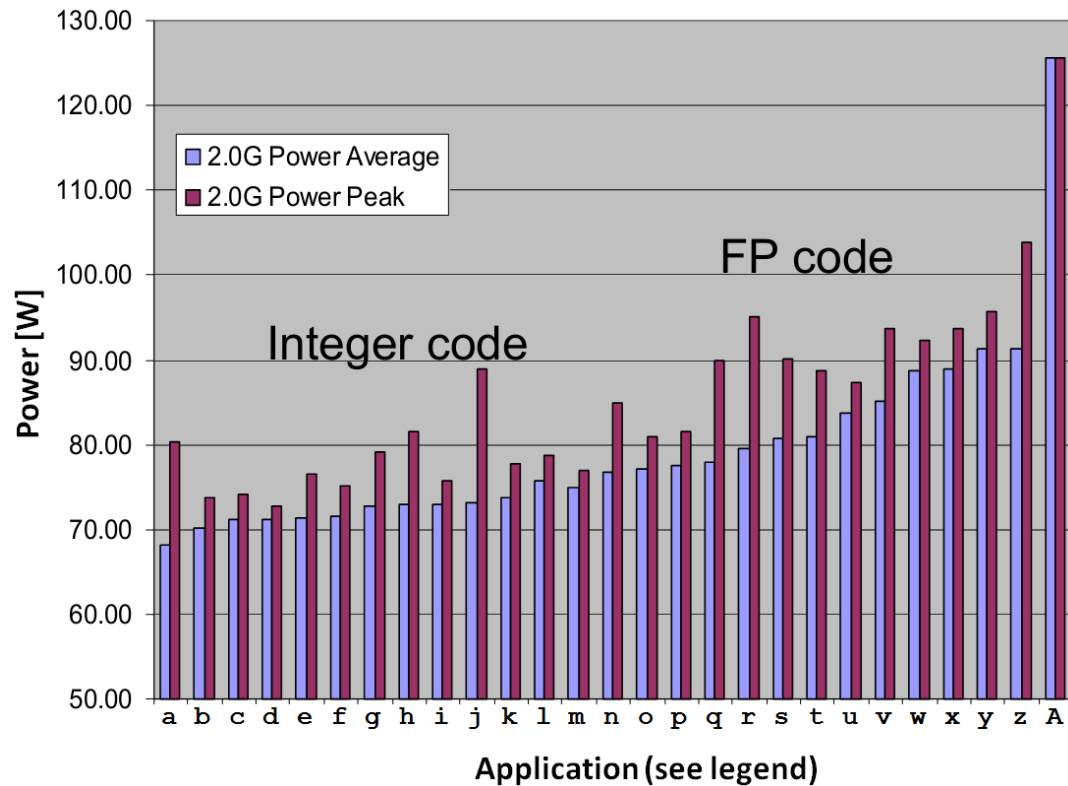


Figure 1: Variability of chip power for Intel Itanium microprocessor family. Figure adapted from [5].



Legend of Application Names							
a	int.181.mcf	h	int.256.bzip2	o	fp.188.ammmp	v	fp.173.applu
b	int.175.vpr	i	int.253.perlbmk	p	fp.187.facerec	w	fp.2000.sixtrack
c	int.186.craft	j	int.176.gcc	q	fp.301.apsi	x	fp.172.mgrid
d	int.300.twolf	k	int.255.vortex	r	fp.191.fma3d	y	fp.171.swim
e	int.254.gap	l	int.252.eon	s	fp.189.lucas	z	fp.178.galgel
f	int.197.parser	m	fp.177.mesa	t	fp.179.art	A	Pathological virus
g	int.164.gzip	n	fp.168.wupwise	u	fp.183.quake		

Figure 2: Variability in power by application for Intel Itanium microprocessor family.
Figure adapted from [5].

Figure 2 shows the broad range of power outputs for various applications. Of particular note, malware such as viruses can cause extreme power excursions; chip thermal management must account for both intended and unintended software applications. Variations in chip thermal properties, which tend to occur over longer

timescales than power variations, can also have a substantial impact on overall thermal management. The thermal interface material (TIM) at the die-lid interface (TIM 1) and the lid-heat sink interface (TIM 2) can vary by chip and can degrade during lifetime use. Chip thermal management must be robust to variations in heat sink performance since the heat sink is typically chosen independently for each application. Heat sink performance can also degrade over the chip lifetime due to fouling and mechanical failure. Thermal management must be robust to variations in heat sink performance since the heat sink is typically chosen independently for each application. Heat sink performance can also degrade over the chip lifetime due to fouling and mechanical failure. There is even evidence of heat sink corrosion caused by contaminated air in datacenters located in regions with very poor air quality.

1.2. Introduction to Dynamic Thermal Management

To address variations in power and thermal properties, the microprocessor industry has widely adopted dynamic thermal management (DTM) schemes to permit more aggressive system utilization. Such schemes require two fundamental components: (1) a temperature measurement or prediction and (2) a controller capable of throttling chip power in response to the temperature.

Intense research has focused on a variety of techniques for throttling chip power, including clock gating [7], Dynamic Frequency Control [8], DVFS [9], SMT thread reduction [10], and activity migration [11]. Combinations of techniques can

be employed to limit the impact on computational performance. For example, McGowen *et al.* [12] noted that the chip core power has the following proportionality:

$$P \propto V^2 F \quad (1)$$

where P is the chip power, V is the voltage, and F is the frequency. By throttling both voltage and frequency, they were able to reduce chip power by 31% with only a 10% reduction in operating frequency.

Compared to the intense research interest in DTM power throttling, relatively little attention has been given to the challenge of obtaining accurate thermal measurements or prediction on the chip while meeting stringent system architecture requirements. Advancements in on-chip thermometry, thermal signal processing, and runtime transient models are critical to minimize temperature uncertainty and maximize reliable chip performance. These topics are the focus of this doctoral research and are explored in detail in the remaining chapters.

1.3. Outline of Doctoral Research

Chapter 2 provides a critical review of research relevant to uncertainty reduction in dynamic thermal management schemes. Several representative multi-core DTM systems are presented to establish a baseline understanding of system integration. Numerous thermal sensor designs are evaluated and considered for their size, process-variation sensitivity, and accuracy. High-resolution laboratory thermography techniques used for sensor calibration and layout are reviewed, and sensor array layout

schemes are explored. Advancement in sensor layout and signal processing are reported with quantitative uncertainties wherever possible.

Chapter 3 presents a technique for determining the accuracy and resolution at which the hotspot heat flux profile can be measured using distributed temperature sensors. Sensor spacing frequencies are varied to represent both embedded temperature sensor and laboratory thermography techniques. The model is based on a novel, computationally-efficient, spatial-frequency domain inverse heat transfer solution. The uncertainty in the calculated heat flux profile is computed for randomized chip heat flux profiles for varying sensor spacing, sensor vertical proximity, sensor error, and chip thermal properties. The inverse solution method decreases mean absolute error in the heat flux profile by more than 30% over a benchmark approach. The results and simulation method can be used to determine the optimal spacing of distributed temperature sensor arrays for hotspot management in chips.

Chapter 4 introduces an original approach for ultra-efficient hotspot prediction for semiconductor devices subjected to arbitrary transient power profiles. The work presented in this chapter was co-authored with Dr. Maxat Touzelbaev; all work was developed by close collaboration between the authors unless otherwise noted. For characterization of the system thermal response, Network Identification Deconvolution (NID) is used; some extensions were made to this characterization methodology by Dr. Touzelbaev and are reported here for completeness. A highly-efficient technique based on digital signal processing is employed to compute the transient thermal response to a power profile. Using an Infinite Impulse Response (IIR)

filtering technique, unnecessary computations can be eliminated while still yielding an accurate response. The technique provides linear scaling of overall computation time with the number of time-steps and dramatically reduces memory requirements. This ultra-efficient algorithm enables the implementation of predictive runtime dynamic thermal management algorithms.

Finally, Chapter 5 offers concluding remarks on the integration of these techniques and opportunities for further research advancements.

CHAPTER 2: THERMOMETRY AND ERROR REDUCTION FOR DYNAMIC THERMAL MANAGEMENT

The development of improved dynamic thermal management schemes requires a thorough understanding of thermal sensors, thermal sensor array design, high-resolution chip thermometry, and thermal signal processing. This chapter presents a review of these topics to identify the state of the art in dynamic thermal management and provide insight into opportunities for further improvements.

In the first section of this chapter, several related review studies are discussed to orient the reader in the literature of dynamic thermal management and thermal sensors. Second, several DTM implementations are presented to provide a broad understanding of system-level design considerations. In the third section, various thermal sensor designs are discussed and key characteristics are examined. Fourth, laboratory thermography techniques used for sensor calibration are examined for temporal and spatial resolution. Finally, techniques for error reduction are examined, including both signal processing and optimized device layout.

2.1. Review Studies Related to Dynamic Thermal Management

In 2004, Blackburn *et al.* [13] conducted a brief review of chip thermometry techniques, including both on-die sensors and laboratory techniques. The review organizes the thermometry techniques into three categories: electrical, optical, and physical contact. The operating principles of each technique are discussed and the spatial resolutions are reported. The review provides an excellent basis for

understanding semiconductor thermometry but does not address temperature error which is critical for DTM applications.

Avenas and Dupont [14] presented a brief review of related thermometry techniques but focused on techniques for power semiconductor devices. Typical non-contact methods used in the power-electronics community are outlined. A series of electrical measurements are proposed and the following characteristics are compared: sensitivity, linearity, accuracy, genericity, and calibration requirements. The authors also discuss the feasibility of characterizing the thermal impedance or temperature during operation in the presence of self-heating. Implications for wide band-gap semiconductors are discussed. While the thermometry techniques proposed for power semiconductor devices are not directly relevant for dynamic thermal management in microprocessors, the authors provide an excellent framework for considering microprocessor thermometry techniques.

In 2007, Naderlinger [15] conducted a brief review of dynamic thermal management techniques (DTM) and outlined three techniques for estimating chip power profile: instruction-level estimation, function-level and macro-modeling estimation, and event counters. Instruction-level estimation uses an energy cost factor for processor instructions to determine overall power consumption. However, instruction energy costs depend on the sequence of instructions, so more detailed calculations are required. Function-level estimation leverages power simulations at the function level to estimate total power. Finally, event counters have been shown to reflect power profiles and can thus be used as an estimation parameter. Naderlinger also mentions

several modeling tools available for thermal simulations. The explanation presented for task migration, however, is incorrect in that it claims silicon is a poor thermal conductor. The thermal conductivity of silicon is amongst the highest of common engineering materials at 148 W/m-K at room temperature [16]. Microprocessors suffer reduced thermal conductivity owing to thermal interfaces and the inclusion of insulating materials. Overall, the review offers a valuable introduction into advances in DTM and power estimation.

Kong *et al.* [17] conducted a review that has yet to be published of various thermal management techniques for chips. The review focuses on techniques related to the chip microarchitecture. The study was organized in 6 categories: “temperature monitoring, microarchitectural techniques, floorplanning, OS/compiler techniques, liquid cooling techniques, and thermal reliability/security”. The review covers thermal sensor types and sensor placement techniques.

2.2. Representative Realizations of Dynamic Thermal Management

Modern microprocessors contain numerous thermal sensors for three main reasons. First, chip dies contain numerous distinct units for which the temperature must be monitored. Second, chip workload can vary greatly resulting in migrating hotspots. Finally, as previously discussed, leakage power is difficult to predict due to both processing-dependencies and temperature-dependencies. [18]

Intel’s “Foxton Technology” which was introduced in their 90-nm Itanium processor family provides a good example of an implemented DTM system. The system

involved four on-chip sensors and an embedded micro-controller to measure and regulate the power and temperature of the system. One temperature sensor was located above each of two floating point units and one sensor was located above each of two cores. The thermal transients in the system were reported to be 60°C per second. A sample rate of 500 Hz was chosen to facilitate response to these gradients without suffering excessive temperature excursion. A micro-controller was used to control system response by throttling chip power once the set threshold temperature was met. Power throttling was achieved by voltage reduction at a rate the authors described as “appropriate for thermal time constants”. The study does not report any power throttling as the system approached the threshold temperature. Figure 3 presents a layout of the microprocessor layout. Figure 4 presents a schematic of the control design.

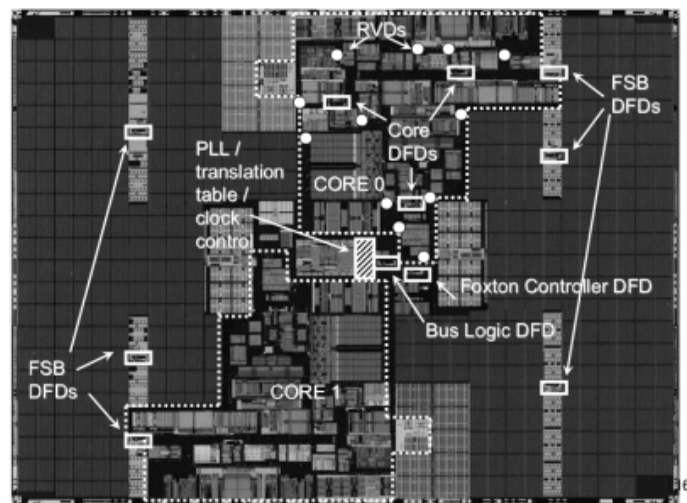


Figure 3: Layout of Foxton Technology Controller used for Intel Itanium microprocessor family [6].

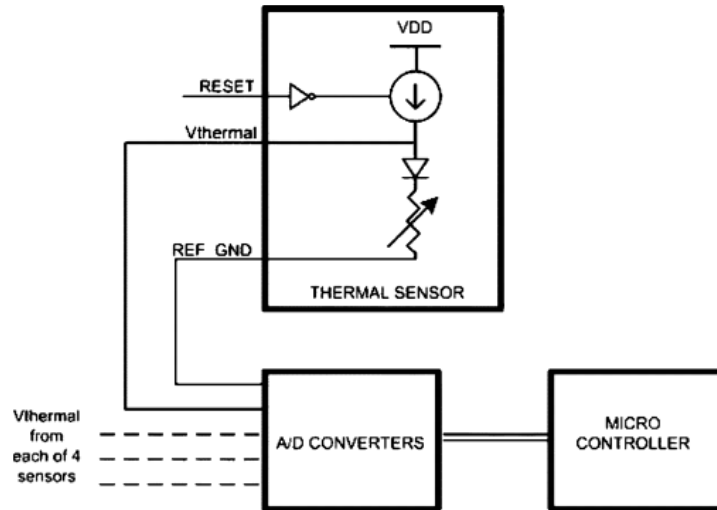


Figure 4: Thermal system block diagram for Intel Foxton Technology Controller [12].

The sensors in the Foxton Technology Controller are calibrated after fabrication and are re-calibrated every 65.536 ms during microprocessor operation. The system requires less than 0.5% of the total die area and consumes less than a watt of power. The average power measurement accuracy is 5% and the temperature accuracy is 3°C.

For an example of an alternative, more modern layout, Dorsey *et al.* [19] provides details of the AMD quad-core processor layout which extends over relatively long distances. Figure 5 shows the processor floorplan including thermal diode locations. The components marked “ThermCenter” represent central hubs for signal processing. The average distance between the sensors and the thermal evaluation circuit, referred to as “TCEN”, is approximately 8mm. The chip has 38 total sensors. The power consumption for the plurality of thermal sensors is estimated to be as high as 10’s of watts.

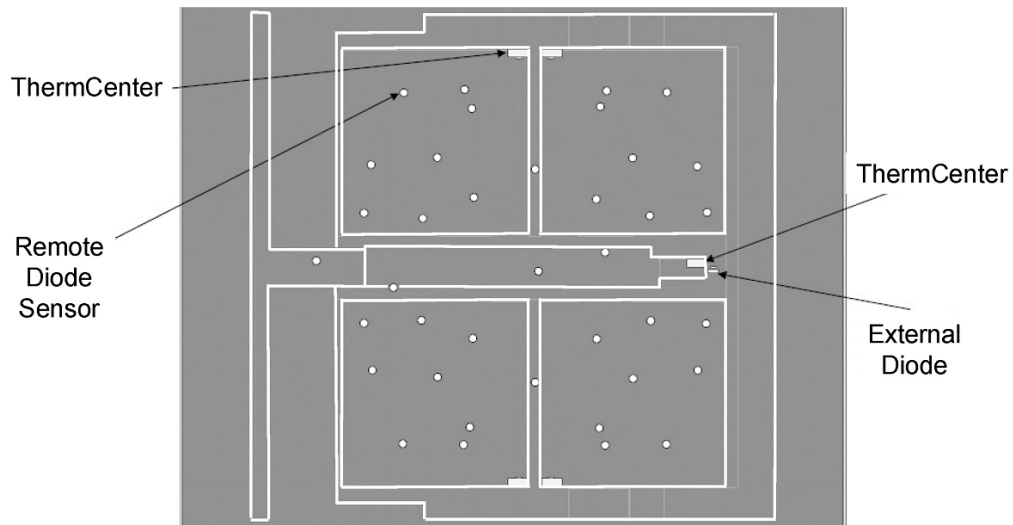


Figure 5: Floorplan for AMD Quadcore Opteron processor showing thermal diodes and thermal processing centers [19].

Implementation of dynamic thermal management has increased significantly over the last decade. In 2002, Krinitsin provided a detailed analysis of dynamic thermal management schemes on AMD Athlon XP and Intel Pentium 4 microprocessors. Typical transient gradients in temperature were reported as 30-50°C/sec for normal operation but could be as high as 70-100°C/sec in the event of cooling system failure. Conventional thermal monitoring at the time involved sampling rates of approximately 5-10 Hz, which was inadequate for response to these transients. In the case of the Intel Pentium 4 processor, a thermal sensor was placed directly over the rapid-integer arithmetic logic unit (ALU) [20].

For the modern, high-performance microprocessors, dynamic thermal management involves off-chip communication for heat sink control. Details of the chip dynamic thermal management scheme are presented in datasheets to enable system integration with heat sinks. One such report for the Quad-core Intel Xeon 5400 series

microprocessor family provides an excellent example of modern industry practices. The chip's four cores are equally divided between two "domains". A Platform Environmental Control Interface (PECI) reports the highest output temperature within each chip domain to an off-chip PECI host, typically to control fan speed. Thermal sensors do not report temperature directly, but instead report PECI counts; each PECI count is equivalent to "approximately 1°C" though linearity "cannot be guaranteed" past 20-30 PECI counts. PECI counts are expressed as negative values relative to the set temperature; positive PECI counts are not reported, which causes a complete signal loss to the PECI host when the set temperature is exceeded. A single control temperature is set for each domain and is not absolute but instead defined relative to the "TCC activation point", which is the starting value of the temperature control circuit (TCC). [21] Further documentation provides details of the three thermal management signals that can be activated on-chip. PROCHOT# is activated for a domain when any temperature sensor in the domain reaches its factory configured trip point which activates the Thermal Control Circuit (TCC) to reduce power output. The chip platform can activate the TCC for all cores by asserting FORCEPR# signal. To prevent damage in the event of a cooling system failure, THERMTRIP# activates complete system shut down independent of processor activity if a temperature is reached that may cause "permanent silicon damage". The signal is tripped within 10 microseconds. The processor core voltage must also be removed. These features enable the chip to operate reliably for cooling solutions capable of removing the "thermal design power" (TDP) specification. The TDP is not the maximum power output of the chip, which may exceed the TDP. For the highest performance Xeon

5400 processor, the TDP is 150 W. According to the report, the dynamic thermal management scheme described above is key to facilitating such aggressive power consumption in this processor family. [22]

Related work on other types of semiconductors can provide insight into the challenges of implementing thermal sensors. Lopez-Buedo *et al.* [23] and Velusamy *et al.* [24] both employed digital thermal sensors on field-programmable gate arrays (FPGAs), and Mondal *et al.* [25] developed a technique for inserting sensors into FPGAs. Mukherjee *et al.* [26] developed an algorithm to use vacant configurable logic blocks as thermal sensors in an FPGA design. Aldrete-Vidrio *et al.* [27] presented two approaches to conduct failure analysis on RF circuits with built-in differential temperature measurements. Finally, Bratek and Kos [28] combined power modules and temperature sensors to detect faults in integrated circuits.

Overall, present day integrated circuit (IC) thermal sensors are accurate though costly; furthermore typical calibration is time-consuming and adds further costs.[29] Sensor calibration is typically conducted by heating the chip and measuring the sensor output [12], [30], [31]. Heating is usually applied by hotplate, high-temperature soak, or hot air jet [29]; further details on industry sensor calibration practices can be found in Schlaepfer [9] and McGowen *et al.* [27].

2.3. Thermal Sensor Design

A range of thermal sensors are has been developed for on-chip thermometry. Analog thermal sensors consist of a thermal diode, a factory-calibrated reference current

source, and a current comparator. When a voltage is applied across the diode, the induced current flow is a function of temperature. Comparison between the diode current and the reference current yields a thermal signal. The accuracy of thermal diodes is limited by the fact that the threshold current varies strongly with processing parameters. For this reason, each thermal diode sensor must be calibrated accordingly. [29] If sensor calibration is not conducted, the resulting error can be very large. The thermal assist unit (TAU) developed for the IBM PowerPC 750 microprocessor produced worst-case readings between 61°C and 109°C for cases where the actual temperature was 95°C . [30]

More sophisticated digital sensors are available but accuracy and stability tends to require additional size. Figure 6 presents a schematic of a digital thermal sensor which integrates a thermal diode and a reference current source. Kaxiras and Xekalakis [32] proposed a 4T-decay sensor design which uses 4T memory cell with a decay counter.

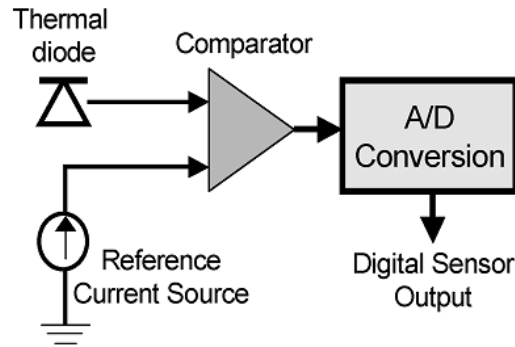


Figure 6: Schematic of digital thermal sensor based on a thermal diode [33].

To achieve $\pm 1^{\circ}\text{C}$ accuracy, the sensor required area of 0.0016 mm^2 and power consumption of 397uW , according to calculations by Long *et al.* [33]. Digital thermal

sensors were used in the IBM Power6 and the Intel Xeon series; in the IBM Power6, ring oscillator sensors were used [29]. The primary parameter for assessing thermal diode sensors is the diode ideality factor which quantifies the deviation from an ideal diode [33].

Aldrete-Vidrio *et al.* [34] presented four differential temperature sensor designs, two of which were active and two of which were passive. Differential temperature sensors are designed to have high sensitivity to local temperature disturbances within the silicon die, but low sensitivity to external temperature variations. The passive sensors were integrated thermopiles. A thermopile is a string of thermocouples connected in-series. The active sensors were differential amplifiers. Lateral parasitic bipolar transistors were used as the temperature transducer devices. The authors defined their figures of merit as “compatibility with [integrated circuit (IC)] technology, used area, power consumption, sensitivity, and linearity”. The authors did not report sensor uncertainty results but did report evidence of sensitivity and linearity, as shown in Figure 7.

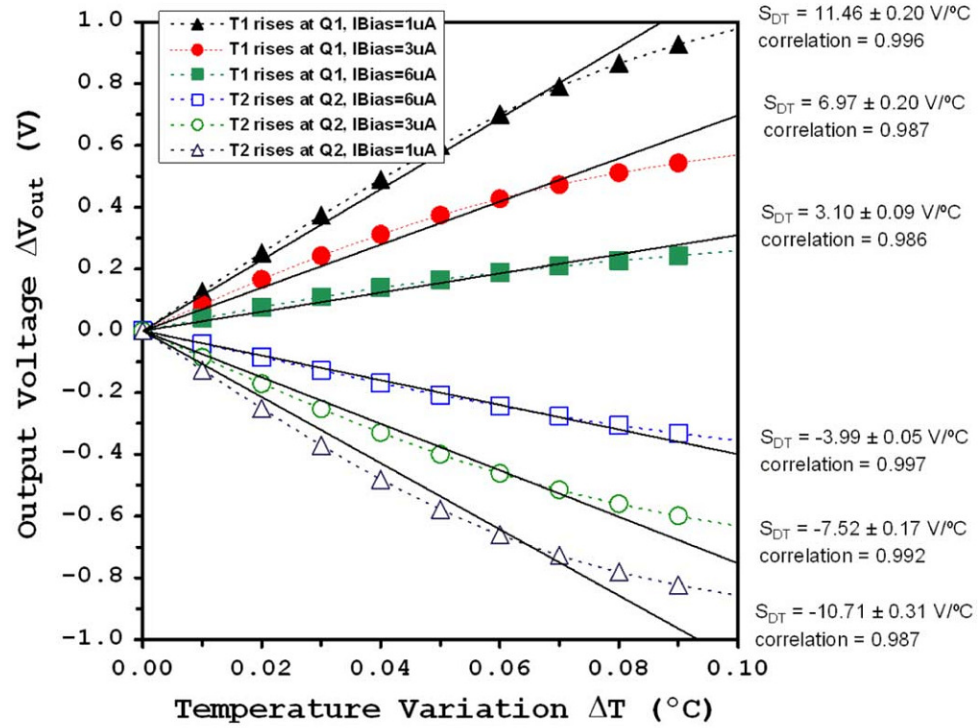


Figure 7: Representative results of linearity and sensitivity of sensors designed proposed by Aldrete-Vidrio *et al.* Solid lines indicate least-squares fit to linear approximation and corresponding parameters are shown on the right. [34]

As before, insights for new thermal sensors for microprocessors may be taken from parallel work on other electronics systems. Barlini *et al.* [35] presented three electrical-based techniques for finding transient average junction temperature in power MOS devices. Two techniques use the temperature-dependence of time-derivative of the drain-source current, $dI_{ds}/dt(T)$, and the other technique uses the temperature dependency of turn-ON delay of device. Simulations of the techniques in PSPICE are compared to two thermal models. The measurement techniques based on time-derivative of the drain-source current demonstrate accuracy within 5°C , as does the technique based on turn-ON delay except at very short timescales. [35]

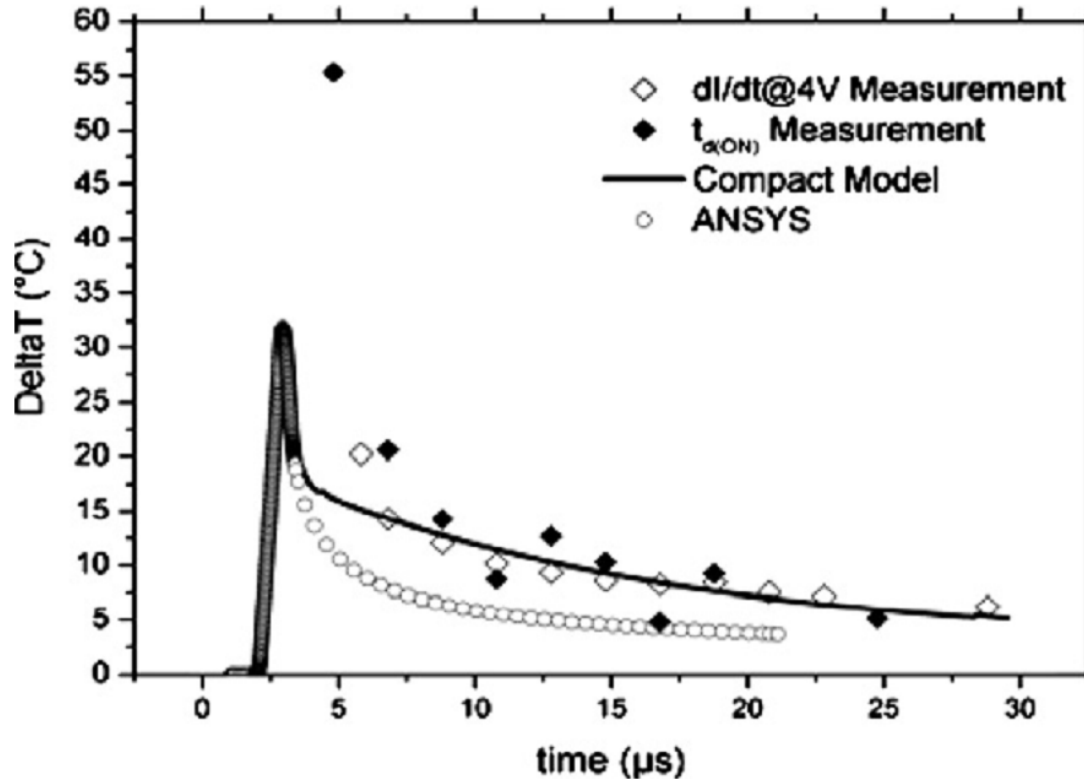


Figure 8: Comparison of proposed thermal measurement technique with simulation tools [35].

Sensors designs that are susceptible to variations in processing parameters, such as doping, require separate calibration for each sensor [29], [33]. Sensor designs have been developed that are process-invariant, but they tend to require large areas [36], [37]. To increase the accuracy of a sensor, the series resistance, which is the resistance of paths leading up to and away from the diode, can be corrected for [33].

Yao *et al.* [29] presented a technique for calibrating on-chip thermal sensors using a local array of five sensors. The temperature drop between neighboring thermal sensors provides a basis for improved thermal measurement, provided an accurate power estimate and thermal conductivity value can be obtained. If a grid of thermal sensors is

available, the redundant sensors in the local sensor array can be eliminated, though error would likely increase. The proposed calibration technique is shown to dramatically reduce sensor error, but overall temperature uncertainty remains greater than 10°C for the majority of cases. Figure 9 shows the error reduction achieved via calibration for a representative case. Figure 10 shows the dependence of temperature error on power estimation. The errors reported remain unacceptably high for dynamic thermal management. [29]

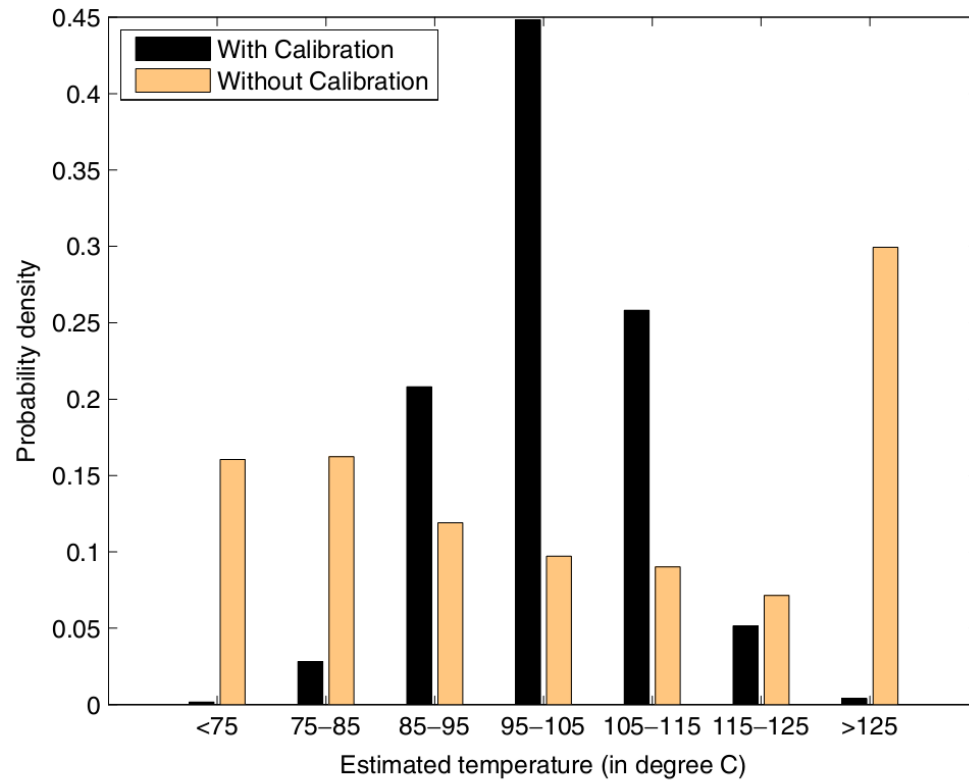


Figure 9: Improvement in thermal measurement due to sensor group calibration [29].

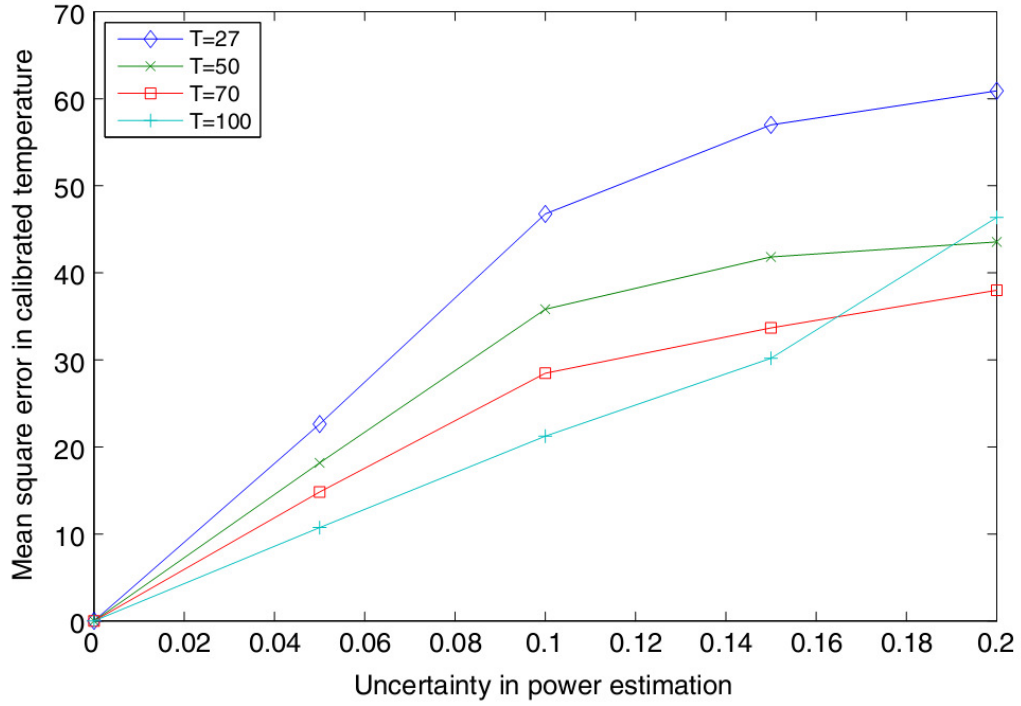


Figure 10: Temperature estimation error as a function of power estimation error [29].

Bharath *et al.* [38] extended the work of Yao *et al.* [29] by developing an alternative technique for calibrating on-chip thermal sensors. The authors noted that the technique presented by Yao *et al.* [29] “does not account for the cyclic dependency between leakage currents and temperature. This dependency comes from the fact that: (i) leakage currents increase with device temperatures and (ii) higher leakage currents result in higher power consumption which in turn increases the device temperatures” [38]. To account for these dependencies, two power inputs are used and the ratio between them is calculated. CMOS chip power comes from dynamic power consumption (P_d) and leakage power consumption (P_l) defined by:

$$P_d = CV^2f \quad (2)$$

$$P_l = V \cdot I(V, T) \quad (3)$$

where C is the capacitive load, V is the supply voltage, f is the frequency, and $I(V, T)$ is the leakage current which is a function of supply voltage and operating temperature. The technique is tested by creating a simulated ring oscillator digital thermal sensor in SPICE and coupling to HotSpot code. [23]

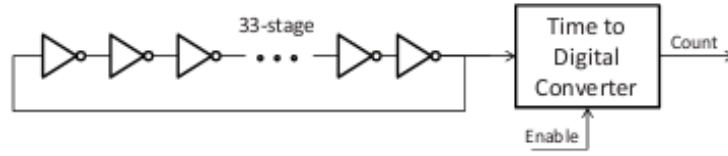


Figure 11: Schematic of simulated ring oscillator digital thermal sensor [38].

The technique was tested for three corner cases, referred to as “slow”, “nominal”, and “fast”. The slow corner had higher threshold voltage and thus lower leakage power, while the fast corner has low threshold voltage and therefore high leakage power. Figure 12 shows the results for the calibrated temperature results for the three corner cases as compared to the results of Yao *et al.* [29]. The techniques are comparably accurate for low leakage power but the technique presented by Bharath *et al.* [38] shows dramatic improvements for the case of elevated leakage power. Despite these improvements in accuracy, it remains unclear if the benefits of this sensor design justify the large sensor footprint requirements.

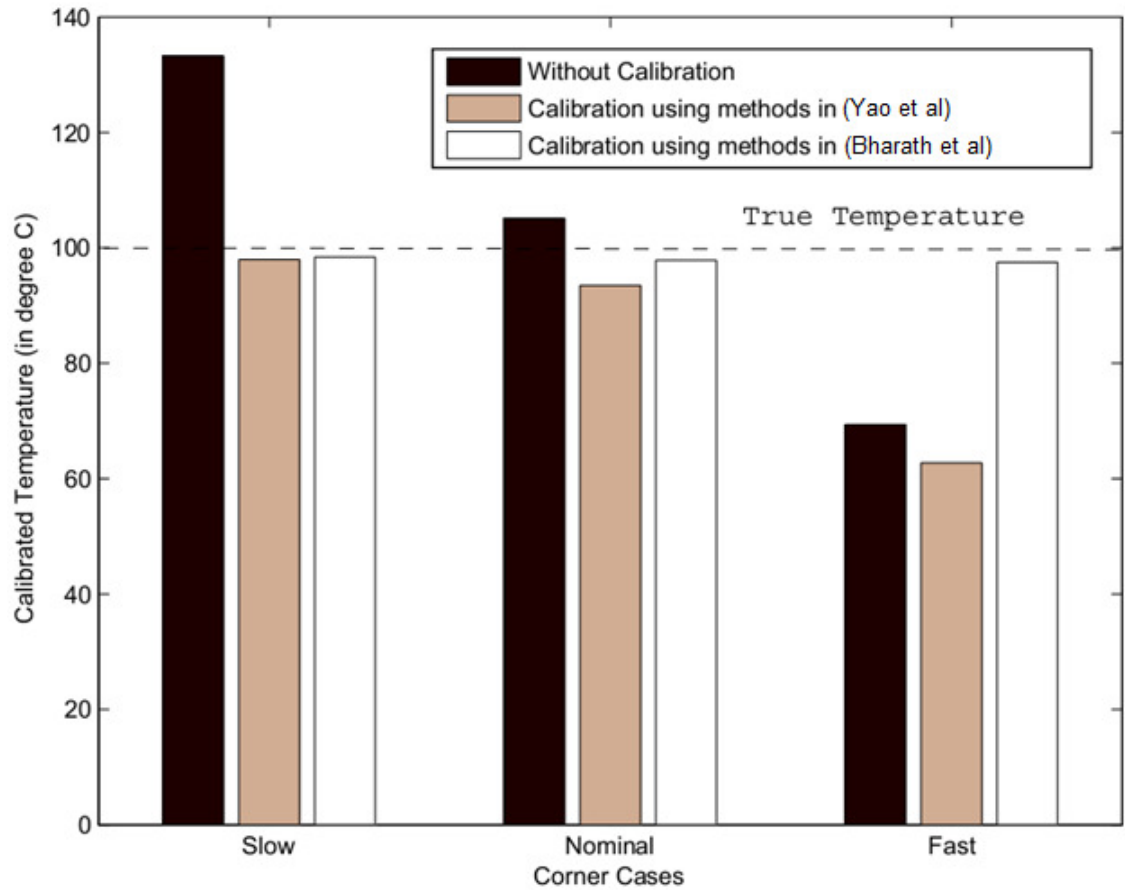


Figure 12: Calibrated sensor response for three corner cases [38].

2.4. Laboratory Thermography Techniques

While on-chip thermal sensors are central to any DTM design, laboratory thermography techniques play an important role in circuit design validation and sensor calibration. Extensive research in high-resolution semiconductor thermography has resulted in numerous available techniques. This section reviews the quantitative spatial resolution, temporal resolution, and accuracy of two thermography techniques capable of providing chip-scale thermography and sensor validation: Raman and infrared (IR) thermography.

Infrared (IR) thermography is one of the most common techniques for chip thermal mapping; the emissivity of the surface of interest is measured by measuring IR radiation counts at a known temperature. After this calibration step, IR radiation from the surface is measured to determine surface temperature. Establishing a uniform, high emissivity surface is beneficial for maximizing the accuracy of IR thermometry.

Salem *et al.* [39] reported the results of various spray coating techniques for IR thermography of power electronics. A FLIR Systems Thermacam SC500 infrared (IR) camera system was used to image spray coatings of high-temperature black spray paint and a boron-nitride spray. The boron-nitride spray is substantially easier to remove than spray paint which makes it appropriate for non-destructive testing. The imaging accuracy is reported to be $\pm 2^{\circ}\text{C}$ based on manufacturer specifications. IR temperature measurements were compared to a thermocouple on the top surface which had an accuracy of $\pm 1.2^{\circ}\text{C}$. An aluminum block was used as a blackbody reference target and additional measurements were demonstrated on an active MOSFET. The measurement error was as high as 30% when the clamping apparatus was within the field of view. Figure 13 shows relatively minor differences in measurement error between the two coatings which depend on the operating temperature.

An important benefit of IR thermometry is the opportunity to integrate electronics cooling while measuring temperature. Hom *et al.* [40] used a microfluidic heat sink to conduct IR thermometry on a chip running realistic traffic patterns. To ensure the heat sink was transparent to the infrared radiation, an IR transparent working fluid was used and a sapphire window was integrated into the top surface of the heat sink. This

approach has also been pursued in at least one research and development lab at a major chip manufacturer.

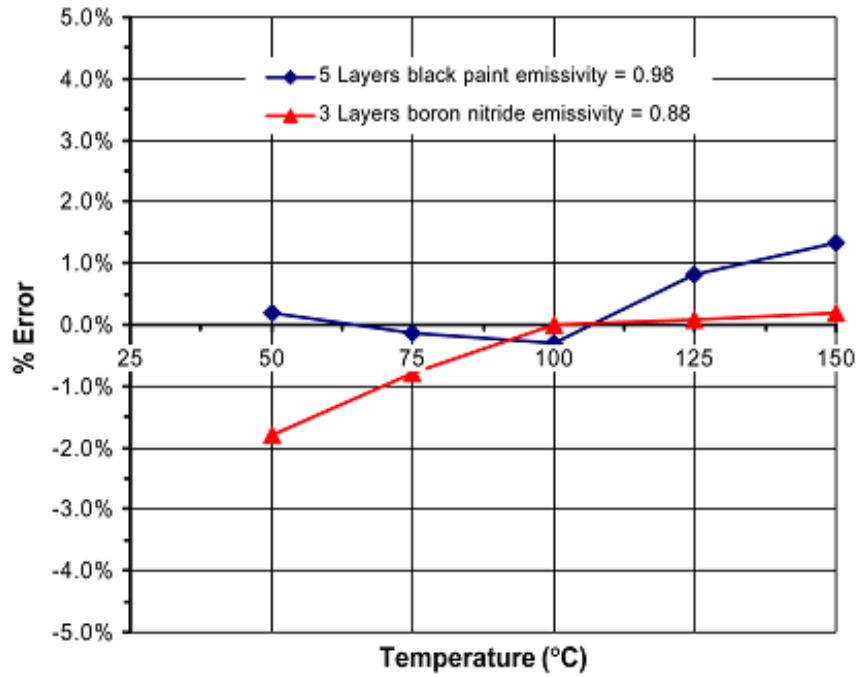


Figure 13: Error in IR thermography for alternative chip coatings [39].

Castellazzi *et al.* [41] demonstrated a novel, transient IR thermal characterization technique using an optical fiber and a coated component. Figure 14 shows a schematic of the optical fiber and corresponding assembly. Figure 15 shows the circuit schematic for signal capture. The system is capable of fast transient operation from microseconds up to milliseconds for semiconductor devices.

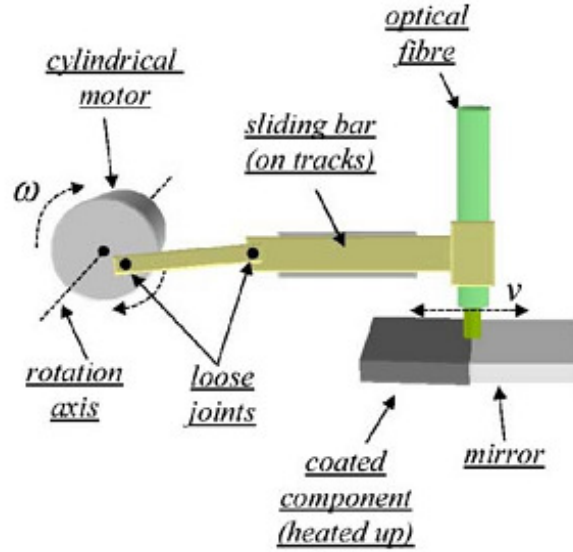


Figure 14: Schematic of novel IR thermography technique [41].

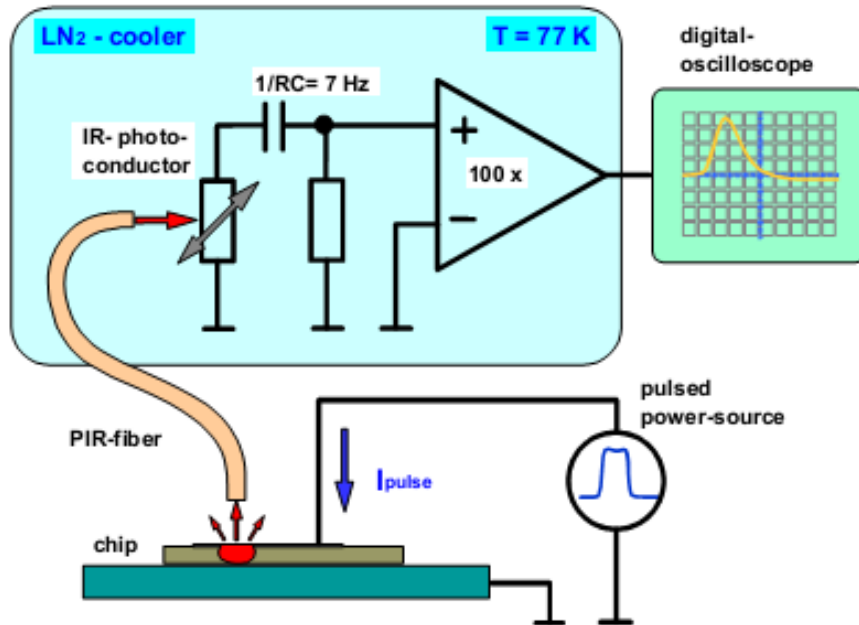


Figure 15: Schematic of fast transient IR thermography technique [41].

Because the spatial resolution of IR techniques is diffraction limited, researchers have turned to Raman thermometry for high resolution chip thermometry. Kuball *et al.* [42] demonstrated the use of Raman thermography on gallium nitride (GaN)

microelectronics and achieved sub-micron spatial resolution and nanosecond time resolution. The results highlighted an important consideration for chip thermography: the spatial resolution of the technique dictates the accuracy at which the hotspot maximum temperature can be measured. For an aluminum gallium nitride/gallium nitride (AlGaN/GaN) heterostructure field effect transistor (HFET), the Raman technique indicated higher peak temperatures than observed by IR thermography. The lower spatial resolution of the IR technique results in averaging the peak temperature with lower temperatures nearby. Figure 16 offers a comparison of the two techniques to demonstrate this effect; spatial averaging alone, however, does not account for the difference in peak temperature measurements so other sources of error must be present. Similar results were found when using micro-Raman spectroscopy to measure multi-finger gallium arsenide (GaAs) pseudomorphic high electron mobility transistor (HEMT) devices [43].

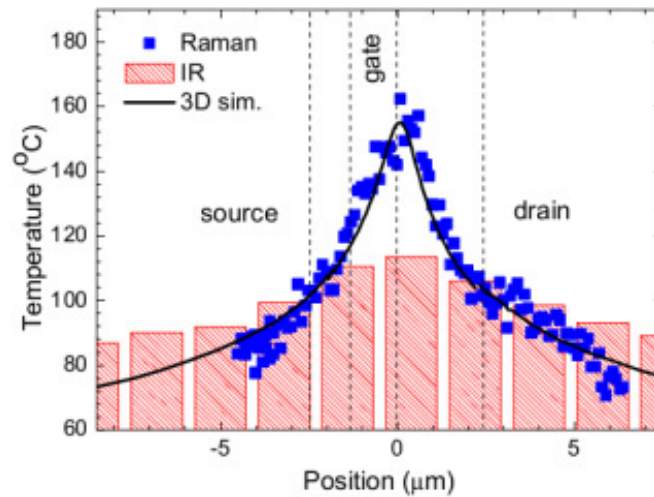


Figure 16: Temperature profile of HFET as measured by Raman spectroscopy and IR thermography. The superior spatial resolution of Raman spectroscopy enables improved peak temperature measurement. [42]

While infrared (IR) thermography has a maximum spatial resolution of 2-5.5 μm [44–46] and a response time as fast as milliseconds [47], [48], Raman is accurate to micron or submicron spatial resolution [49–54] but is limited to materials with appropriate phonon characteristics and therefore has limited applicability for metals or plastics. Unlike IR thermography in which a two-dimensional (2D) thermal map is captured by an imaging sensor, Raman requires rastering to produce a thermal map.

A third thermometry technique known as transient interferometric mapping (TIM) is used in semiconductor applications. The technique uses temperature-induced changes in silicon refractive index to determine the chip temperature. Typically transient interferometric mapping is applied to the chip backside. Bychikhin *et al.* [55] demonstrated the use of transient interferometric mapping to measure temperature distribution in DMOS devices subjected to repetitive stress. Results were compared to measurements from built-in temperature sensors and 3D thermal simulations and showed good agreement. A two-dimensional holographic interferometry technique with 10 μs time resolution and a scanning heterodyne interferometer with 3 ns time resolution were used for the thermal mapping. Heer *et al.* [56] developed an automated experimental set-up that used transient interferometric mapping and IR mapping to measure DMOS degradation mechanisms. Haberfehlner *et al.* [57] employed a compact transient interferometric mapping system to spatially resolve thermal runaway onset in a smart power DMOS. The technique takes two thermal images with a delay ranging from 100 μs to a few milliseconds. The phase measurements are made using superluminescent diodes and focal plane array cameras. The field of view ranged

between 250um by 300um to 2.5mm x 3mm. Figure 17 shows representative phase shift results corresponding to a shift in temperature measured by a built-in sensor. Figure 18 provides a schematic of the transient interferometric mapping set-up used. Blaho *et al.* [58] also employed backside laser interferometric thermal mapping technique on a double-diffused metal–oxide–semiconductor (DMOS) transistor.

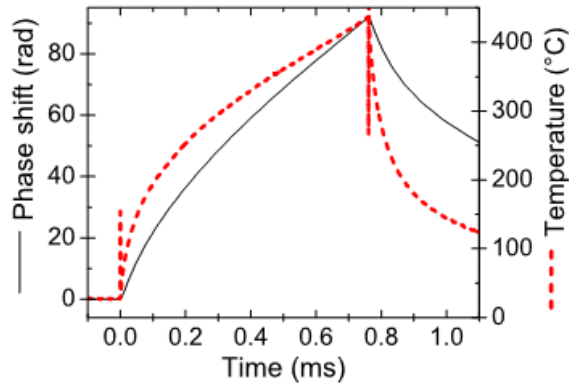


Figure 17: Temperature response measured by built-in sensor and corresponding phase shift response for transient interferometric mapping [57].

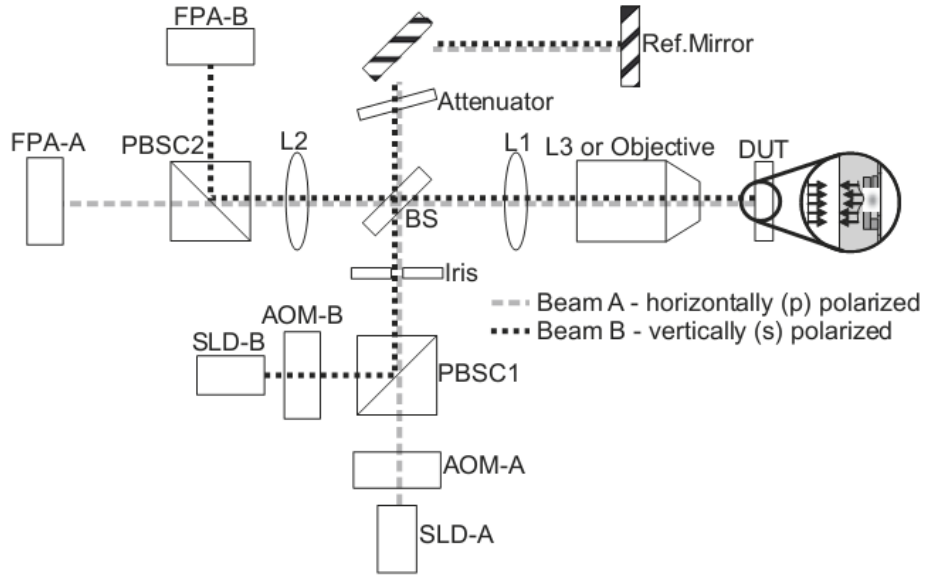


Figure 18: Schematic of transient interferometric mapping set-up. Abbreviations are: SLD: superluminescent diode; AOM: acousto-optical modulator; PBSC: polarizing beam splitter cube; BS: beam splitter; L: lens; DUT: device under test; FPA: focal plane array camera [57].

2.5. Thermal Sensor Arrays and Signal Processing

Regardless of sensor type and accuracy, on-chip thermal measurements are significantly affected by sensor placement. Improvements in thermal measurements can be made by optimizing sensor placement within the constraints imposed by circuit architecture and system integration. Optimized sensor placement can reduce error associated with hotspot migration and can improve sensor network requirements.

Hotspots on multicore processors migrate during chip operation, resulting in differences between the maximum measured temperature on the chip and the actual maximum chip temperature. Deviations as large as 12.6°C can occur even with 16 sensors per core. [60] Long *et al.* [33] also noted that measurement inaccuracy is

introduced when thermal sensors are located away from the hotspot, but did not attempt to directly quantify this error. The authors referred to the practice of correcting for these deviations as “remote sensing”. The practice of remote sensing is valuable both for space allocation and to reduce temperature-dependencies in sensor signal lines.

Lee *et al.* [61] demonstrated an analytical model for determining maximum temperature drop between a hotspot and a region on the chip. The model is a superposition of the exponential temperature decays of each hotspot. An activity factor is used to scale the contribution from each hotspot for power outputs between zero and maximum power. The authors implement a run-time thermal model which facilitates the use of “virtual sensors” and examine two benchmark cases of concentrated thermal stress.

Gunther *et al.* [62] suggested opportunities for optimized sensor placement but did not present concrete methods. Mukherjee and Memik [26] developed a thermal sensor allocation scheme for single-core microprocessors. This scheme was used by Long *et al.* [59] as a basis for a multi-core sensor allocation and sampling strategy. The authors proposed three techniques to create sensor infrastructures to improve hotspot monitoring in a multi-core system. First, they propose an improved interpolation scheme. The technique uses a low-resolution sensor array to estimate the location of a hotspot. As shown in Figure 19, localized, higher-resolution sensor array is then activated in the vicinity of the hotspot to determine its precise location and magnitude.

This technique dramatically reduces signaling requirements but still requires a high-resolution array of sensors throughout the chip.

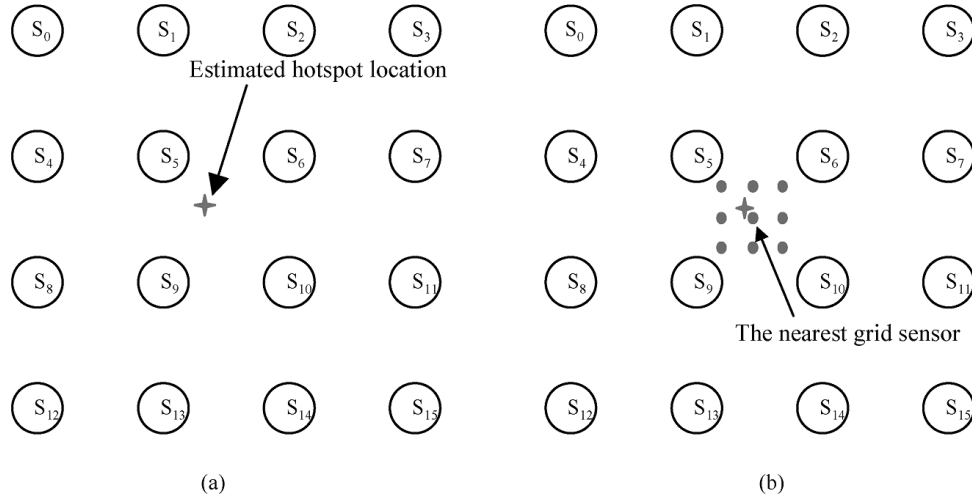


Figure 19: Schematic of two-stage hotspot interpolation technique proposed by Long *et al.* [59]. A low-resolution sensor array (a) is used to estimate the hotspot location and then a local high-resolution array (b) is activated in the vicinity of the hotspot.

Several studies have taken advantage of the way in which heat dissipation in microprocessors provides smooth temperature gradients. This results in reductions in the thermal signal bandwidth in the spatial frequency domain. Cochran and Reda [18] presented improved thermal sensor array interpolation techniques. The authors prepared several approaches based on Fourier spectral analysis and Nyquist-Shannon sampling theory. The proposed methods could tolerate both uniform and non-uniform sensor placements. The error associated with the full thermal profile and the hotspot temperature were reported and compared to a nearest neighbor benchmark approach. For the best case, the average absolute error was 0.6% and the overall performance was deemed to exceed that of grid-based interpolation [59] and geostatistical Kriging estimators [63].

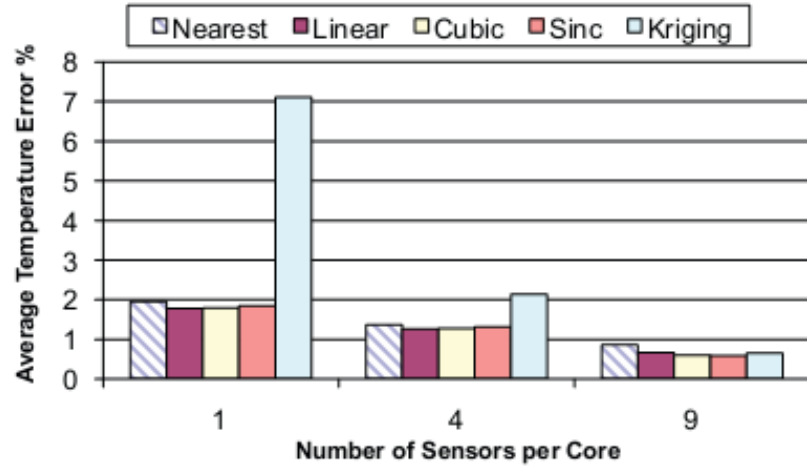


Figure 20: Error results for full thermal characterization using various thermal sensor array interpolation schemes [18].

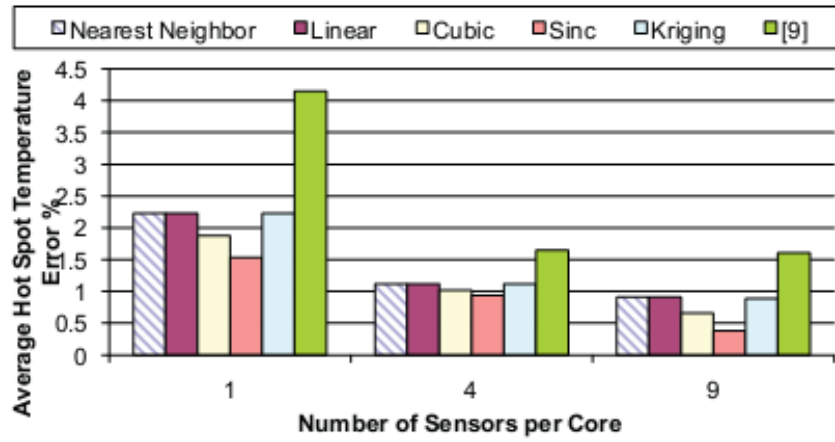


Figure 21: Error results for hotspot estimation using various thermal sensor array interpolation schemes [18].

Sharif and Rosing [64] presented two techniques for accurate on-chip temperature sensing as shown in Figure 22. Their first technique was developed for sensor calibration and sensor allocation during chip layout design. This technique was shown to reduce the number of sensors needed for a particular accuracy level by 16% on average. The authors also proposed a technique to determine the temperature at

arbitrary locations on the die using noisy temperature readings from sensors located elsewhere on the chip. The technique based on a Kalman filter is designed to operate in runtime and is shown to reduce the standard deviation and maximum value of temperature error by an order of magnitude. [64]

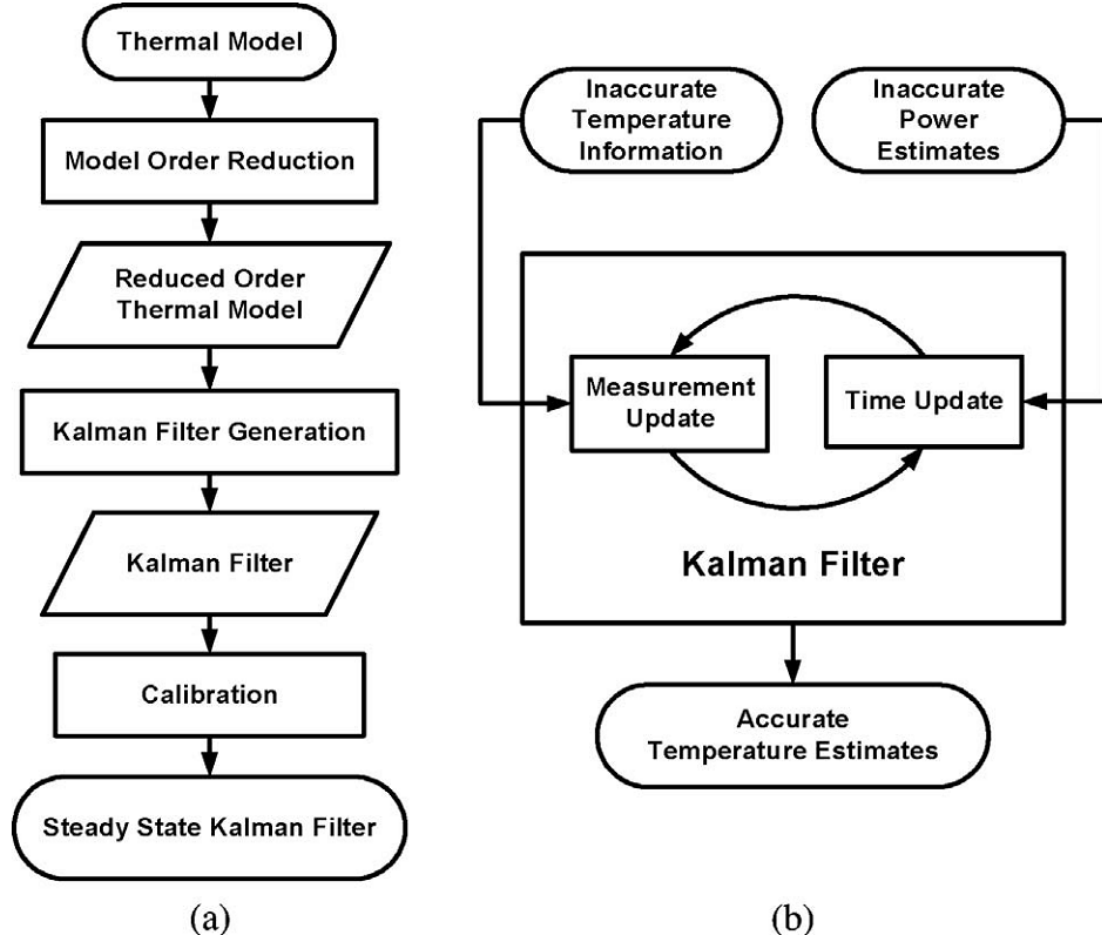


Figure 22: Schematic of two techniques proposed for (a) off-line sensor set-up and (b) runtime temperature estimation [64].

All of the aforementioned techniques require monitoring of signals from distributed thermal sensors, introducing possible signal noise and imposing additional power consumption. Ituero *et al.*[60] developed improvements to sensor monitoring that

reduce the power consumption and improve the accuracy of the array. Figure 23 provides a schematic of interconnect lines required for monitoring a single microprocessor core, shown overlaid on the full microprocessor in Figure 24. The authors proposed a new monitoring network paradigm for communicating between sensors and the controller, and demonstrated the technique for both a single-core and an eight-core system.

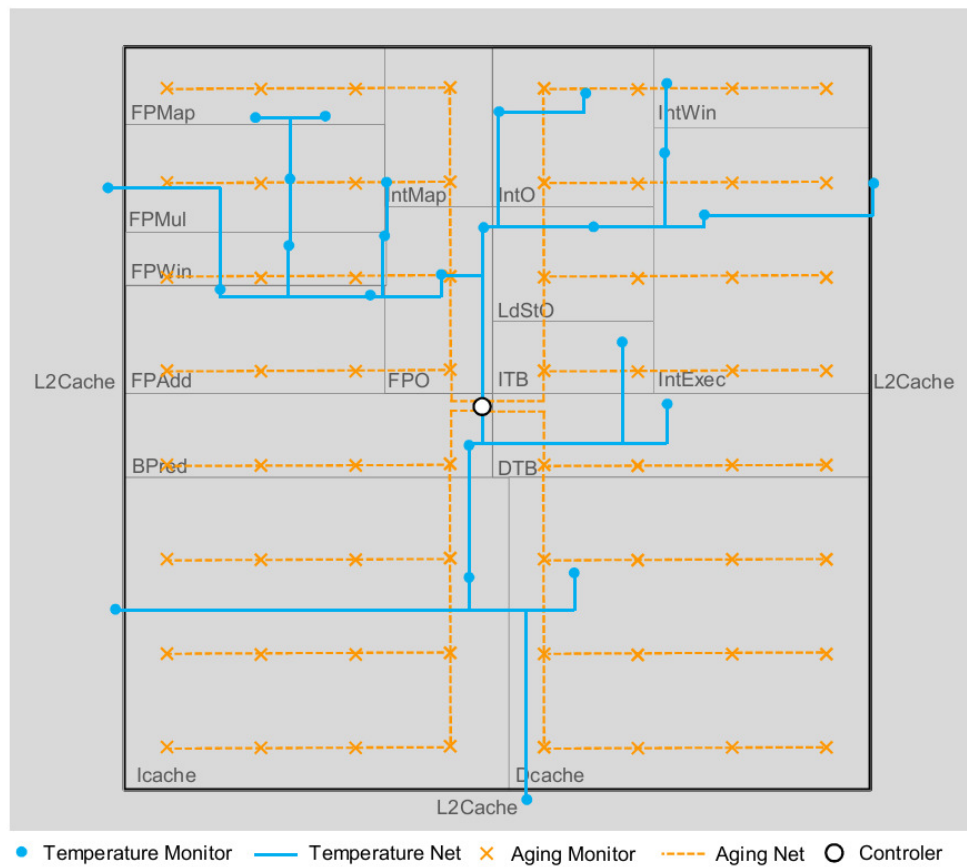


Figure 23: Schematic of interconnect lines required for monitoring a single core [60].

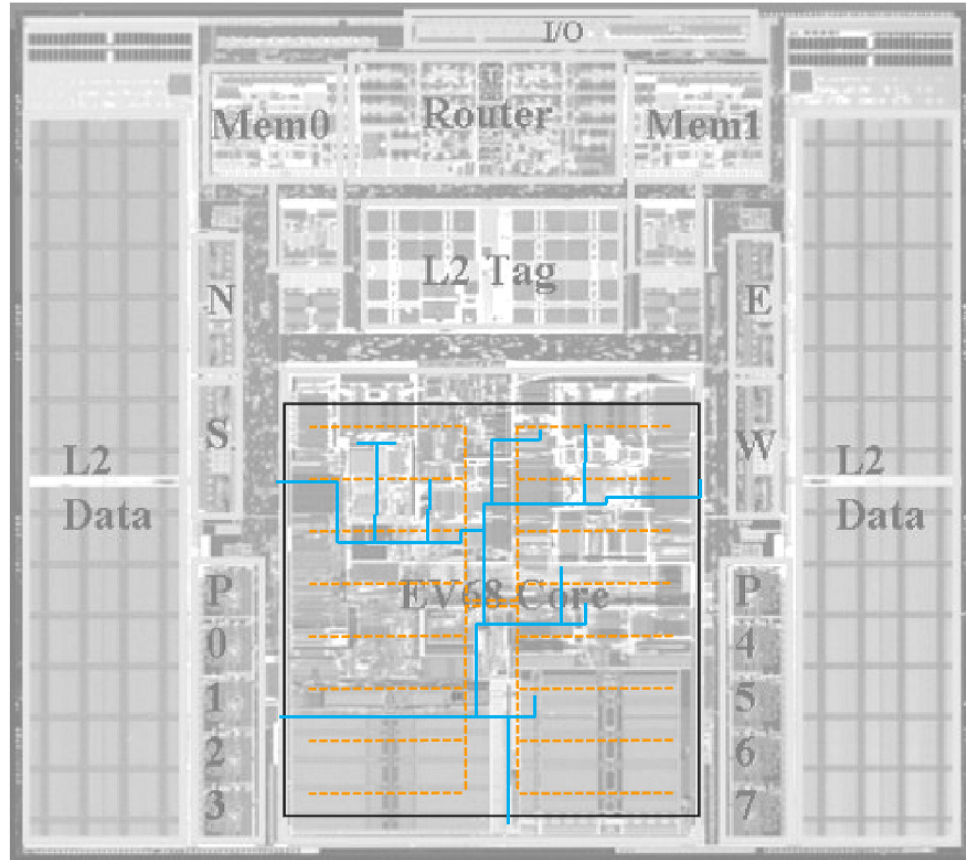


Figure 24: Layout of interconnect lines for single core monitoring [60].

These advances in sensor design, laboratory thermography, sensor placement, and thermal signal processing provide a basis for improved dynamic thermal management implementations as well as opportunities for further research. A discussion of new research directions is presented in Chapter 5.

CHAPTER 3: UNCERTAINTY IN HOTSPOT DETECTION

3.1. Introduction to Hotspot Detection

As microprocessor manufacturers have adopted multi-core circuit architectures, the detection and management of temporal hotspots have become increasingly important for chip reliability and performance. While much attention has been given to increases in the overall chip power, hotspot heat fluxes are increasing even more rapidly for many applications [62]. Active portions of a microprocessor can produce as much as 20 times as much heat as inactive regions [11]. These high heat fluxes can cause elevated junction temperatures leading to electromigration and subsequent circuit failure. Furthermore, temperature non-uniformities in the chip can cause severe thermo-mechanical stress on the package leading to system failure. These challenges will be exacerbated in future processors that are expected to include many more processor cores integrated in three-dimensional geometries.

To date, chip cooling alone does not seem capable of addressing these challenges. Most cooling solutions are best suited to address relatively slow thermal phenomena occurring over large regions of the chip. It is especially difficult to directly address highly localized, dynamic hotspots with cooling solutions implemented in chip packaging. Thermal engineers are forced to overdesign the cooling solution to satisfy worse-case scenario conditions for a hotspot region. This can be both difficult and expensive, particularly because the cost of cooling solutions increases rapidly as a function of maximum local heat flux [62]. Various methods of dynamic, localized

cooling (e.g. use of Peltier devices) are being investigated to address these difficult thermal requirements, but none have been adopted to date

An alternative overall approach to managing chip hotspots is to regulate the chip power output to maintain device temperature within specified limits. Such techniques are referred to as dynamic thermal management (DTM) and have been a subject of intense investigation since being introduced by Brooks and Martonosi [65].

All dynamic thermal management techniques fundamentally involve two steps: (1) interpreting temperature data from the chip and (2) responding to that data by reducing power. The majority of research has focused on the latter problem for DTM, specifically on finding innovative ways to locally regulate chip power. Proposed DTM techniques involve clock gating [7], Dynamic Frequency Control [8], Dynamic Voltage and Frequency Scaling (DVFS) [9], simultaneous multithreading (SMT) thread reduction [10], and activity migration [11]. Much less attention has been given to designing temperature sensor arrays and interpreting the resulting thermal signals. Two important sources of uncertainty need to be considered for DTM applications. First, the thermal sensors used for DTM feedback are subject to error. Most DTM studies do not consider the effect of this error and thus provide overly optimistic results. Skadron *et al.* [66] demonstrated that sensor error can cause significant performance reductions due to incorrect DTM triggering and reduced DTM threshold levels.

Discussions of uncertainty in DTM studies are typically limited to sensor error, but additional attention should be paid to the uncertainty caused by sensor placement. Because thermal sensors are not necessarily located at the chip hotspot, a DTM scheme must account for the temperature difference between the sensor location and the actual hotspot. Skadron *et al.*[66] used an estimated spreading factor within a core to try to account for this discrepancy as an additional source of error. In their study, the spreading factor contributed an additional 2°C error in the temperature signal. To attempt to account for uncertainty in hotspot location and intensity, DTM methods are currently designed to be conservative, which causes reduced system performance.

To reduce the uncertainties associated with thermal sensing for DTM, a challenging optimization problem must be considered. Circuit designs with high circuit density but low sensor density suffer from increased uncertainty in the thermal profile. Increased uncertainty about hotspot location and magnitude requires more cautious DTM control algorithms, which diminishes performance metrics. Increasing sensor resolution improves DTM control algorithms but also reduces circuit density, ultimately reducing computational power. An optimization approach is required to find a design that maximizes computational power while maintaining the chip in reliable operating conditions.

This study endeavors to help address this challenging optimization problem by quantifying the uncertainty that should be accounted for in a DTM scheme given a particular thermal sensor array. We consider the generalized case of a grid-array of thermal sensors located some distance above an arbitrary heat flux profile. In order to

better represent real applications, the thermal sensors are not necessarily located directly above known heat flux peaks. The chip heat flux profile is considered unknown, and the purpose of the thermal sensor array is to detect the regions of the chip that require dynamic power control.

We introduce a novel, computationally efficient, inverse heat transfer solution method and determine the accuracy to which it resolves the underlying heat flux profile. We consider cases with varying numbers of thermal sensors located with varying proximity to the circuitry level of the chip. Sensor error is also introduced to determine its effect on the estimated heat flux profile. For certain cases, the inverse solution method is shown to be susceptible to temperature sensor error. The results of these tests are compared to the uncertainty that results from treating the unprocessed thermal signal as a representation of the heat flux profile.

The approach taken here also has implications for the use of discrete thermal data in resolving the source of a hotspot. DTM schemes need not consider this uncertainty because the standard response to a hotspot is to throttle all activity in the vicinity. In chip development and production, however, thermal measurements are used to characterize the power distribution of the circuit design. For these tests, high resolution thermometry can be used (e.g. infrared microscopy [40]). The maximum spatial resolution at which these techniques can resolve neighboring hotspots is dictated by the resolution of the applied thermometry technique, the extent of thermal spreading in the chip, and the measurement error. The present study simulates the case of distinguishing two similar hotspot sources using discrete temperature

measurements. For a given measurement error and chip configuration, there is a minimum spatial sampling frequency required to correctly resolve the source of a hotspot.

Section II of this paper presents the overall methodology used to simulate chip heat flux profiles and determine the uncertainty associated with a particular thermal sensor array. Section III presents the inverse heat transfer solution method derived for this study. The uncertainties in the heat flux profile associated with direct temperature interpretation and inverse solution method are presented in Section IV. Section V provides concluding remarks.

3.2. Simulation Methodology

3.2.1. Overall Simulation Methodology

The present study is based on a simplified conduction model for the chip. Figure 25 shows the model geometry. The chip is modeled as an isotropic, single-layer structure. The isotropic condition can be relaxed by transformation of the thermal conductivity and chip thickness [67]. The boundary condition on the top surface is convective heat transfer with a uniform heat transfer coefficient. The boundary conditions on the four sidewalls are adiabatic. On the bottom surface, an arbitrary heat flux profile boundary condition is applied. The chip is 1 cm by 1cm and its thickness and thermal conductivity is varied in the simulations. The system operates in steady-state. This simplified model of the chip facilitates the generalized simulation methodology taken in this study which would be impractical with a highly discretized chip model.

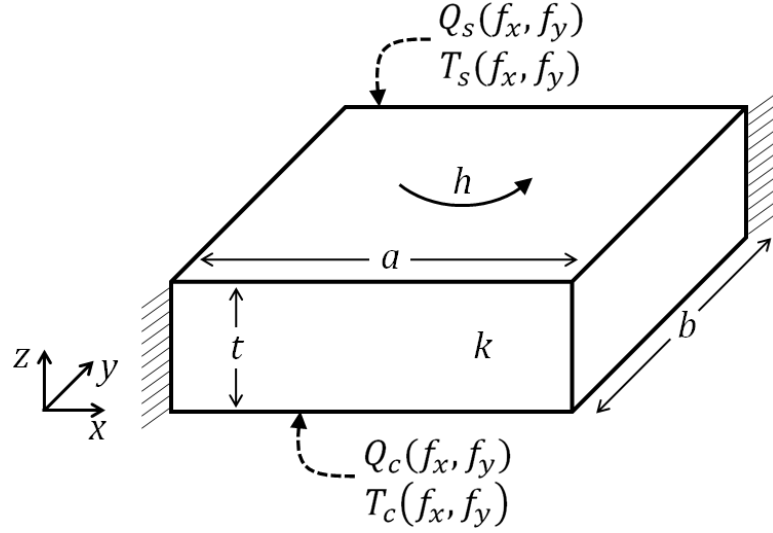


Figure 25: Schematic of model geometry. An arbitrary heat flux profile is applied on the bottom boundary. The boundary condition on all sidewalls is adiabatic; the boundary condition on the top surface is uniform heat transfer.

Figure 26 shows the four main steps involved in the overall simulation methodology.

The simulation begins by defining the geometry and system parameters and generating a randomized heat flux profile. The forward solution method is used to resolve the sensor-level, full-resolution temperature profile, $T_{s,f}$ (Figure 26.b), based on the circuit-level, full resolution heat flux profile, $Q''_{c,f}$ (Figure 26.a). A set of low-resolution temperature profiles, $T_{s,l}$ (Figure 26.c), is created by interpolating the full-resolution temperature profile, $T_{s,f}$, at various spatial frequencies. Each low-resolution temperature profile represents the temperature profile that would be measured by a temperature sensor array of a particular spatial frequency. For example, for a temperature sensor spatial sampling frequency of 1000 m^{-1} (equivalent to nominal sensor spacing of 1mm), the low-resolution temperature profile, $T_{s,l}$, is a 10x10 grid on a 1cm by 1cm chip.

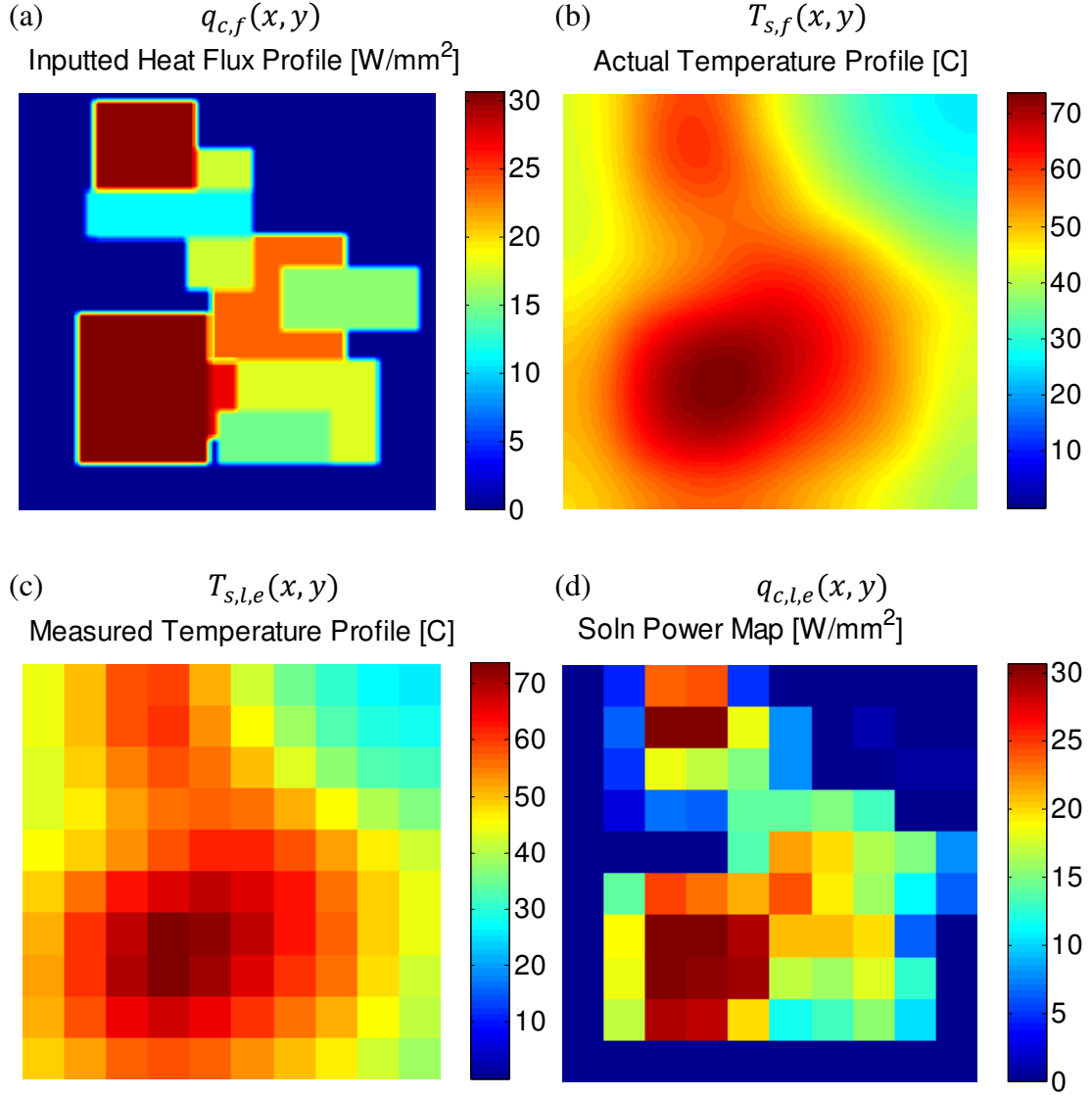


Figure 26: Representative images of each of the four main steps in the simulation methodology. The inputted heat flux profile (a) is used as a reference for determining the error in (d) the calculated heat flux profile.

Random error is added to the low-resolution temperature profile, $T_{s,f}$, to simulate the measurement error introduced by real temperature sensors. The sensor error, T_{error} , is normally-distributed about the interpolated temperature value with a standard

deviation that is specified relative to the maximum interpolated temperature. The sensor error at each index is calculated as:

$$T_e(i, j) = \sigma_{relative} T_{s,l,max} \Gamma \quad (4)$$

where $\sigma_{relative}$ is the standard deviation of the relative sensor error, $T_{s,l,max}$ is the maximum measured temperature, and Γ is a random number with a mean value of zero and a standard deviation of unity.

This study shows the results for standard deviations in the relative sensor error of 0, 0.5, and 1 percent. The case of 0 percent standard deviation in the relative sensor error is equivalent to no measurement error.

The sensor-level, low-resolution temperature profile with error, $T_{s,l,e} = T_{s,l} + T_e$, is used to calculate the circuit-level, heat flux profile, $Q''_{c,l,e}$ (Figure 26.d), using a spatial sampling frequency domain, inverse heat transfer solution, described in detail in the next section. Because the inputted temperature profile is low resolution, the resulting heat flux profile, $Q''_{c,l,e}$, is also low resolution. To calculate the error resulting from the solution method, the low-resolution heat flux profile is interpolated to full resolution. The mean absolute error (MAE) is calculated by finding the difference between the correct profile and the calculated heat flux profile:

$$MAE = \frac{1}{N_x N_y} \sum_{i=1}^{N_x} \sum_{j=1}^{N_y} |Q''_{c,f}(i, j) - Q''_{c,l,e}(i, j)| \quad (5)$$

where N_x and N_y are the total number of indices in the x and y directions, respectively.

The mean absolute error is normalized by the average heat flux:

$$\text{Normalized MAE} = \frac{\frac{1}{N_x N_y} \sum_{i=1}^{N_x} \sum_{j=1}^{N_y} |Q''_{c,f}(i,j) - Q''_{c,l,e}(i,j)|}{\frac{1}{N_x N_y} \sum_{i=1}^{N_x} \sum_{j=1}^{N_y} Q''_{c,f}(i,j)} \quad (6)$$

$$\text{Normalized MAE} = \frac{\sum_{i=1}^{N_x} \sum_{j=1}^{N_y} |Q''_{c,f}(i,j) - Q''_{c,l,e}(i,j)|}{\sum_{i=1}^{N_x} \sum_{j=1}^{N_y} Q''_{c,f}(i,j)} \quad (7)$$

Figure 27 provides a block diagram of the simulation procedure. The procedure is repeated for numerous randomly-generated heat flux profiles and the results are averaged. The average MAE is plotted against the thermal sensor spatial sampling frequency.

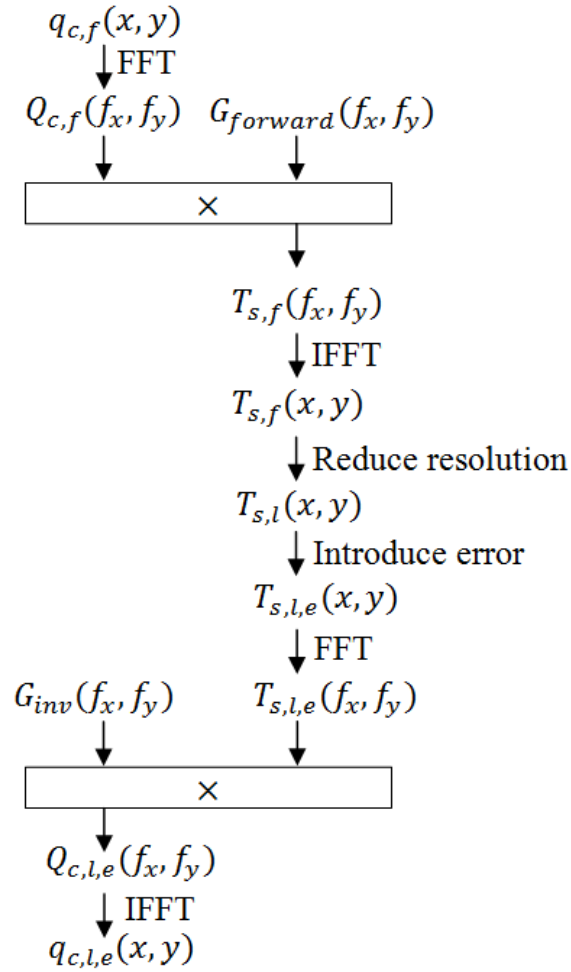


Figure 27: Block diagram of numerical approach used for determining hotspot detection accuracy. FFT and IFFT refer to the Fast Fourier Transform and the Inverse Fast Fourier Transform, respectively.

In practice, an inverse heat transfer technique is not always used to interpret measured temperature profiles. Instead, the measured temperature profile is assumed to be representative of the chip heat flux profile. This technique is equivalent to treating the measured temperature profile as directly proportional to the heat flux profile:

$$Q_c''(i,j) = \frac{T_s(i,j)}{R_{th}''} \quad (8)$$

where R_{th}'' is the chip vertical thermal resistance for unit area:

$$R_{th}'' = \frac{t_0}{k} \quad (9)$$

This simplification results in additional uncertainty in the heat flux profile, the magnitude of which depends on the chip properties and boundary conditions. In this paper, this approach is referred to as “direct interpretation” of the temperature profile and is compared to the inverse solution method in the results section.

3.2.2. Randomization of Heat Flux Profiles

The uncertainty in the calculated heat flux is dependent on the characteristics of the heat flux profile. Simple, well-spaced heat flux profiles are easier to resolve than overlapping, complicated heat flux profiles. To represent the most general case, the simulation is conducted over a set of heat flux profiles that contain varying degrees of complexities. The heat flux profiles are randomly generated to include between 1 and 15 hotspots which can vary in laterals dimension between 273 μm (equivalent to 7 grid cells) and 4.18 mm (equivalent to 107 grid cells). For reference, the chip is 1 by 1cm. The hotspots are created with soft edges; the edge of the hotspot spans 156 μm (equivalent to 4 grid cells) and has a linear slope from the value of the background heat flux to value of the hotspot heat flux. The background heat flux is 1 W/cm^2 and the maximum possible hotspot heat flux is 320 W/cm^2 . Hotspots are permitted to

overlap with each other but not with the edge of the chip. For the first set of simulations, the hotspots have a random heat flux value between the background and the maximum heat flux. This is referred to as “variable heat flux”. For the second set of simulations, all hotspots have the maximum heat flux, referred to as “binary heat flux”. The case of binary heat flux represents a core that is either active or inactive. The case of variable heat flux represents a core for which the amount of activity is unknown. Since the variable heat flux case is most challenging from an uncertainty perspective, only select results are presented for binary heat flux cases.

The key result of each simulation is the mean absolute error (MAE) in the calculated heat flux profile. Because conduction through the chip is linear, the results are generalized by normalizing the error in the heat flux profile by the input heat flux profile. Thus only the relative magnitude of the heat flux as compared to the background heat flux is relevant for consideration.

3.2.3. Resolution Study

A second study was conducted to quantify the ability of the inverse solution method to resolve a single hotspot from a group of neighboring hotspots. Two circuit-level heat flux profiles are created; the first heat flux profile, referred to as Case I, consists of a single hotspot in the center while the second heat flux, referred to as Case II, profile consists of 9 closely packed hotspots in the center. The average heat flux is the same in both cases. Figure 28 shows the two heat flux profiles. The circuit-level temperature profile resulting from the single-hotspot heat flux profile is calculated using the

forward solution. The temperature profile is sampled at reduced spatial sampling frequency to simulate the signal from a thermal sensor array, as before.

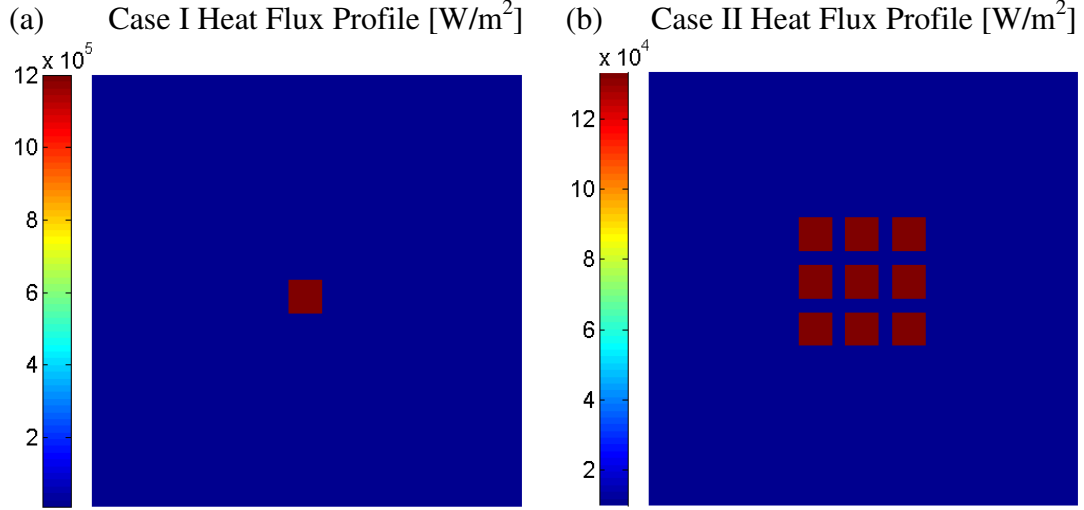


Figure 28: Heat flux profiles used for resolution study, referred to as Case I and II respectively. Both heat flux profiles have equivalent average heat flux and produce similar temperature response profiles. The solution methods are tested for their ability to correctly resolve these heat flux profiles.

Each solution method is used to deduce which of two possible heat flux profiles yielded the measured temperature profile. To do so, the inverse solution method is used to calculate the circuit-level heat flux profile. The results are compared to the two possible inputted heat flux profiles by calculating the mean absolute error. The profile resulting in the lower MAE represents the solution chosen by the inverse solution method. For example, if the MAE between the calculated heat flux profile and the single-hotspot heat flux profile is lower than the MAE between the calculated heat flux profile and the multi-hotspot heat flux profile, the inverse solution method chooses the single-hotspot heat flux profile. If the choice correctly corresponds to the

actual inputted heat flux profile, the inverse solution method is correct. This procedure is conducted for all sensor spatial frequencies, and is also conducted for the direct interpretation method.

3.3. Spatial Frequency Domain Inverse Heat Transfer Solution

3.3.1. *Inverse Heat Transfer Solution Method*

To conduct the forward and inverse solutions needed for the overall simulation methodology, an analytical, spatial-frequency domain heat transfer analysis has been developed. This approach is more computationally efficient than finite-difference methods and thus facilitates rapid multi-parameter design optimization and possible integration into DTM schemes.

The thermal profile in the model geometry is defined by the heat diffusion equation. For each layer in the stack, the solution to the heat diffusion equation is given by:

$$\begin{aligned}
 T(x, y, z) = & A_0 + B_0 z \\
 & + \sum_{m=1}^{\infty} [A_m \cosh(\lambda_m z) + B_m \sinh(\lambda_m z)] \cos(\lambda_m x) \\
 & + \sum_{n=1}^{\infty} [A_n \cosh(\gamma_n z) + B_n \sinh(\gamma_n z)] \cos(\gamma_n y) \quad (10) \\
 & + \sum_{m=1}^{\infty} \sum_{n=1}^{\infty} [A_{mn} \cosh(\beta_{mn} z) \\
 & + B_{mn} \sinh(\beta_{mn} z)] \cos(\lambda_m x) \cos(\gamma_n y)
 \end{aligned}$$

where:

$$\beta_{mn} = \sqrt{\lambda_m^2 + \gamma_n^2} \quad (11)$$

$$\lambda_m = \frac{m\pi}{a} \quad (12)$$

$$\gamma_n = \frac{n\pi}{b} \quad (13)$$

For the boundary conditions imposed in this model, Etessam-Yazdani [67] demonstrated a technique of representing this conduction problem as a two-port terminal network. The technique has been shown to be both accurate and fast for the forward heat transfer solution [68] and is adapted in this study for the inverse problem.

Figure 29 presents a schematic of the two-port terminal network for this system. The two-dimensional Fourier transforms of the heat flux profiles at the circuit and sensor levels are Q_c'' and Q_s'' , respectively.

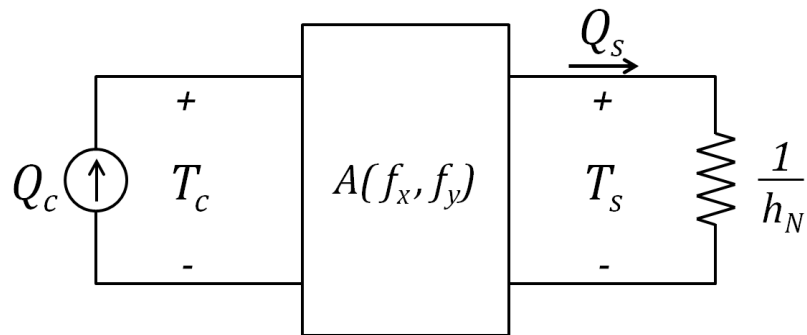


Figure 29: Schematic of two-port terminal network [67].

Similarly, T_c and T_s are the two-dimensional Fourier transforms of the temperature profiles at the circuit and sensor levels, respectively. The matrix A is a 2x2 matrix that relates T_c and Q_c'' to T_s and Q_s'' :

$$\begin{bmatrix} T_c(f_x, f_y) \\ Q_c''(f_x, f_y) \end{bmatrix} = A(f_x, f_y) \begin{bmatrix} T_s(f_x, f_y) \\ Q_s''(f_x, f_y) \end{bmatrix} \quad (14)$$

For radial spatial frequency $f_r > 0$:

$$A(f_x, f_y) = \begin{bmatrix} \cosh(2\pi f_r t_o) & \frac{\sinh(2\pi f_r t_o)}{2\pi f_r k} \\ 2\pi f_r k \sinh(2\pi f_r t_o) & \cosh(2\pi f_r t_o) \end{bmatrix} \quad (15)$$

And for $f_r = 0$:

$$A(f_x, f_y) = \begin{bmatrix} 1 & t_o/k \\ 0 & 1 \end{bmatrix} \quad (16)$$

where the radial spatial frequency f_r is defined as:

$$f_r = \sqrt{f_x^2 + f_y^2} \quad (17)$$

Further details on the derivation of the two-port terminal analysis are provided in [67].

Etessam-Yazdani *et al.* [67] used the two-port terminal analysis to solve for the temperature as a function of the heat flux on the same level of the geometry, which represents the forward solution. In this study, the solution was modified to determine the heat flux profile on the circuit plane, Q_c'' , using the temperature profile on the

sensor plane, T_s , which represents the inverse solution. From the two-port terminal analysis, the equation for Q_c'' is:

$$Q_c'' = A_{21}T_s + A_{22}Q_s'' \quad (18)$$

Applying the top boundary condition, $Q_s = hT_s$, and substituting the appropriate values of A_{ij} , the result for cases where $f_r > 0$ is:

$$Q_c = (2\pi f_r k \sinh(2\pi f_r t_o) + h * \cosh(2\pi f_r t_o))T_s \quad (19)$$

For which the inverse solution transfer function $G_{inv}(f_r)$ can be defined such that:

$$Q_c = G_{inv}(f_r) * T_s \quad (20)$$

and

$$G_{inv}(f_r) = 2\pi f_r k \sinh(2\pi f_r t_o) + h * \cosh(2\pi f_r t_o) \quad (21)$$

For cases where $f_r = 0$, the transfer function reduces to equal the heat transfer coefficient, h , and the equation is given as $Q_c = hT_s$.

3.3.2. *High-Frequency Filtering*

A filtering technique based on the forward solution transfer function is employed to reduce error in the inverse solution method. As shown by [69], the forward solution to the conduction problem yields a transfer function in the frequency domain that acts as

a low pass filter. Physically this represents the attenuation of high spatial frequency components of the thermal signal via heat spreading in the chip.

The inverse transfer function has the form of a high-pass filter, as shown in Figure 30.

The minimum of the transfer function occurs at $f_r = 0$ and increases rapidly as a function of f_r , thus amplifying the high frequency components of the temperature profile. The components of the temperature profile that are greater than the -3dB frequency of the forward solution transfer function, however, represent sensor noise.

A filtering method has been developed to prevent this noise from propagating to the calculated heat flux profile. A low-pass filter is applied to the inverse transfer function with a filter cut-off frequency at the -3dB frequency of the forward solution transfer function. The filter has a soft roll-off. Figure 30.c shows the filtered transfer function. This filtering technique dramatically improves the performance of the inverse solution method by decreasing sensitivity to high-frequency noise.

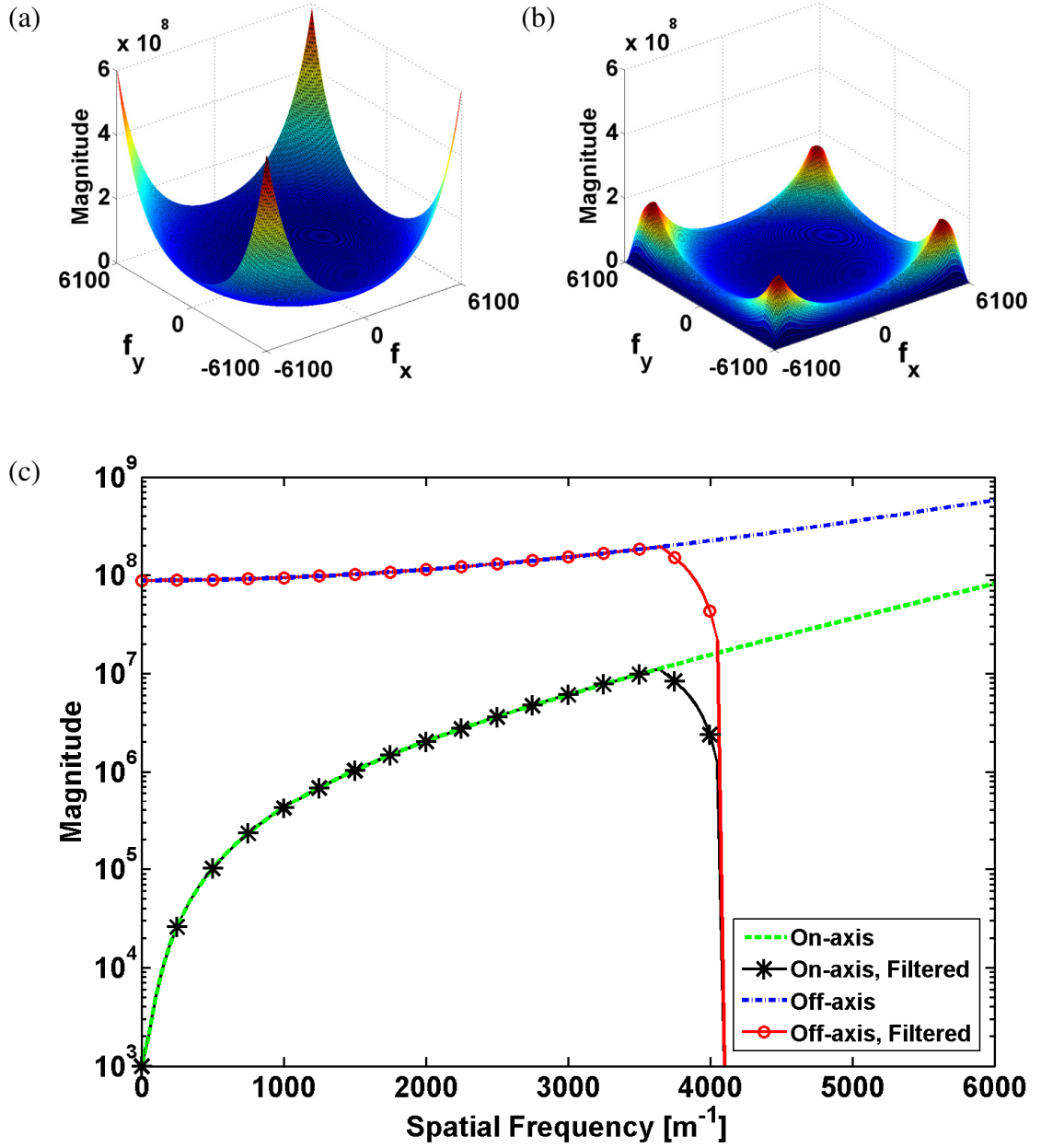


Figure 30: Representative plots of inverse solution transfer function. Plots show two-dimensional shape of the transfer function (a) without filtering and (b) with filtering. (c) Values of the transfer function for varying x-direction spatial frequency and for y-direction frequency of zero (labeled “on-axis”) as well as for maximum y-direction frequency (labeled “off-axis”). The filter roll-off occurs at approximately 4000 m^{-1} .

3.3.3. Solution Validation

The solution method was validated by comparison to COMSOL Multiphysics software using representative simulation parameters. The heat transfer coefficient was $10 \text{ W/m}^2\text{-K}$ and the thermal conductivity was 148 W/m-K . The simulated chip was 1 cm by 1 cm in lateral dimensions and 100 microns in thickness. A representative heat flux was applied in the COMSOL model and the temperature profile was resolved. The temperature profile was used as an input to the inverse solution method and the applied heat flux was calculated. The calculated heat flux matched the COMSOL heat flux at greater than 0.01% accuracy.

Additional testing was conducted to ensure the results for average heat flux error are independent of the number of random heat flux maps tested, N . Figure 31 shows the results for varying number of randomly generated heat flux maps for both varying heat flux and binary heat flux. The results are shown to be N -independent (i.e. independent of the number of random heat flux profiles) after 50 randomly generated heat maps. For all of the reported results, data was averaged for 50 heat maps ($N = 50$).

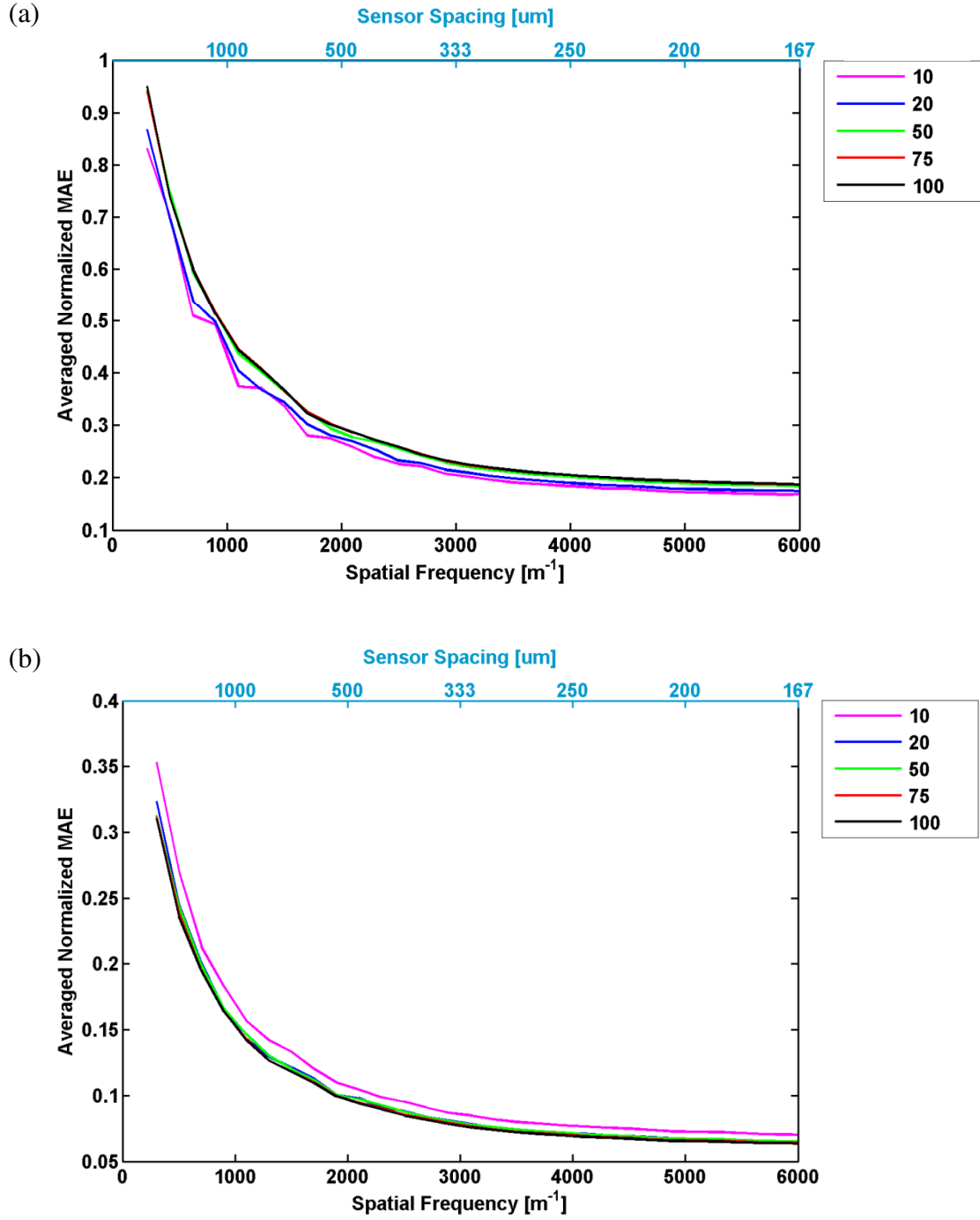


Figure 31: Average mean absolute error (MAE) for varying numbers of randomized heat flux profiles for (a) variable heat flux and (b) binary heat flux. Results for both cases are independent of the number of heat flux profiles for more than 50 heat flux profiles.

3.4. Simulation Results

Figure 32 shows representative distributions of mean absolute error for 50 randomized heat flux profiles. Results are reported for the case of variable heat flux and binary heat flux. These results provide a basis for understanding the effects of sensor spatial frequency on the calculated heat flux profile. Simulation parameters are typical of chip applications: the distance from the sensor array is 100 μm , the conductivity is 148 W/m-K and the heat transfer coefficient is 10,000 $\text{W/m}^2\text{-K}$. The sensor error is zero for this case. The mean value (shown in bold black) follows the expected trend of increased accuracy at higher spatial sampling frequency. A sampling frequency of 2000 m^{-1} (approximately 500 μm sensor spacing) is required to achieve an average mean absolute error (MAE) below 25% for the variable heat flux case. At lower resolutions, the average MAE is dramatically higher. Significant deviations from the mean value are caused by variations between the randomized heat flux profiles.

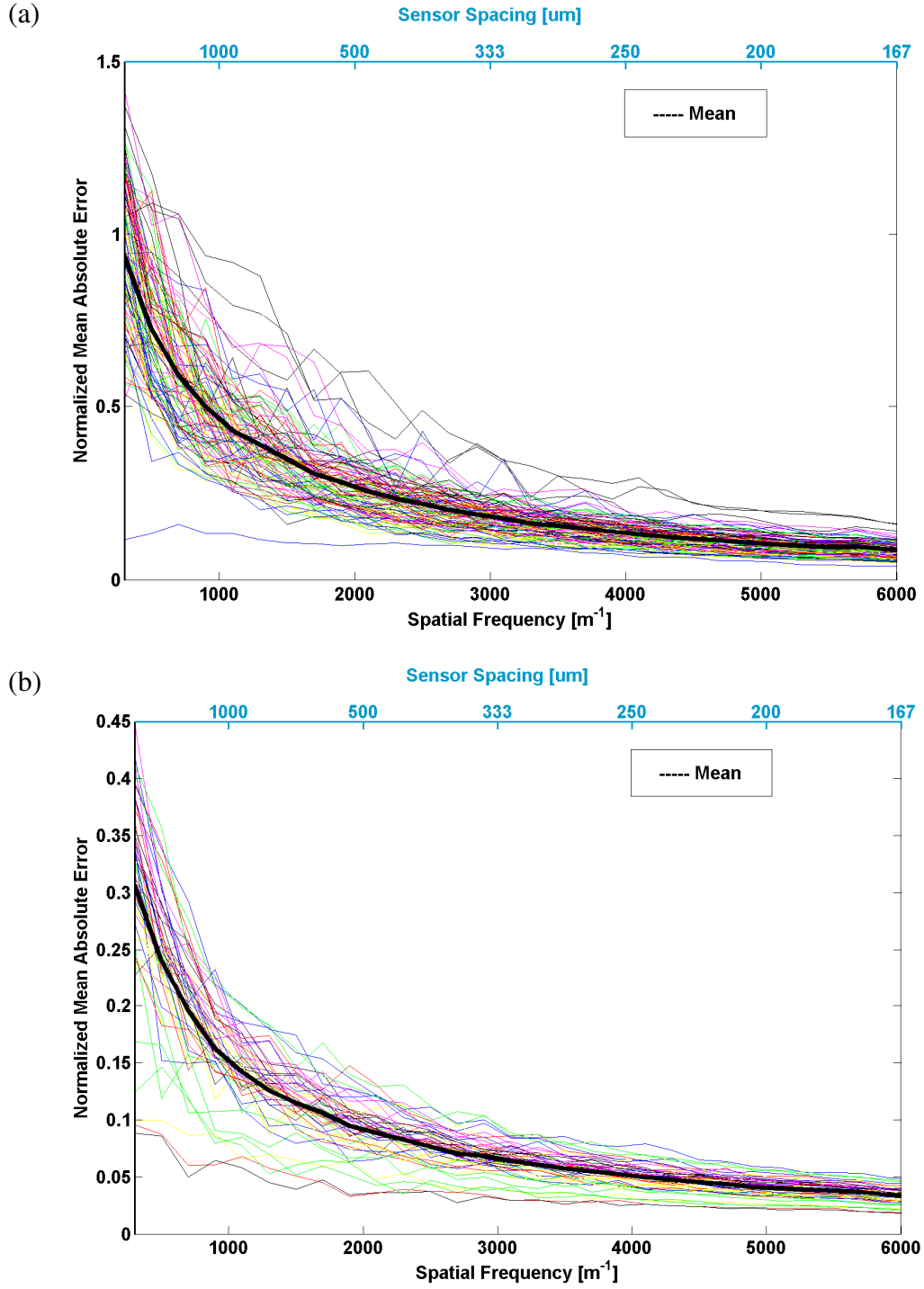


Figure 32: Demonstration of the averaging technique for (a) variable heat flux and (b) binary heat flux. Results for 50 heat flux profiles are shown. The bold black line indicates the average value.

The average MAE error is dependent on the heat transfer coefficient, the sensor error, and the proximity between thermal sensors and the circuit level. These effects are discussed in more details below. For clarity, only the average MAE is shown. The solid curves and dotted curves represent the average MAE for the inverse solution method and the direct temperature interpretation method, respectively.

The average MAE of the inverse solution method is dependent on whether the inputted heat flux profile is binary or variable. Figure 33 shows an approximately 65% drop in average MAE if the input heat flux is binary rather than variable. Since the heat flux cannot always be assumed to be binary, the remaining plots show results for variable heat flux.

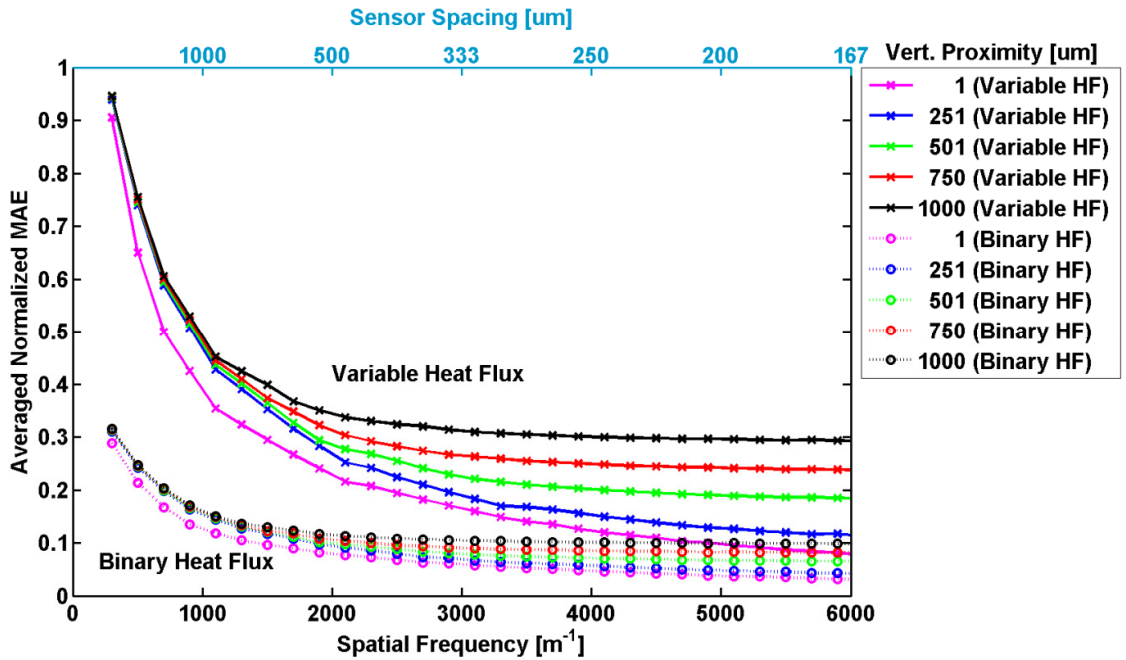


Figure 33: Effects on uncertainty of variable versus binary inputted heat flux profile for varying vertical proximity between sensor and circuit level. The binary heat flux profile results in substantially lower MAE.

Figure 34 shows the performance of the inverse solution method for varying heat transfer coefficients for variable heat flux with no sensor error. As expected, the average MAE for both methods is reduced by increasing heat transfer coefficient values. The direct interpretation method performs poorly at low heat transfer coefficients but makes significant improvements as the heat transfer coefficient is increased. The inverse method produces significantly lower average MAE and is less sensitive to changes in the heat transfer coefficient.

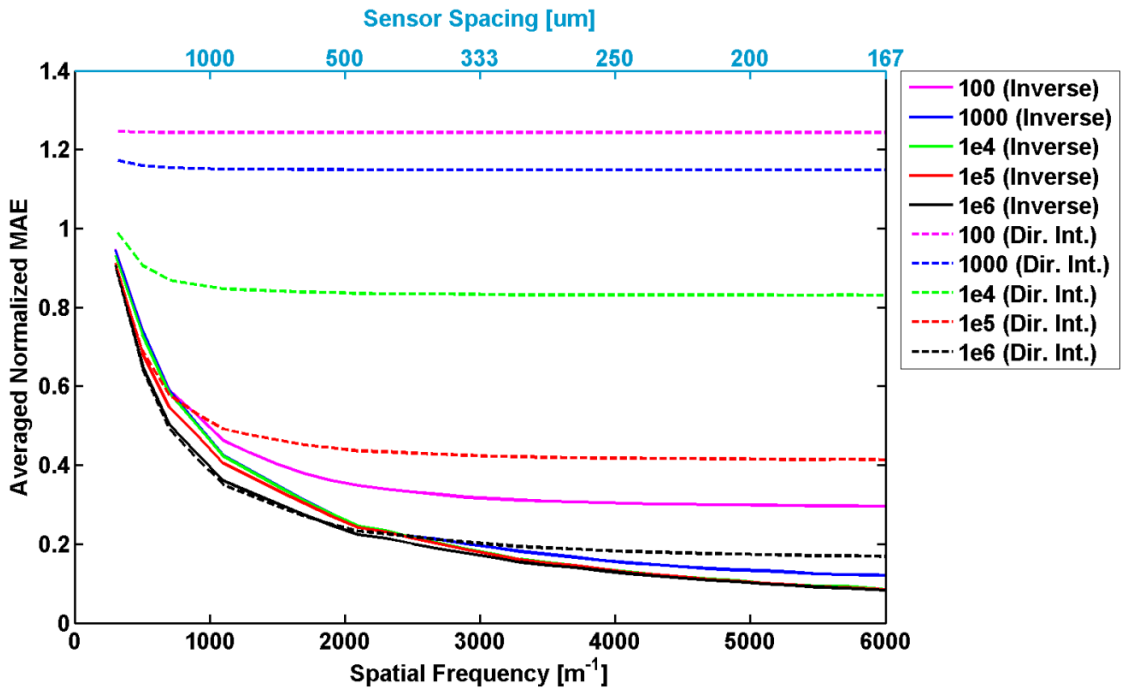


Figure 34: Uncertainty in calculated heat flux profile for varying convective heat transfer coefficient. The inverse solution method is much less sensitive to heat transfer coefficient than the direct interpretation method.

Figure 35 illustrates the difficulty of calculating the heat flux profile from temperature profiles containing sensor error. For the ideal case of zero sensor error, the inverse solution method outperforms the direct method by up to 50% MAE for variable heat

flux. However, measurement error causes the inverse method to diverge from the solution. For a case of 0.5% measurement error, the inverse solution is slightly better than the direct method for spatial frequencies up to about 3000 m^{-1} , at which point it diverges rapidly. For the case of 1% standard deviation in the sensor error, the direct interpretation method is superior for sensor spatial frequencies greater than 2000 m^{-1} . Similar trends are observed for the case of binary heat flux profiles as well.

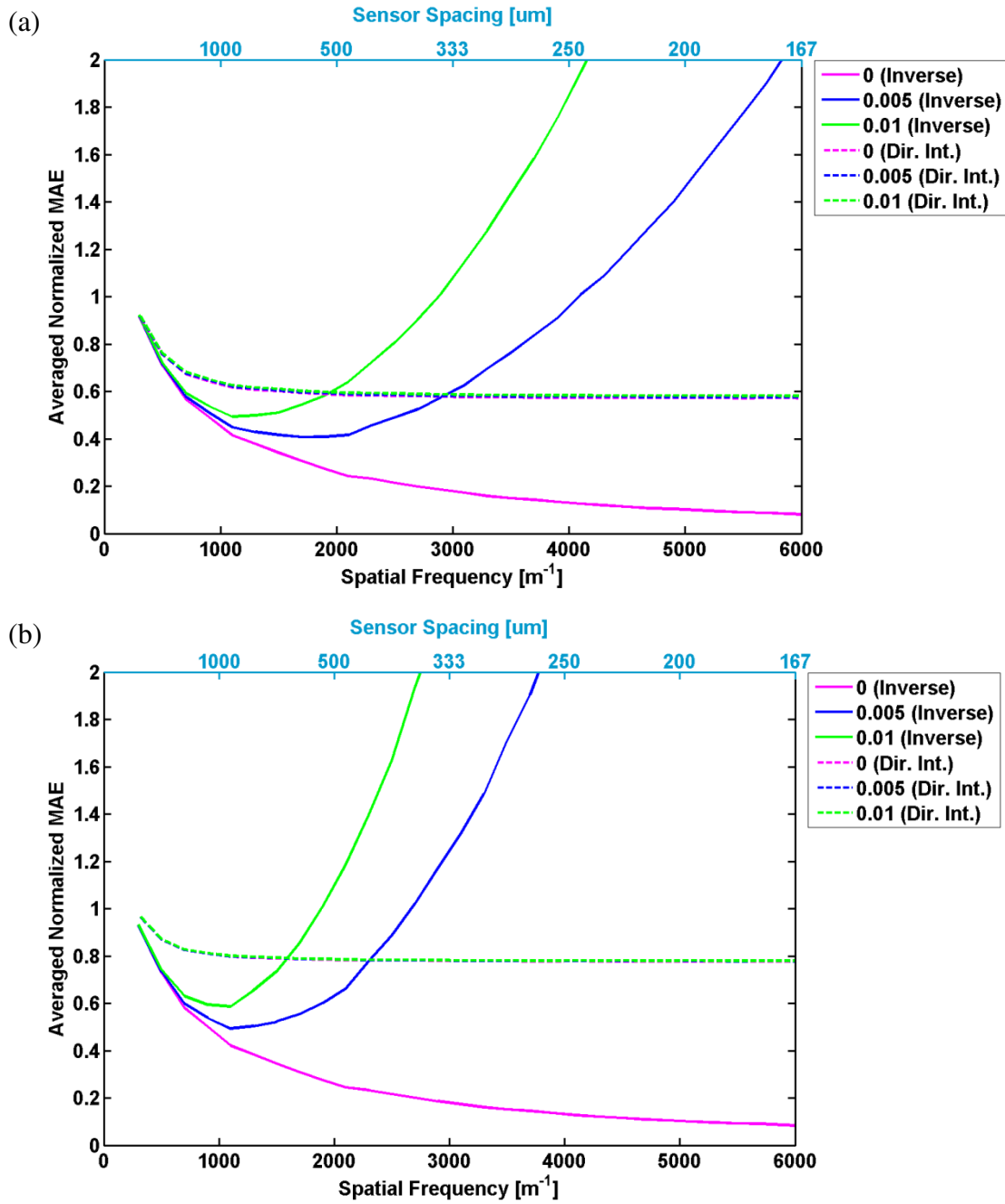


Figure 35: Uncertainty in calculated heat flux profile for varying sensor error at a vertical proximity of (a) 2.575 μm and (b) 7.53 μm . The inverse solution method is susceptible to sensor error at high spatial frequency. The MAE for the direct interpretation method is not affected by varying sensor error.

Figure 36 presents the effect of vertical proximity between the sensor level and the circuit level. Average MAE results are shown for vertical distances between 1 μm and 1 mm for variable heat flux profiles. As the vertical proximity is reduced, modest improvements in MAE are observed for both the inverse and direct interpretation techniques with the exception of the 1 μm case where improvements in the direct interpretation method are approximately 0.8 normalized averaged MAE. For the extreme case of 1 μm of vertical proximity, the inverse and direct interpretation methods are comparable, but for all other cases the inverse solution significantly outperforms the direct interpretation method.

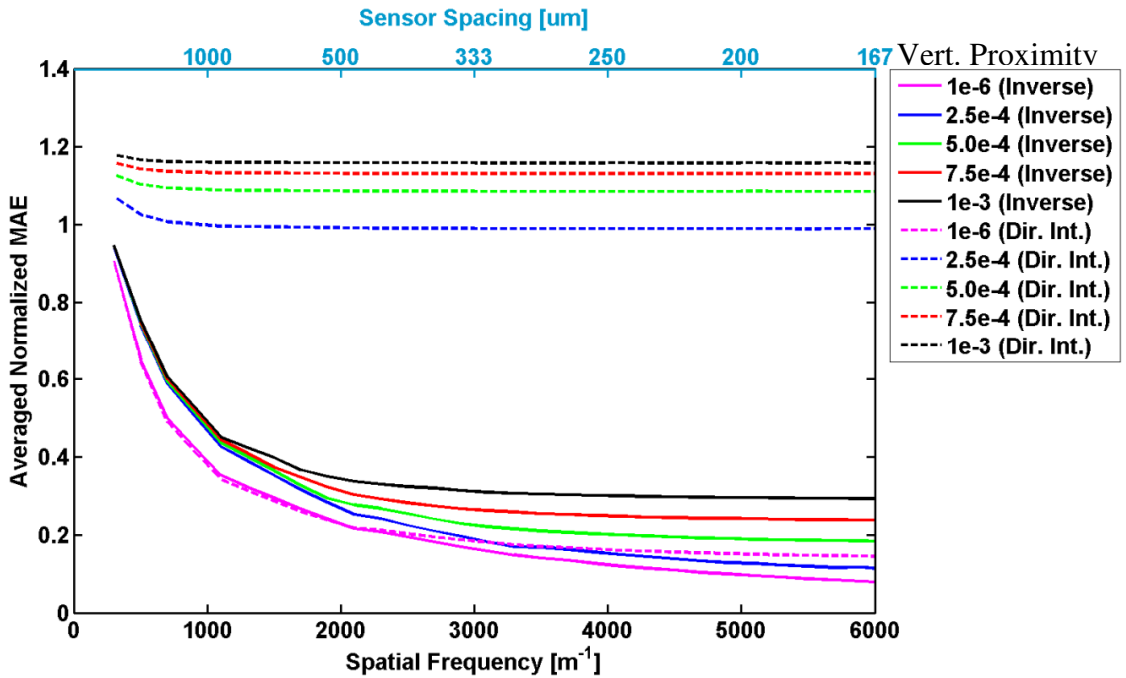


Figure 36: Uncertainty in calculated heat flux profile for varying vertical proximity between the sensor and circuit levels for zero sensor error. For most cases, large changes in vertical proximity yield modest improvements in heat flux uncertainty.

Figures 37 and 38 show the performance of the inverse solution method in resolving neighboring hotspots. The figures show the minimum sensor spatial frequency required to correctly differentiate between a single hotspot and a group of equivalent neighboring hotspots. The results are presented as a function of vertical proximity between the distributed thermal sensor array and the circuit plane, and a moving-average smoothing function is applied to remove discretization artifacts. The gray region of the plot shows the domain in which the inverse solution can correctly identify the underlying heat flux profile. A relatively low sensor spatial frequency is adequate when positioned in close proximity to the hotspot. Increasing the separation between the sensor array and the hotspot requires an increase in the sensor spatial frequency. Figures 37 and 38 show results for convective heat transfer coefficients of 10,000 and 50,000 W/m²-K, respectively. For a convective heat transfer coefficients of 10,000 W/m²-K at distances greater than approximately 240 μ m, the inverse solution method is unable to resolve the hotspot. Figure 38 shows that the limit of the inverse solution can be extended by increasing the heat transfer coefficient. For this case, the inverse solution method produces the correct results up to 300 μ m. For all cases shown, the direct interpretation method failed to correctly identify the single hotspot. The inverse technique is shown to be superior to the direct interpretation method for resolving neighboring hotspots. These results provide insight into the optimization of sensor vertical proximity and sensor spatial frequency for resolving neighboring hotspots.

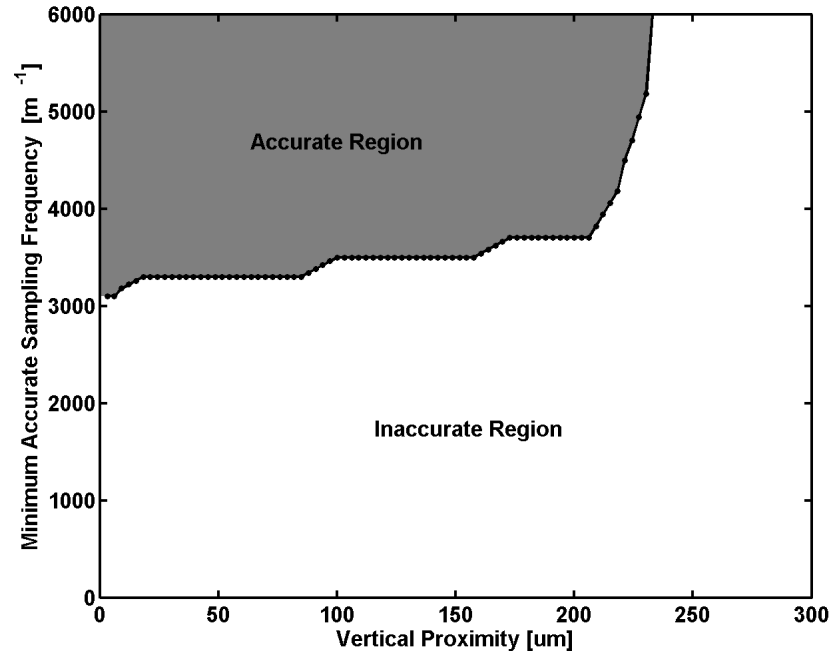


Figure 37: Plot of minimum accurate sampling frequency as a function of vertical proximity between chip and sensor level for heat transfer coefficient of $10^4 \text{ W/m}^2\text{-K}$. The inverse solution method is accurate in the shaded region. The direct interpretation technique is inaccurate across the entire domain.

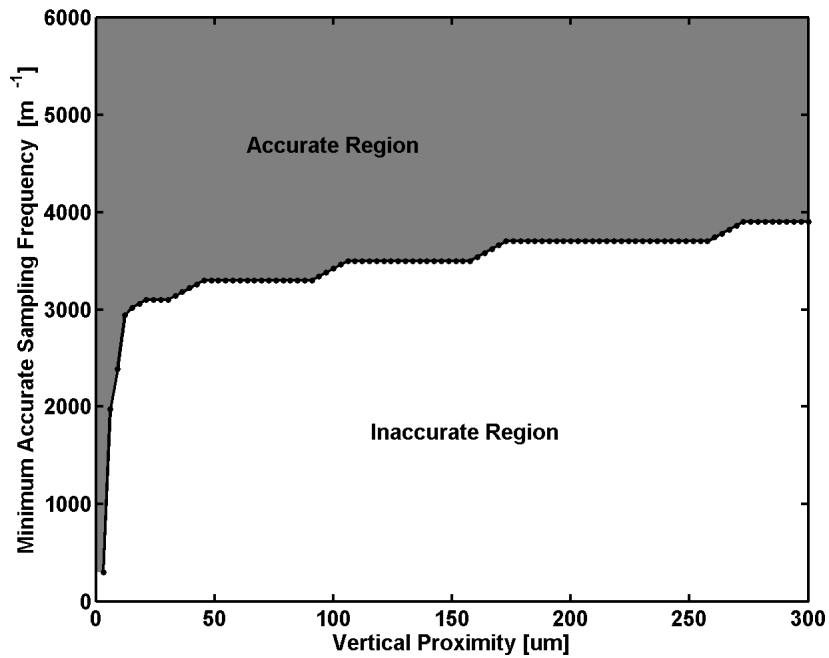


Figure 38: Plot of minimum accurate sampling frequency as a function of vertical proximity between chip and sensor level for heat transfer coefficient of $10^5 \text{ W/m}^2\text{-K}$.

The inverse solution method is accurate in the shaded region. The direct interpretation technique is inaccurate across the entire domain.

This study investigates uncertainty and error propagation in distributed thermal sensor arrays in microprocessors. A novel, inverse heat transfer solution methodology is developed to provide a computationally efficient method for determining the heat flux profile at a remote level in a chip. The inverse solution method is used to determine the expected mean absolute error of the calculated heat flux profile in a chip. Several key conclusions are drawn.

For systems with relatively low sensor spatial frequency such as typical microprocessors, large improvements in the accuracy of the calculated heat flux can be made by making relatively small improvements in the resolution of the sensor array. As the sensor array increases resolution, the uncertainty in the calculated heat flux is much reduced. For cases of very low sensor error, the proposed inverse solution technique more accurately calculates the heat flux profile than direct interpretation of the temperature profile.

Depending on the system configuration and the magnitude of the sensor error, the inverse solution method can be inaccurate. This inaccuracy is mitigated by the proposed filtering method, but nonetheless represents a fundamental limitation of this technique. Direct interpretation of the temperature signal is shown to result in significant error in the calculated heat flux profile. Accounting for these errors in DTM techniques causes decreased computational performance and should therefore be considered during overall system design.

These conclusions regarding the nature of error propagation from distributed thermal sensor arrays can provide a basis for considering the difficult system-level optimization required for integrated circuit design. Sensor error, sensor spatial frequency, proximity between a sensor array and hotspots, and signal processing all affect hotspot uncertainty as well as circuit design. Each of these parameters can help improve DTM accuracy but can also pose costs for the performance of the circuit. Careful optimization of these parameters is necessary to maximize computational performance while ensuring reliable thermal conditions.

CHAPTER 4: FAST CALCULATION OF TEMPERATURE EVOLUTION IN ELECTRONIC DEVICES

4.1. Introduction to Transient Hotspot Modeling

Thermal management of high-performance integrated-circuit chips has become one of the most critical design challenges throughout all integrated-circuit architectural and manufacturing communities. As transistor feature size in semiconductor devices continues to shrink delivering increased performance, the corresponding power densities and operating frequencies have increased rapidly. Multi-core circuit architectures further exacerbate the problem by creating highly localized transient heat fluxes. These high heat fluxes can cause temperature excursions that have major adverse impacts on device performance, reliability and power efficiency.

To date cooling solutions have been unable to address these highly localized, transient hotspots; furthermore, cost constraints in industrial applications suggest that cooling alone will not be able to resolve these thermal challenges. Instead, attention has shifted to dynamic thermal management (DTM) [65], as well as accurate evaluation of thermal behaviors under long power traces for thermal-aware optimization applications during the early stages of architecture-level designs [70]. Developing such models is challenging because the device thermal response strongly depends on the temporal pattern of input power, the disparate thermal time constants of the components, and varying boundary conditions. Furthermore, implementing these models in chip-level runtime temperature regulation requires highly efficient model computation.

Widely available numerical modeling software (e.g. finite element models) is inappropriate for runtime applications due to extensive requirements on computational resource. These models instead serve as inputs for the model compression technique used in this work and provide a reference for model validation.

Improvements in modeling efficiency can be made by recognizing the fact that the thermal response is only needed for specific regions of the chip, typically at junction or hot spot locations. Hence, various dynamic compact thermal models with a decreased number of model parameters have been developed for rapid calculation of transient temperature responses.

The methods for constructing dynamic compact thermal models can be divided into two general categories: thermal RC network approach and thermal RC ladder approach. The first approach constructs an equivalent thermal RC network that accurately describes dynamics of the thermal system. This can be achieved by transforming the spatially-discretized system matrices of the governing equation in finite element/volume models into a thermal circuit network consisting of thermal resistance element interconnecting neighboring nodes and heat capacity element to the reference thermal ground [71]. The thermal circuit network is used to formulate a system of ordinary differential equations (ODE):

$$[C] \dot{\vec{T}} + [R]^{-1} \vec{T} = \vec{F} \vec{P} \quad (22)$$

where $[C]$ and $[R]$ are the thermal capacitance and thermal resistance matrices, \vec{T} is the vector of node temperatures, $\dot{\vec{T}}$ is the time-derivate of the vector of node temperatures, \vec{F} is the input power select matrix that maps the power source vector \vec{P} onto the nodes. While the ODE system in Equation (22) can be directly solved in circuit simulation software such as SPICE and Hotspot [72], its dimension is proportional to the number of nodes which makes it poorly suited for runtime applications due to high computational requirements. To improve the computational efficiency, several model order reduction methods are developed to transform the high-dimensional system to a low-dimensional one for a faster calculation [73], [74]. Figure 39 shows a schematic of a representative thermal network model.

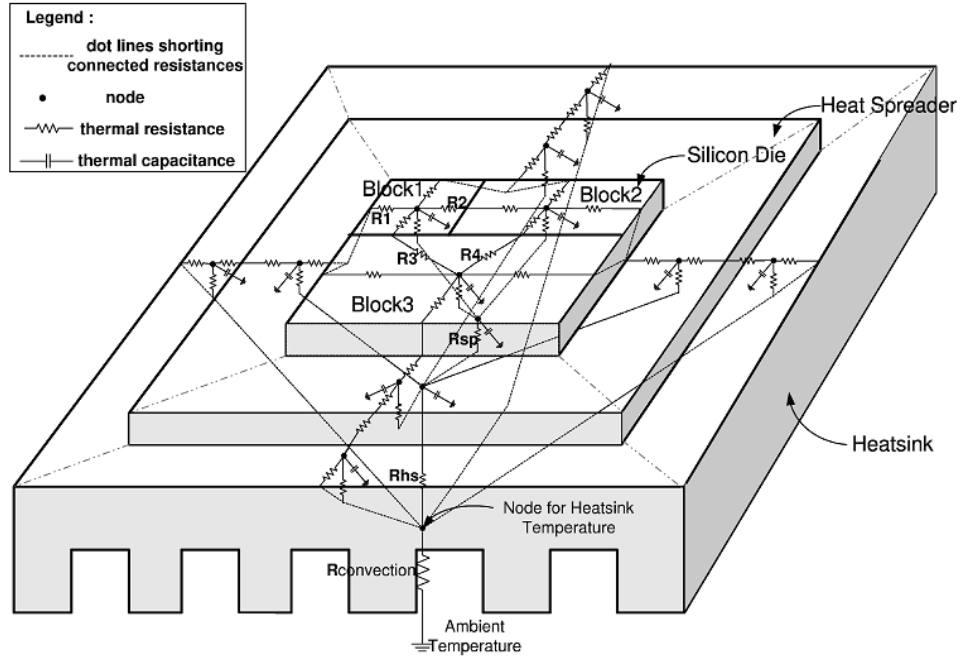


Figure 39: Example of a network circuit model of a chip die on a heat spreader attached to a heat sink. This type of model is too computationally intensive for runtime implementation [66].

An alternate approach using an RC ladder model for the location of interest is formulated for better computational efficiency. The thermal response of the system, subjected to a step-function power pulse, is recorded with appropriately resolved timescales. A suitable extraction technique is then used to define an RC ladder model with an equivalent response [75]. The needed step-function response $Z(t)$ can be modeled numerically or obtained directly from a measurement. Assuming model linearity, the temperature response $T(t)$ to an arbitrary power trace input $P(t)$ can be computed by the convolution integral between input power and the time derivative of $Z(t)$:

$$T(t) = T_0 + \int_0^t P(\tau) \cdot \dot{Z}(t - \tau) d\tau \quad (23)$$

where T_0 is the temperature at $t = 0$.

The previously mentioned RC network model can be used for describing a multi-input-multi-output (MIMO) thermal system. To do so requires a spatial discretization of the governing equations over a complete model domain to obtain a dynamic thermal response in Equation (23). The model complexity often requires extensive linear algebra manipulations for reducing the number of unknowns in the studied system. A model order reduction method must be carried out with caution in order to ensure numerical stability. In addition, the RC network approach relies on the discrete numerical models representing the thermal system and cannot directly be based on the experimental results. In contrast, the RC ladder approach can take either simulation or experimental input for model formulation. The RC ladder model resolves only a single

conduction path; for modeling a MIMO system such as a chip with multiple hotspots, the results of several models must be linearly superimposed. Therefore, the RC ladder approach is more appropriate for thermal systems with a limited number of points of interest, such as hotspots. The technique described in the present study employ an RC ladder model.

Within the category of RC ladder models, there are two types of models to consider: Foster ladder models and Cauer ladder models, as shown in Figure 40. Cauer ladder models provide a better physical description the heat flow path in the system, while Foster ladder models only capture the thermal behavior but have no physical equivalent. Use of the Cauer-ladder network is not straightforward due to its complicated mathematical representation. Fortunately, a Foster-ladder network can be easily be transformed to a Cauer ladder network which provides the same step-function thermal response $Z(t)$. For this reason, step-function responses are typically characterized by the Foster-type RC ladder network, and such an approach is taken in this paper.

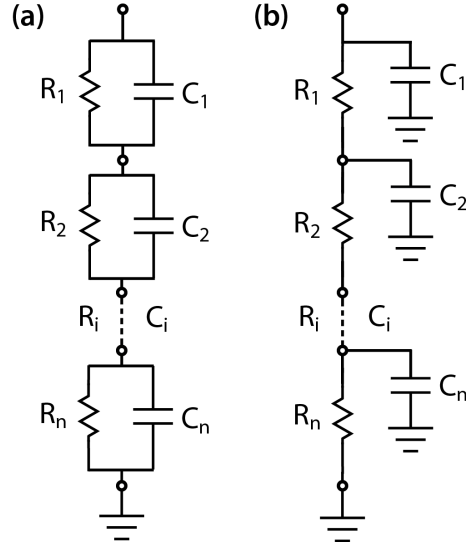


Figure 40: Foster and Cauer RC ladder network representation of thermal system.

To implement an RC ladder model for the system, first a Foster RC ladder network must be derived to represent the thermal response of the system; the second challenge – and indeed the primary contribution of this paper – is to develop a technique for bypassing the direct evaluation of convolution integral in Equation (23). In the next section, a brief review of existing methods for deriving the Foster RC ladder and evaluating the convolution integral is presented and their limitations are discussed. We then propose an improved method for calculating temperature evolution with arbitrary power traces which bypasses direct evaluation of the convolution integral. It employs a recursive infinite impulse response (IIR) digital filter on sampled power traces, which is a well-developed and commonly used approach in digital signal processing community. A derivation of the IIR digital filter coefficients based on the parameterized thermal RC ladder network is presented. Section 4.3 discusses model validation results and demonstrates best achievable scaling of runtime computations

with the required execution time by comparing the proposed technique to the existing convolution methods.

4.2. Thermal Modeling Approach

4.2.1. Determining the Step Response Function of the Foster RC ladder Network

Using a Foster RC ladder model, the time dependent thermal impedance can be written as:

$$Z(t) = \sum_{i=1}^n R_i \left(1 - e^{-\frac{t}{R_i C_i}} \right) \quad (24)$$

where R_i (K/W) and C_i (J/K) form the i^{th} stage of RC ladder network (Figure 40).

Various methods have been developed to determine the discrete elements in a lumped ladder network shown in Figure 40, e.g., through least square fitting to the simulated/measured heating curves in time-domain [76] and frequency-domain [77]. A method preferred in this work was proposed by Szekely and Van Bien [78] and is based on computing time-constant spectrum of distributed network from measured thermal transient response in time domain. Recently, for an experimental implementation, this method has been extended by applying similar identification procedure on measured impulse response spectrum in frequency-domain [79]. For a typical electronic power device, a very accurate approximation of $Z(t)$ can be achieved with less than 10 pairs of R 's and C 's in Equation (24). As previously stated, the Foster RC ladder can be converted to an equivalent Cauer RC ladder to provide increased

physical insight, if desired. A summary of this technique with added observations follows next.

The step response for unit power can be generalized by considering continuous time-constant spectrum

$$Z(t) = \int_0^{\infty} R(\tau) \left(1 - e^{-\frac{t}{\tau}}\right) d\tau \quad (25)$$

Introducing variables

$$z = \ln(t); \psi = \ln(\omega); \xi = \ln(\tau) \quad (26)$$

one obtains

$$Z(z) = \int_{-\infty}^{\infty} R(\xi) \left(1 - e^{-e^{z-\xi}}\right) d\xi \quad (27)$$

$$Q_t(z) = \frac{dZ(z)}{dz} = \int_{-\infty}^{\infty} R(\xi) e^{z-\xi-e^{z-\xi}} d\xi = R(z) \otimes e^{z-e^z} \quad (28)$$

Similarly, in frequency domain, the following expressions are found

$$Q_r(\psi) = -\frac{d \operatorname{Re}[\theta(\psi)]}{d\psi} = R(-\psi) \otimes w_r(\psi) \quad (29)$$

$$Q_i(\psi) = -\operatorname{Im}[\theta(\psi)] = R(-\psi) \otimes w_i(\psi) \quad (30)$$

Expressions (28), (29), and (30) are convolutions of the desired system response with the functions:

$$w_t(x) = e^{x-e^x} \quad (31)$$

$$w_r(x) = \frac{2 e^{2x}}{(1 + e^{2x})^2} \quad (32)$$

$$w_i(x) = \frac{e^x}{1 + e^{2x}} \quad (33)$$

Szekely [80] discussed the possibility of deconvolving the time-constant spectrum from the measured or calculated system responses. Since numerical Fourier transforms of functions Equations (31), (32), and (33) were found to impact the deconvolution accuracy, we instead use exact expressions for the spectrum of these functions:

$$W_t(k) = \Gamma(1 - 2 i \pi k) \quad (34)$$

$$W_r(k) = \pi^2 k \operatorname{csch}(\pi^2 k) \quad (35)$$

$$W_i(k) = \frac{\pi}{2} \operatorname{sech}(\pi^2 k) \quad (36)$$

where k is the variable in the Fourier domain complementary to the logarithm of the time constant, the gamma function Γ (analytically continued into a complex plane) is

used for the time-domain deconvolution method, and hyperbolic functions are used for deconvolution using the associated frequency-domain techniques.

To reveal the computational aspects of NID methodology, it is useful to examine a system consisting of one single time constant. In this case, a transform of a Q_t function, given by $R_0 \times e^{z-\xi_0-e^{z-\xi_0}}$, can be written out as

$$Q(k, \xi_0) = R_0 \times e^{-2\pi k i \xi_0} \times W_t(k) \quad (37)$$

Dividing this expression by W_t in principle yields a spectrum of a single pole located at ξ_0 . In practice, W_t attenuates sharply with increasing k . If the transform is calculated numerically, as opposed to the analytical form in Equation (37), the deconvolution accuracy will be limited by errors in calculating $Q(k, \xi_0)$. The errors with magnitudes comparable to $|R_0 \times W_t(k)|$ at sufficiently high k -values will render these components unusable. Filtering this part of the spectrum is then necessary, which in turns leads to broadening of an identified peak.

The time constant spectrum can then be found by the inverse transform:

$$R_{t,r,i} \begin{pmatrix} z_t \\ -z_r \\ -z_i \end{pmatrix} = F^{-1} \left[\frac{G_{t,r,i}(k)}{W_{t,r,i}(k)} F(Q_{t,r,i}) \right] \quad (38)$$

where F and F^{-1} is the Fourier transform pair and $G_{t,r,i}$ is the appropriately chosen filter suppressing discretization and numerical round-off errors at high k values.

4.2.2. Existing Methods for Computing the Convolution Integral

The convolution integral can be calculated using fast Fourier transforms (FFT) by simple multiplication in frequency domain of the calculated power spectrum and the known response followed by inverse operation. This approach leads to difficulties with computational efficiency since a separate transform is required for every time step.

An alternate approach is to directly compute the convolution integral in time domain. With a piecewise linear approximation of $P(t)$ that is sufficiently accurate for small time intervals, a semi-analytical formula of Equation (23) is obtained in [81] by applying the superposition principle:

$$T(t) = \sum_{i=1}^n R_i \cdot \sum_{j=1}^m \left\{ [P_j + a_j \cdot (t_{j+1} - t_j - \tau_i)] \times e^{-\frac{t-t_{j+1}}{\tau_i}} - [P_j - a_j \cdot \tau_i] \times e^{-\frac{t-t_j}{\tau_i}} \right\} \quad (39)$$

where:

$$a_j = \frac{P_{j+1} - P_j}{t_{j+1} - t_j} \quad (40)$$

For a more general case of a piecewise constant power input P_j between t_j and t_{j+1} :

$$T(t) = \sum_{j=1}^m P_j [Z(t - t_j) - Z(t - t_{j+1})] \quad (41)$$

with $t = t_{m+1}$ and which does not require the model fitting of $Z(t)$. A similar approach with interpolated $Z(t)$ is also proposed in [82].

4.2.3. *Recursive Digital Filtering Technique for Computing the Convolution*

Integral

A technique for computing the convolution integral, presented in this work, is developed by recognizing the fact that a transfer function of a digital filter can be constructed to approximate a response of a modeled thermal system due to its linear time-invariant (LTI) property. A model reduction is based on Foster RC ladder network shown in Figure 40 as a series of n -stage first-order low-pass filters in series, acting on input power and outputting temperature responses in the continuous time domain. These filters in the continuous time domain can be transformed to digital format in the discrete time domain, which would then be applied to the sampled input power to achieve the desired temperature response.

The use of a digital filter instead of a series of continuous-time filters is advantageous for numerous reasons: it is easy to design and implement; it can handle large dynamic range; it has extremely stable performance; and it is programmable to adapt to input signals. The general representation of digital filter in discrete time domain is the difference equation:

$$y(n) = \sum_{i=0}^M b_i \cdot x(n-i) - \sum_{i=0}^N a_i \cdot y(n-i) \quad (42)$$

where $x(n)$ is the input signal at instant n , $y(n)$ is the output signal at instant n , and constants $a_i, i = 0, 1 \dots N$, and $b_j, j = 0, 1 \dots M$, are feedback and feedforward coefficients, respectively. Infinite impulse response (IIR) filters have non-zero feedback coefficients ($a_i \neq 0$), as opposed to the finite impulse response (FIR) filters ($a_i = 0$). The IIR filter has an impulse response that is non-zero over an infinite length of time, which is desirable property for modeling physical thermal response.

For a system represented by the Foster RC ladder network, the impulse response function $\dot{Z}(t)$ has properties of an IIR filter due to its exponentially decaying terms. To determine the feedback a_i and feedforward b_j coefficients for the difference equation, the continuous transfer function $H(s)$ of the analog multi-stage filter is transformed into the discrete transfer function $H(z)$ of its approximate IIR digital equivalent. An inverse z-transform is then applied to calculate filter coefficients for the difference equation. There are numerous available methods to transform from $H(s)$ to $H(z)$; in this work, bilinear transformation method is preferred due to the absence of frequency aliasing distortions.

In the continuous frequency domain, the complex impedance of the Foster thermal RC ladder network is given by:

$$H_{th}(s) = \sum_{i=1}^K \frac{R_{th,i}}{1 + s \cdot \tau_i}; \quad \tau_i = R_i C_i \quad (43)$$

where s is the complex frequency, and K is the number of stages in RC ladder network. A bilinear transformation is carried out by performing the substitution of s in $H_{th}(s)$ with:

$$s = \frac{2}{\Delta t} \frac{z - 1}{z + 1} \quad (44)$$

where Δt is the sampling interval. Using inverse z -transform one then obtains the desired difference equation:

$$T(n) = \sum_{i=1}^K T_i(n) \quad (45)$$

$$= \sum_{i=1}^K \left\{ \frac{2\tau_i - \Delta t}{2\tau_i + \Delta t} T_i(n-1) + \frac{R_{th,i} \Delta t}{2\tau_i + \Delta t} [P(n) + P(n-1)] \right\}$$

where $T(n)$ and $P(n)$ are the discrete temperature output and power input at time $n \times \Delta t$. As readily seen, this method bypasses the direct convolution of the integral in Equation (23), and instead recursively calculates the transient temperature response at any time step by using the temperature output at the previous time step while applying trapezoidal rule for integrating power input within the time interval.

The derivation presented here provides a model for a system subjected to a single heat source. For a thermal system subjected to multiple heat sources, the transient temperature response at any location is the superposition of the responses from multiple power excitations. Due to non-linear frequency mapping, the bilinear

transformation method cause frequency warping in $H_{th}(z)$ that does not preserve the frequency characteristics of corresponding $H_{th}(s)$. This effect can be either eliminated by pre-warping before transformation, or minimized by reducing sampling interval Δt according to the following condition:

$$\frac{\Delta t}{\tau_{min}} \approx \tan\left(\frac{\Delta t}{\tau_{min}}\right) \quad (46)$$

where τ_{min} is the minimum thermal time constant of the stages in Foster RC ladder network. For $\Delta t/\tau_{min} = 0.2$, the amount of frequency warping is approximately 1.4% and has negligible effect on transient response in the time domain.

To summarize, the general procedure for obtaining digital filter coefficients followed in this work involves identification of time constant spectrum, its discretization into Foster's network, mapping of a transfer function into z-domain and finally an inverse z-transform to obtain feedback coefficients for the part of the filter acting recursively on previous temperature history and feed-forward coefficients for the part of the filter acting on power input into the identified Foster's network.

4.3. Model Verification and Applications

The technique is verified against the analytical solution for one-dimensional conduction in a Cartesian geometry. The solution in complex frequency domain for a semi-infinite geometry yields:

$$\theta(s) = \frac{1}{\gamma \sqrt{s/\alpha}} \quad (47)$$

where $\theta(s)$ is the temperature response per unit heat flux, γ and α are thermal conductivity and diffusivity, respectively. For the finite geometry restricted to a slab of thickness L , with zero temperature boundary condition at the side opposing to the entering heat flux, the solution is:

$$\theta(s) = \frac{\tanh(L\sqrt{s/\alpha})}{\gamma \sqrt{s/\alpha}} \quad (48)$$

The temperature response at a time constant τ can be directly calculated using real negative axis in s -plane to obtain the spectrum as [83]:

$$R(\tau) = \frac{1}{\pi} \text{Im } \theta(-\tau^{-1}) \quad (49)$$

The time-constant spectrum for the distributed case of semi-infinite media is then

$$R(\tau) = \frac{\sqrt{\alpha \tau}}{\pi \gamma} \quad (50)$$

while for the case with restricted geometry the spectrum is discontinuous,

$$\tau_n = \frac{4}{(\pi n)^2} \frac{L^2}{\alpha} \quad (51)$$

$$R_n = \frac{8}{(\pi n)^2} \frac{L}{\gamma} \quad (52)$$

for $n = 1, 3, 5 \dots$

with zero response at all other time constant values besides τ_n listed above.

The solutions in time domain for the slab geometry is [84]:

$$T(t, L) = \frac{2 \sqrt{\alpha t}}{\gamma} \sum_{n=0}^{\infty} \left\{ (-1)^n \left[\operatorname{ierfc} \frac{n L}{\sqrt{\alpha t}} - \operatorname{ierfc} \frac{(n+1)L}{\sqrt{\alpha t}} \right] \right\} \quad (53)$$

Figure 41 shows time constant spectrum identified by the NID procedure with all three deconvolution methods using numerical Fourier transforms in the Equation (38). The thermal properties are that of silicon ($\gamma=150 \text{ W m}^{-1} \text{ K}^{-1}$, $\alpha=8.47 \cdot 10^{-5} \text{ m}^2 \text{ s}^{-1}$) and thickness of the slab is $L=15 \times 10^{-3} \text{ m}$. All extracted time constants collapse onto one curve due to the use of an identical Gaussian filter function in the complementary domain. First several system poles, given by Equations (51) and (52) are identified rather accurately by the deconvolution procedure, as shown further in Figure 42. The identified response converges to the distributed semi-infinite limit given by Equation (50); the time constants at these values can be lumped into Foster network by integrating sections of a continuous spectrum.

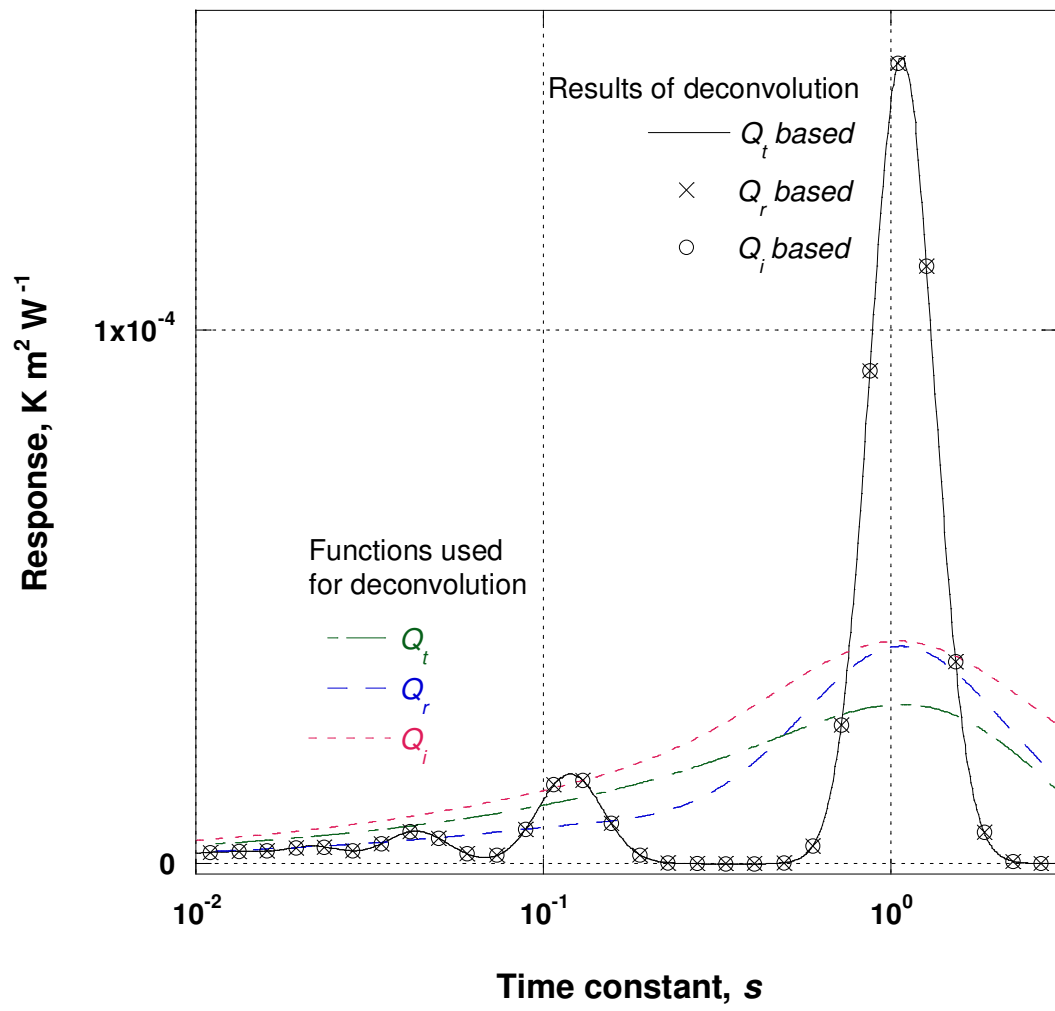


Figure 41: Network deconvolution using responses in time and frequency.

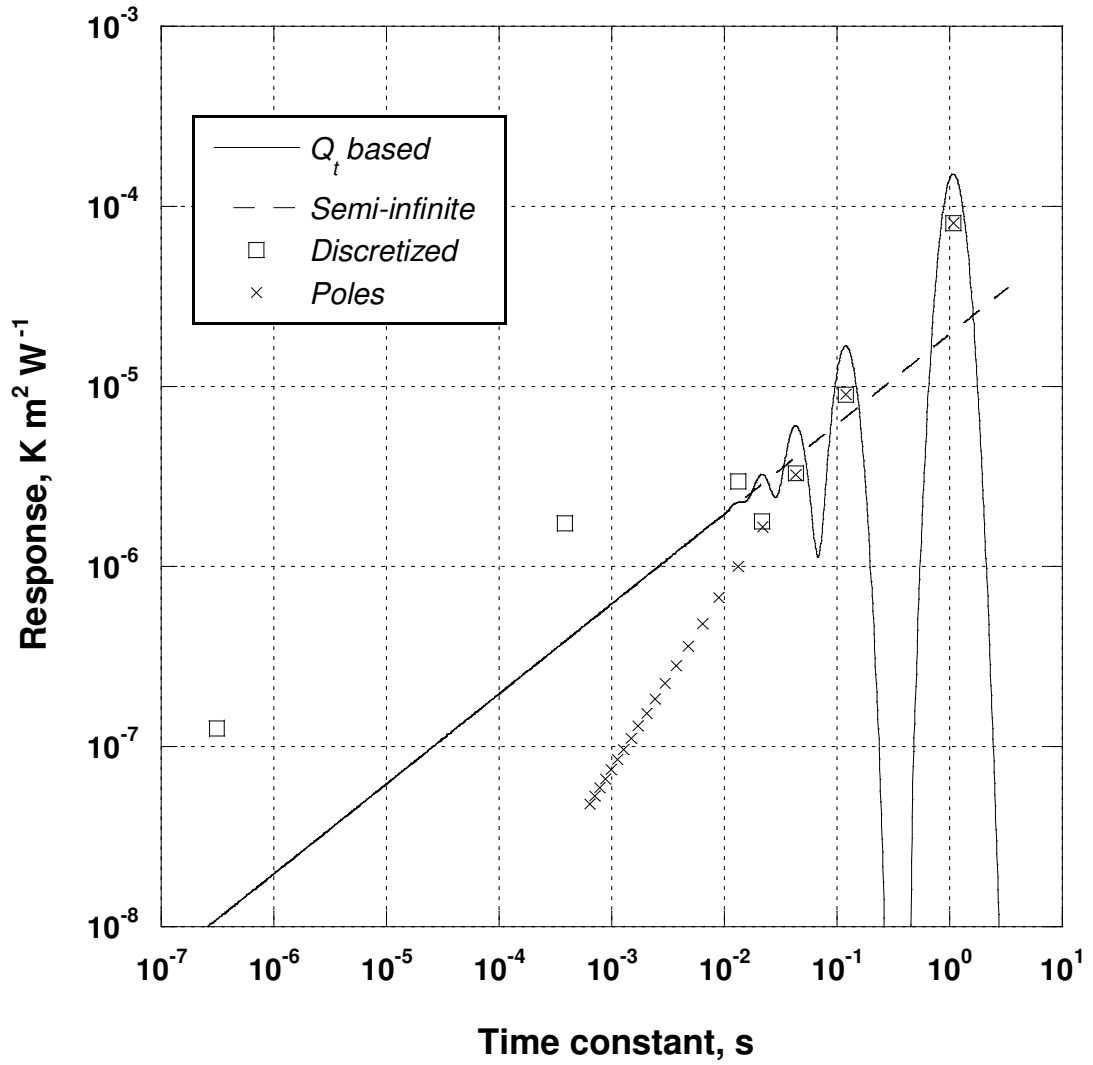


Figure 42: Identification of system poles and semi-infinite limit.

Figure 43 shows transfer functions constructed from the identified network. Using full spectrum of identified response provides near-identical match with the exact form given by Equation (48) at $s = i \times \omega$. A network with seven or more elements provides a good approximation to the system thermal behavior at frequencies with relevant amplitude responses. Figure 44 shows response of the constructed IIR filter to step in power for seven-stage identified network with varying time steps. The frequency

warping becomes pronounced at large time steps but the filter output is stable and always converges to the expected steady state value.

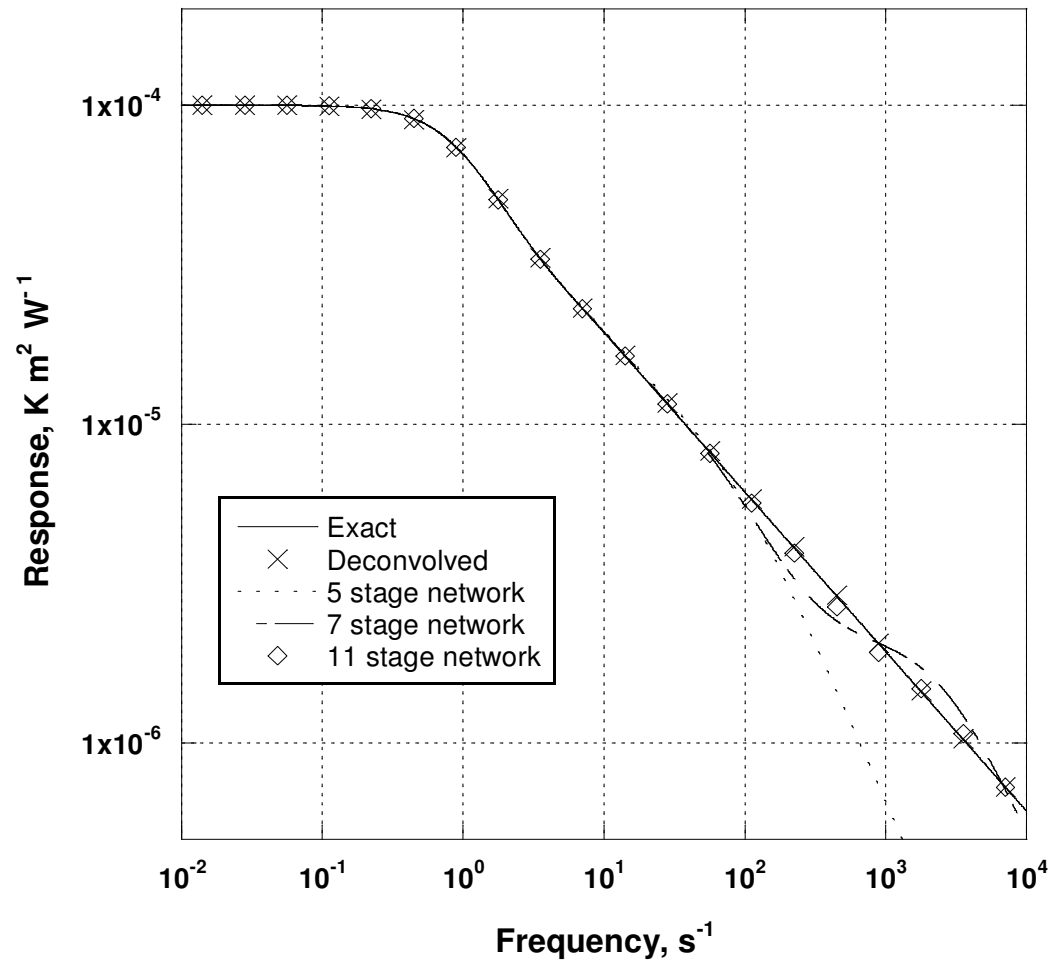


Figure 43: Transfer function identification.

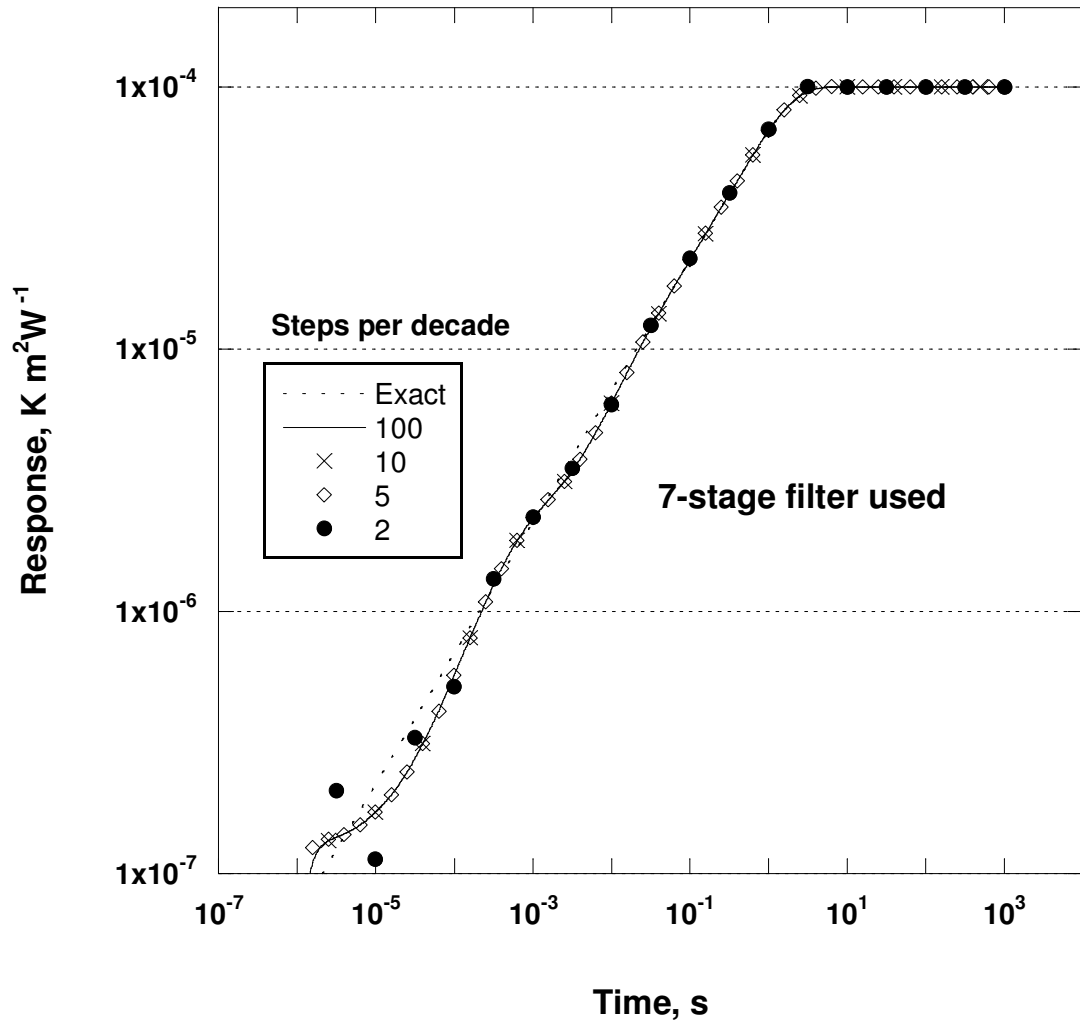


Figure 44: IIR Response for varying inputted time discretization.

Figure 45 depicts a schematic of a common thermal system to demonstrate the proposed model generation using a commercially available solver. A representative two-layer configuration consists of a chip with properties of silicon and a heat sink with properties of copper subjected to forced convective cooling. A thickness of $10 \text{ mm} \times 10 \text{ mm}$ silicon chip is 0.5 mm and a thickness of $30 \text{ mm} \times 30 \text{ mm}$ copper spreader is 1 mm , with a uniform convection coefficient of $10^4 \text{ W m}^{-2} \text{ K}^{-1}$. A uniform power of 100 W was applied at the top surface of the silicon layer. The representative

system was modeled using COMSOL[®] Multiphysics software. Three different mesh sizes ranged from approximately 1,100 to 10,500 mesh points. For all cases, the results are self-consistent.

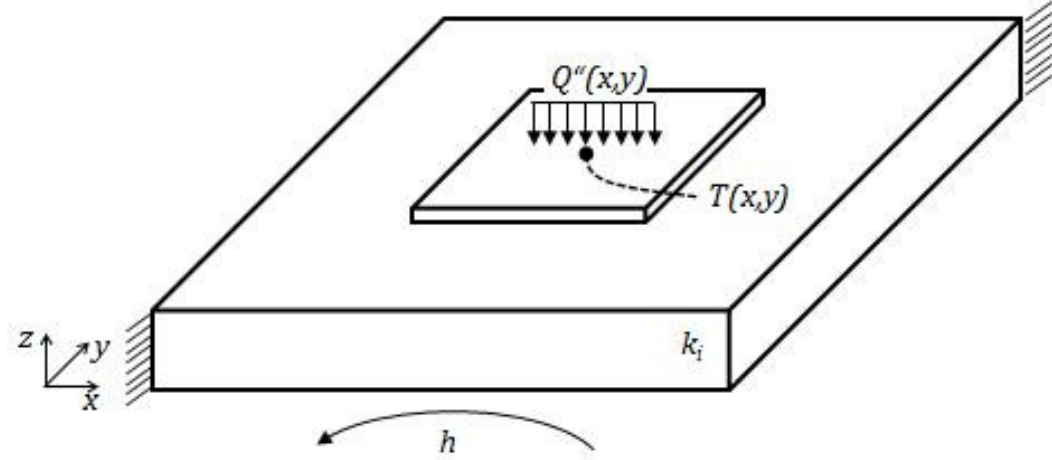


Figure 45: Schematic of chip system used for numerical simulation and proposed modeling method.

Figure 46 shows results of the network extraction using power step response, as discussed earlier, with logarithmically spaced time steps. The IIR filter, generated with 11-element identified network, was used to process the same power step at identical time samples. The agreement is excellent with less than 0.4 % maximum transient error with respect to the steady-state response value. Figure 47 shows a good agreement achieved between the simulation results with linearly-spaced steps and the filter output subjected to the same sequence of power steps.

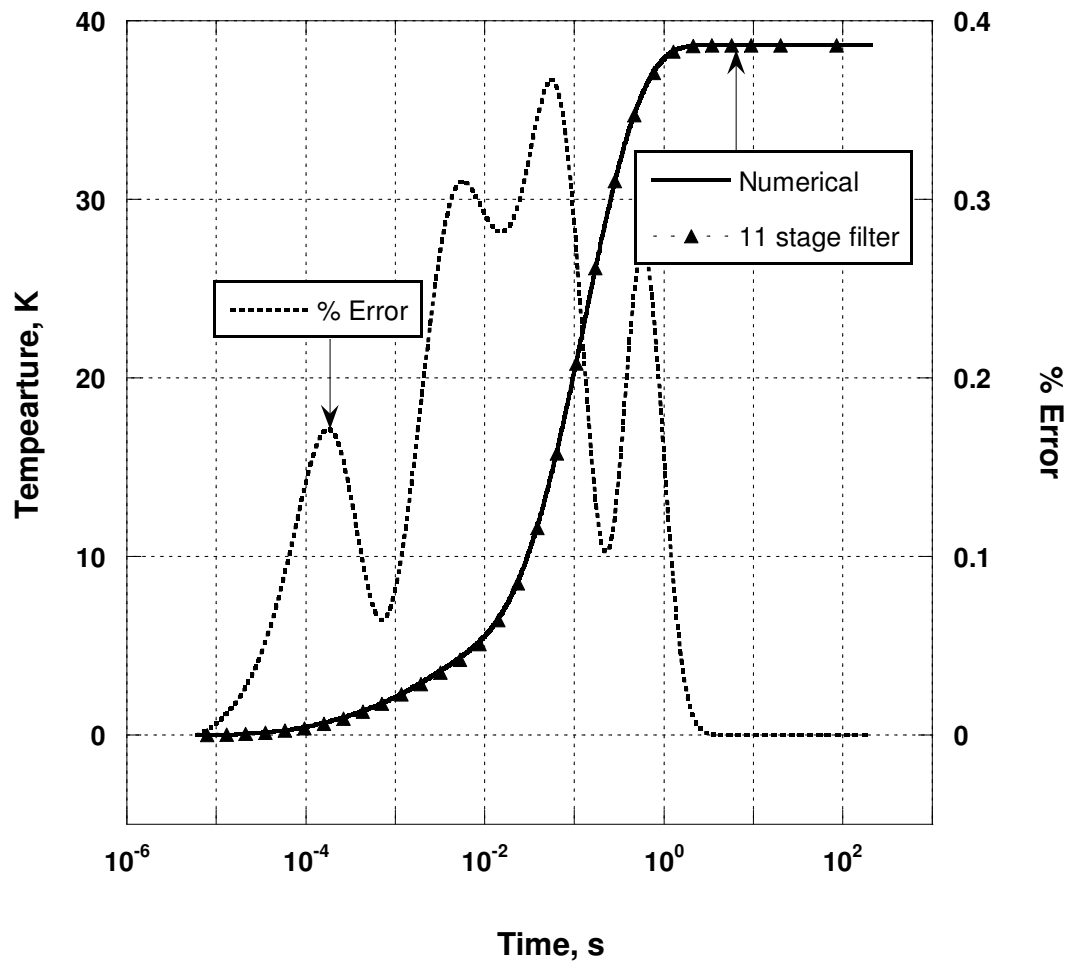


Figure 46: Thermal modeling of chip-spreader geometry shown on Figure 45. Compared are the results of simulations using commercial solver with the output of an IIR filter based on 11 stage network. The power step is at 100 W. The maximum transient errors are less than 0.4% of the steady state response.

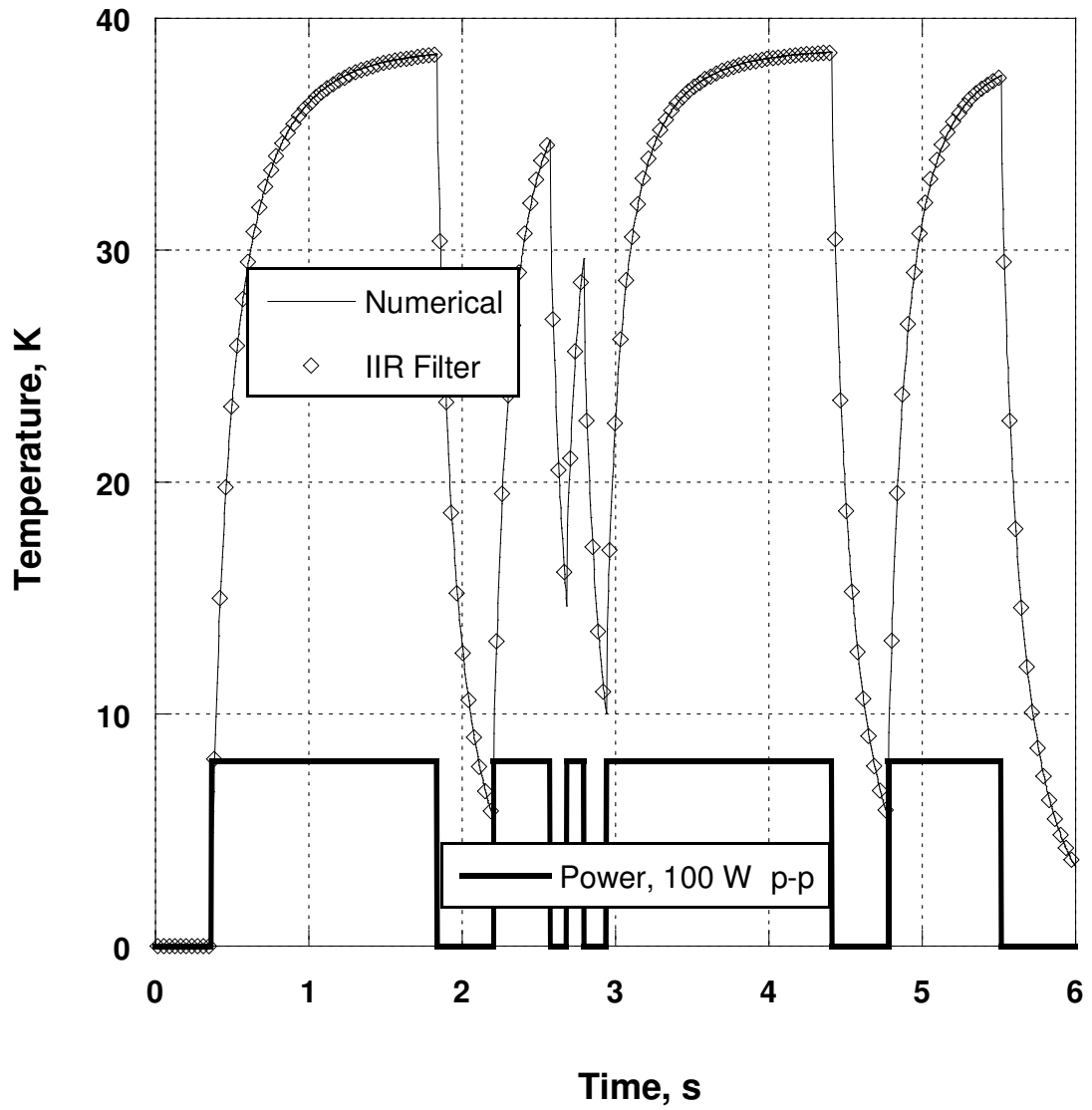


Figure 47: Comparison between numerical model and IIR digital filter output subject to square-wave input power excitations.

Figure 48 compares execution times of the IIR filter with that of the existing techniques. The widths of the summation intervals in Equations (39) and (41) are proportional to the time step index m . This results in linear with number of past time steps demand for compute power for each new time increment, which in turn increases the overall computation load proportionally to m^2 . Memory requirements for the

storage of time and power traces also become a consideration. For the reasons stated, the existing algorithms do not provide adequately efficient calculations of runtime computational temperature given real-time transient power input. These inefficiencies make implementations particularly difficult in most embedded platforms.

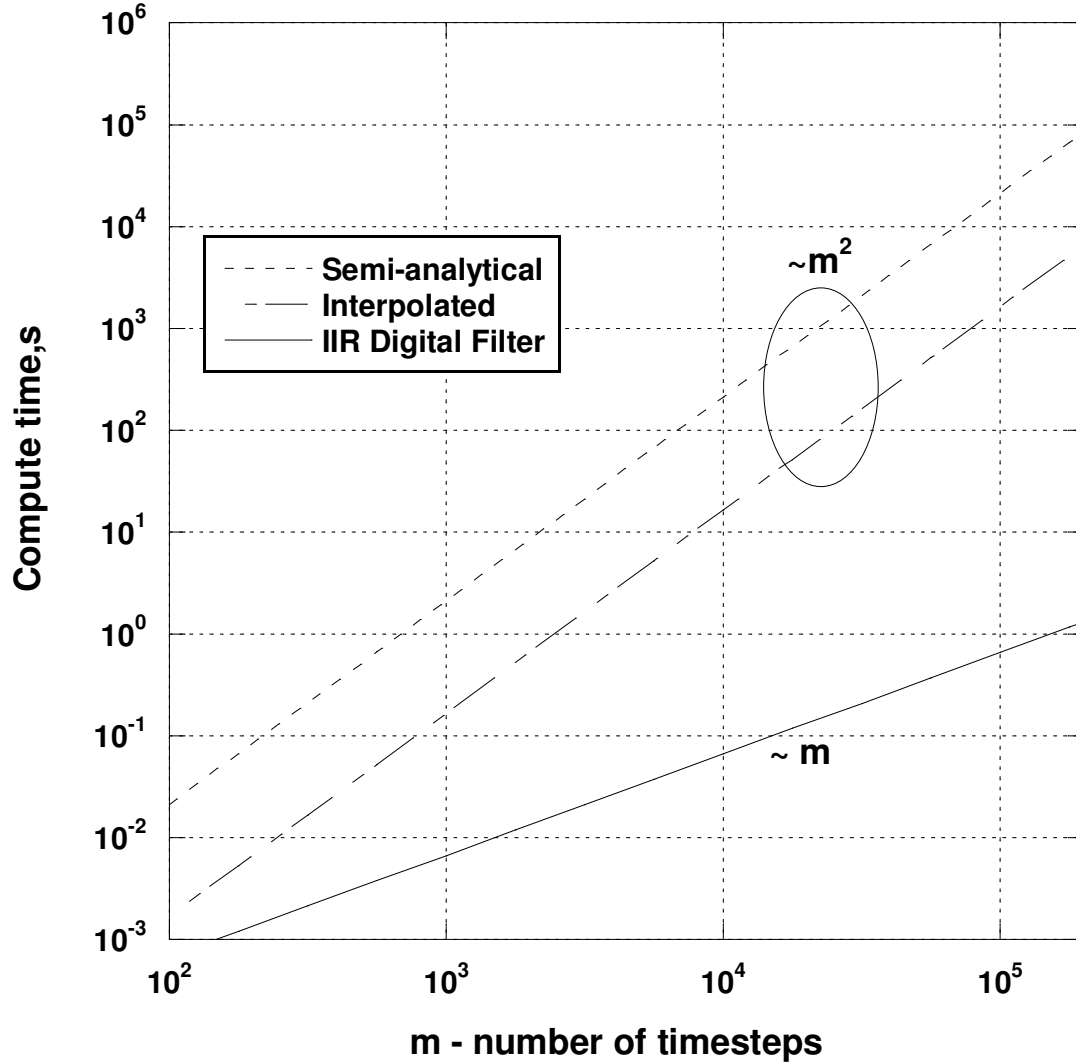


Figure 48: Comparison in computational efficiency of different methods for evaluation of convolution integrals. The recursive IIR digital filter is superior to existing convolution methods and achieves best possible scaling due to its constant computation overhead for each time step.

The IIR digital filter based technique provides the best possible scaling, linear scaling, for the overall computational load. Each new time step results in exactly the same number of processing operations as the previous. Additionally, memory requirements are much less demanding compared to existing techniques, making the proposed methodology ideally-suited for runtime temperature calculations.

4.4. Summary

This work presents a novel approach for predicting temperature evolution in electronic devices subjected to transient heat sources. It is based on modeling dynamic behavior of a thermal system with an identified network. We revisit the model reduction by network identification (NID) and present an extension of a method to obtain time-constant spectrum of a thermal network based on analytical form of convolving functions, while providing new insights to limitations of the technique. We verify the model extraction procedure using analytical solution and demonstrate correct identification of known system poles and convergence of the extracted time constant spectrum to the limiting case.

We then present IIR digital filters suited for run-time evaluation of convolution integral in discrete time-domain. A simple formulation of recursive digital filters makes the algorithm well-suited for run-time temperature predictions. The resulting recursive algorithm yields temperature calculation at a given time instant using very limited depth of recorded temperature history. A numerical model of semiconductor device is created to generate time-domain temperature responses to step-function

power excitation; excellent accuracy of the filter output is confirmed when compared to simulations.

Comparison with conventional integral-based convolution methods also indicates a dramatic improvement in computational efficiency compared to existing techniques.

The achieved scaling is best possible, linear, with the number of temperature evaluations, a feature enabled by the use of a DSP technique. This improvement allows implementation of sophisticated runtime dynamic thermal management algorithms for all high-power architectures and expands the application range to embedded platforms for implementations in a pervasive computing environment.

CHAPTER 5: CONCLUSIONS

Neither chip cooling nor dynamic thermal management will be able to address microprocessor thermal management challenges alone. Advances in cooling, particularly for high performance microprocessors, are required to push the power limits of future generation microprocessors. Low-cost, small-form factor cooling advancements will be needed for mobile applications. While developments on cooling are critical, it is unlikely they will resolve the hotspot challenges that occur over the various length and time scales of the processor. At the package length scale, minimizing temperature non-uniformities is critical for mitigating thermo-mechanical stress on the chip and package. At the length scales of the active processor regions, temperature variations in space and time can lead to local temperature excursions above critical temperatures posing system failure risks. The temperature reported from any particular thermal sensor does not capture the temperature of neighboring regions, which may be higher or lower than the measured location. Similarly, rapid transient fluctuations in temperature due to processor activity must be accounted for when specifying reliable operating temperatures. By throttling chip power in response to thermal signals, dynamic thermal management (DTM) can directly address these multi-scale spatial and temporal temperature fluctuations. Overly conservative DTM schemes, however, cause unnecessary computational performance degradation; DTM techniques must be optimized for chip thermal management and chip performance. Understanding the magnitude and effects of measurement uncertainty in DTM

schemes and developing techniques for uncertainty reduction is central to the task of designing robust, efficient DTM schemes.

The present work offers novel computational tools for characterizing and minimizing hotspot detection uncertainty and predicting transient hotspot response. A novel inverse heat transfer solution is introduced that leverages analytical, spatial-frequency domain analysis of heat transfer in a chip. The solution technique is implemented across randomized chip heat flux profiles to demonstrate the generalized limits of hotspot detection using spatially discretized thermal signal from either laboratory thermography or on-chip sensors. Under particular test conditions, the inverse technique reduces the normalized mean absolute error in the calculated heat flux by as much as 30% as compared to direct interpretation of the thermal data. Parametric studies of sensor vertical proximity, sensor measurement error, convective boundary conditions, and chip thermal conductivity provide regime maps of performance improvements.

To address transient hotspot fluctuations and sensor placement limitations, an ultra-efficient transient model for chip hotspots is developed. The technique employs network identification deconvolution (NID) for characterization of chip thermal response using step-function response from either finite element modeling or experimental results. Based on the thermal response characterization, a digital signal processing technique is used to rapidly compute the chip response to an arbitrary transient power profile. An infinite impulse response (IIR) digital filter provides highly accurate results with the best possible computational scaling and reduced

memory requirements. This improvement in computational efficiency facilitates the on-chip, run-time prediction of transient hotspot behavior, enabling numerous possible improvements in dynamic thermal management schemes.

The improvements in hotspot detection and prediction offered by the present study contribute to an emerging field of research focused on the optimization of dynamic thermal management schemes. Research efforts to improve dynamic thermal management schemes include improved sensor design, high-resolution laboratory thermography, and sensor signal processing. Past research on thermal sensor design has yielded higher accuracy sensors, some of which demonstrate robustness to fabrication process variations. However, these designs remain too large for most applications. Advances in the design of smaller, accurate sensors that are largely independent of process variations would enable new levels of precision in DTM schemes. Smaller sensor size would allow integration of more sensors throughout the chip, both in the plane of the chip and also in the vertical dimension as manufacturers move to three dimensional integrated circuit (3D-IC) architectures. Control algorithms could also be designed more aggressively with more accurate sensors. Improvements in thermal sensing can be made by leveraging arrays of sensors. Research reviewed in this area showed computationally-efficient techniques for deconvolving spatially discretized thermal signals as well as networking techniques to minimize sensor signal traffic. Further work is needed to extend these techniques to 3D-IC applications where heat is dissipated from stacked components. Recent demonstrations of high-resolution laboratory thermography techniques, especially infrared (IR) and micro-Raman

thermography, were presented. While higher resolution thermal sensing techniques are known (e.g. scanning thermal microscopy), techniques that can be implemented with a high-heat flux cooling solution are of particular interest because they permit measurements with realistic chip traffic conditions. Though extremely challenging, the development of laboratory thermography techniques capable of resolving thermal profiles throughout a 3D-IC would be especially powerful for characterizing 3D-IC dynamic thermal management techniques. These new research directions offer exciting opportunities as state-of-the-art thermometry and signal processing are extended for next generation dynamic thermal management.

Well-established trends for microprocessor development indicate that the challenges ahead for thermal management are enormous. Both cooling and dynamic thermal management will need to play central roles in effective thermal management solutions. Improvements in hotspot detection and prediction, which constitute the central contribution of the present study, will be critical for managing the multi-scale, spatial and temporal temperature fluctuations in microprocessors. As these advances are made, unprecedented computational performance will become possible.

BIBLIOGRAPHY

- [1] R. Viswanath, V. Wakharkar, A. Watwe, and V. Lebonheur, "Thermal Performance Challenges from Silicon to Systems," *Intel Technology Journal*, vol. 4, no. 3, pp. 1-16, 2000.
- [2] D. Sylvester, "ElastIC: An Adaptive Self-Healing Architecture for Unpredictable Silicon," *IEEE Design & Test of Computers*, vol. 23, pp. 484-490, Nov. 2006.
- [3] M. Sabry, A. Sridhar, and D. Atienza, "Towards thermally-aware design of 3D MPSoCs with inter-tier cooling," in *Proceedings of the Design, Automation & Test in Europe Conference & Exhibition (DATE)*, 2011, pp. 1-6.
- [4] F. Alfieri, M. K. Tiwari, I. Zinovik, D. Poulikakos, T. Brunschweiler, and B. Michel, "3D Integrated Water Cooling of a Composite Multilayer Stack of Chips," *Journal of Heat Transfer*, vol. 132, no. 12, p. 121402, 2010.
- [5] J. H. Lau and T. G. Yue, "Thermal management of 3D IC integration with TSV (through silicon via)," in *Proceedings of the 59th Electronic Components and Technology Conference (ECTC)*, 2009, pp. 635-640.
- [6] S. Naffziger, "Foxton Technology." Intel Corp., HotChips, 2005.
- [7] E. Kursun, G. Reinman, S. Sair, A. Shayesteh, and T. Sherwood, "Low-overhead Core Swapping for Thermal Management," in *Power-Aware Computer Systems (PACS)*, 2004, pp. 46-60.
- [8] A. Cohen, F. Finkelstein, A. Mendelson, R. Ronen, and D. Rudoy, "On Estimating Optimal Performance of CPU Dynamic Thermal Management," *Computer Architecture Letters*, vol. 2, no. 1, pp. 6-6, 2003.
- [9] A. K. Coskun, R. Strong, D. M. Tullsen, and T. S. Rosing, "Evaluating the Impact of Job Scheduling and Power Management on Processor Lifetime for Chip Multiprocessors," in *Proceedings of SIGMETRICS/Performance*, 2009, pp. 169-180.
- [10] M. Goma, M. D. Powell, and T. Vijaykumar, "Heat-and-Run: Leveraging SMT and CMP to Manage Power Density through the Operating System," in *ACM SIGARCH Computer Architecture News*, 2004, vol. 32, no. 5, pp. 260-270.
- [11] S. Heo, K. Barr, and K. Asanovic, "Reducing Power Density through Activity Migration," in *Proceedings of the International Symposium on Low Power Electronics and Design (ISLPED)*, 2003, no. C, pp. 217-222.

- [12] R. McGowen et al., "Power and Temperature Control on a 90-nm Itanium Family Processor," *IEEE Journal of Solid-State Circuits*, vol. 41, no. 1, pp. 229-237, 2006.
- [13] D. L. Blackburn, "Temperature measurements of semiconductor devices - a review," *Twentieth Annual IEEE Semiconductor Thermal Measurement and Management Symposium*, pp. 70-80, 2004.
- [14] Y. Avenas, L. Dupont, and Z. Khatir, "Temperature Measurement of Power Semiconductor Devices by Thermo-Sensitive Electrical Parameters - A Review," *IEEE Transactions on Power Electronics*, no. 99, pp. 1-12, 2011.
- [15] A. Naderlinger, "A Survey of Dynamic Thermal Management and Power Consumption Estimation," in *Software Systems Seminar*, 2007.
- [16] F. Incropera, D. Dewitt, T. Bergman, and A. Lavine, *Introduction to Heat Transfer*, 5th ed. Hoboken, New Jersey: John Wiley & Sons, 2007, p. 842.
- [17] J. Kong, S. Chung, and K. Skadron, "Recent Thermal Management Techniques for Microprocessors," *Submitted to ACM Computing Surveys*, Nov. 2011.
- [18] R. Cochran and S. Reda, "Spectral Techniques for High-Resolution Thermal Characterization with Limited Sensor Data," in *Proceedings of the IEEE/ACM Design Automation Conference (DAC)*, 2009, pp. 478-483.
- [19] J. Dorsey et al., "An Integrated Quad-Core Opteron Processor," *2007 IEEE International SolidState Circuits Conference Digest of Technical Papers*, vol. 50, pp. 102-103, 2007.
- [20] V. Krinitzin, "Pentium 4 and Athlon XP: Thermal Conditions." [Online]. Available: <http://ixbtlabs.com/articles/pentium4athlonxpthermalmanagement/>. [Accessed: 17-May-2012].
- [21] Intel Corp, "Quad-Core Intel Xeon Processor 5400 Series Thermal/Mechanical Design Guidelines," 2007.
- [22] Intel Corp, "Quad-Core Intel Xeon Processor 5400 Series Datasheet," 2008.
- [23] S. Lopez-Buedo, J. Garrido, and E. I. Boemo, "Dynamically inserting, operating, and eliminating thermal sensors of FPGA-based systems," *IEEE Transactions on Components and Packaging Technologies*, vol. 25, no. 4, pp. 561-566, Dec. 2002.

- [24] S. Velusamy, W. Huang, J. Lach, M. Stan, and K. Skadron, "Monitoring temperature in FPGA based SoCs," in *Proceedings of the International Conference on Computer Design*, 2005, pp. 634-637.
- [25] S. Mondal, R. Mukherjee, and S. O. Memik, "Fine-Grain Thermal Profiling and Sensor Insertion for FPGAs," in *Proceedings of the IEEE International Symposium on Circuits and Systems*, 2006, pp. 4387-4390.
- [26] R. Mukherjee and S. O. Memik, "Systematic Temperature Sensor Allocation and Placement for Microprocessors," in *Proceedings of the IEEE/ACM Design Automation Conference (DAC)*, 2006, vol. 184.
- [27] E. Aldrete-Vidrio et al., "Strategies for built-in characterization testing and performance monitoring of analog RF circuits with temperature measurements," *Measurement Science and Technology*, vol. 21, no. 7, p. 075104, Jul. 2010.
- [28] P. Bratek and A. Kos, "Temperature sensors placement strategy for fault diagnosis in integrated circuits," in *Symposium on Semiconductor Thermal Measurement and Management*, 2001, pp. 245-251.
- [29] C. Yao, K. K. Saluja, and P. Ramanathan, "Calibrating On-chip Thermal Sensors in Integrated Circuits: A Design-for-Calibration Approach," *Journal of Electronic Testing*, vol. 27, no. 6, pp. 711-721, Sep. 2011.
- [30] "Calibrating the Thermal Assist Unit in the IBM25PPC750L Processors," vol. IBM PowerP. pp. 1-10, 2001.
- [31] E. Schlaepfer, "External Temperature Sensor Calibration for the MAX16031/MAX16032 System Monitors." Maxim Application Note 4284, pp. 1-7, 2008.
- [32] S. Kaxiras and P. Xekalakis, "4T-Decay Sensors: A New Class of Small, Fast, Robust, and Low-Power, Temperature/Leakage Sensors," in *Proceedings of the International Symposium on Low Power Electronics and Design (ISLPED)*, 2004, pp. 108-113.
- [33] J. Long, S. O. Memik, G. Memik, and R. Mukherjee, "Thermal monitoring mechanisms for chip multiprocessors," *ACM Transactions on Architecture and Code Optimization*, vol. 5, no. 2, pp. 1-33, Aug. 2008.
- [34] E. Aldrete-Vidrio, D. Mateo, and J. Altet, "Differential Temperature Sensors Fully Compatible With a 0.35-um CMOS Process," *IEEE Transactions on Components and Packaging Technologies*, vol. 30, no. 4, pp. 618-626, 2007.

- [35] D. Barlini, M. Ciappa, M. Mermet-guyennet, and W. Fichtner, "Measurement of the transient junction temperature in MOSFET devices under operating conditions," *Microelectronics Reliability*, vol. 47, no. 9–11, pp. 1707-1712, Sep. 2007.
- [36] Q. Chen, M. Meterelliyo, and K. Roy, "A CMOS thermal sensor and its applications in temperature adaptive design," in *Proceedings of the 7th International Symposium on Quality Electronic Design (ISQED)*, 2006.
- [37] S. Remarsu and S. Kundu, "On process variation tolerant low cost thermal sensor design in 32nm CMOS technology," in *Proceedings of the 19th ACM Great Lakes Symposium on VLSI (GLSVLSI)*, 2009, p. 487.
- [38] K. Bharath, C. Yao, N. S. Kim, P. Ramanathan, and K. K. Saluja, "A Low Cost Approach to Calibrate On-Chip Thermal Sensors," in *12th International Symposium on Quality Electronic Design*, 2011, pp. 1-5.
- [39] T. E. Salem, D. Ibitayo, and B. R. Geil, "Validation of Infrared Camera Thermal Measurements on High-Voltage Power Electronic Components," *IEEE Transactions on Instrumentation and Measurement*, vol. 56, no. 5, pp. 1973-1978, 2007.
- [40] L. Hom, A. Durieux, J. Miler, M. Asheghi, K. Ramani, and K. E. Goodson, "Calibration Methodology for Interposing Liquid Coolants Infrared Thermography of Microprocessors," in *Proceedings of ITherm*, 2012.
- [41] A. Castellazzi, M. Honsberg-Riedl, and G. Wachutka, "Thermal characterisation of power devices during transient operation," *Microelectronics Journal*, vol. 37, no. 2, pp. 145-151, Feb. 2006.
- [42] M. Kuball et al., "Time-Resolved Temperature Measurement of AlGaIn/GaN Electronic Devices Using Micro-Raman Spectroscopy," *IEEE Electron Device Letters*, vol. 28, no. 2, pp. 86-89, 2007.
- [43] A. Sarua, A. Bullen, M. Haynes, and M. Kuball, "High-Resolution Raman Temperature Measurements in GaAs p-HEMT Multifinger Devices," *IEEE Transactions on Electron Devices*, vol. 54, no. 8, pp. 1838-1842, 2007.
- [44] J. D. McDonald and G. C. Albright, "Microthermal imaging in the infrared," *Electronic Cooling*, p. 26, 1997.
- [45] G. C. Albright, J. A. Stump, C. Li, and H. Kaplan, "Emissivity-corrected infrared thermal pulse measurement on microscopic semiconductor targets," in *Thermosense XXIII*, 2001, vol. 4360, no. 1, pp. 103-111.

- [46] K. Decker, S. Ko, and D. Rosato, "Thermal characterization of gallium arsenide FETs," *J. High Density Interconnect*, vol. 3, p. 26, 2000.
- [47] G. C. Albright, J. A. Stump, J. D. McDonald, and H. Kaplan, "True temperature measurements on microscopic semiconductor targets," in *SPIE Thermosense XXI*, 1999, p. 245.
- [48] D. D. Griffin, "Infrared techniques for measuring temperature and related phenomena of microcircuits," *Applied Optics*, vol. 3, p. 1749, 1968.
- [49] A. Hefner, D. Berning, D. Blackburn, C. Chapuy, and S. Bouche, "A high-speed thermal imaging system for semiconductor device analysis," in *Seventeenth Annual IEEE Symposium on Semiconductor Thermal Measurement and Management.*, 2001, pp. 43-49.
- [50] M. Kuball et al., "Measurement of temperature distribution in multifinger AlGaIn/GaN heterostructure field-effect transistors using micro-Raman spectroscopy.," *Applied Physics Letters*, vol. 82, no. 1, p. 124, Jan. 2003.
- [51] M. Kuball et al., "Measurement of temperature in active high-power AlGaIn/GaN HFETs using Raman spectroscopy," *Electron Device Letters, IEEE*, vol. 23, no. 1. pp. 7-9, 2002.
- [52] I. Ahmad et al., "Self-heating study of an AlGaIn/GaN-based heterostructure field-effect transistor using ultraviolet micro-Raman scattering.," *Applied Physics Letters*, vol. 86, no. 17, p. 173503, Apr. 2005.
- [53] G. Abstreiter, "Micro-Raman spectroscopy for characterization of semiconductor devices," *Applied Surface Science*, vol. 50, no. 1-4, pp. 73-78, Jun. 1991.
- [54] M. S. Liu and L. A. Bursill, "Temperature dependence of Raman scattering in single crystal GaN films.," *Applied Physics Letters*, vol. 74, no. 21, p. 3125, May 1999.
- [55] S. Bychikhin, G. Haberfehlner, J. Rhayem, D. Vanderstraeten, R. Gillon, and D. Pogany, "Investigation of smart power DMOS devices under repetitive stress conditions using transient thermal mapping and numerical simulation," *Microelectronics Reliability*, vol. 50, no. 9-11, pp. 1427-1430, Sep. 2010.
- [56] M. Heer et al., "Automated setup for thermal imaging and electrical degradation study of power DMOS devices," *Microelectronics Reliability*, vol. 45, no. 9-11, pp. 1688-1693, Sep. 2005.

- [57] G. Habermann et al., "Thermal imaging of smart power DMOS transistors in the thermally unstable regime using a compact transient interferometric mapping system," *Microelectronics Reliability*, vol. 49, no. 9–11, pp. 1346–1351, Sep. 2009.
- [58] M. Blaho, D. Pogany, E. Gornik, M. Denison, G. Groos, and M. Stecher, "Study of internal behavior in a vertical DMOS transistor under short high current stress by an interferometric mapping method," *Microelectronics Reliability*, vol. 43, no. 4, pp. 545–548, Apr. 2003.
- [59] J. Long, S. O. Memik, G. Memik, and R. Mukherjee, "Thermal monitoring mechanisms for chip multiprocessors," *ACM Transactions on Architecture and Code Optimization*, vol. 5, no. 2, pp. 1–33, Aug. 2008.
- [60] P. Ituero, M. López-vallejo, M. Ángel, S. Marcos, and C. G. Osuna, "Light-Weight On-Chip Monitoring Network for Dynamic Adaptation and Calibration," *IEEE Sensors Journal*, vol. 12, no. 6, pp. 1736–1745, 2012.
- [61] K.-J. Lee, K. Skadron, and W. Huang, "Analytical model for sensor placement on microprocessors," in *Proceedings of the 2005 IEEE International Conference on Computer Design*, 2005, pp. 24–27.
- [62] S. H. Gunther, F. Binns, D. M. Carmean, and J. C. Hall, "Managing the Impact of Increasing Microprocessor Power Consumption," *Intel Technology Journal*, no. Q1, pp. 1–9, 2001.
- [63] F. Liu, "A General Framework for Spatial Correlation Modeling in VLSI Design," in *IEEE/ACM Design Automation Conference (DAC)*, 2007, pp. 817–822.
- [64] S. Sharifi and T. S. Rosing, "Accurate Direct and Indirect On-Chip Temperature Sensing for Efficient Dynamic," *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, vol. 29, no. 10, pp. 1586–1599, 2010.
- [65] D. Brooks and M. Martonosi, "Dynamic Thermal Management for High-Performance Microprocessors," in *Proceedings of the 7th International Symposium on High-Performance Computer Architecture*, 2001, no. C.
- [66] K. Skadron, M. R. Stan, K. Sankaranarayanan, W. E. I. Huang, S. Velusamy, and D. Tarjan, "Temperature-Aware Microarchitecture: Modeling and Implementation," *Architecture*, vol. 1, no. 1, pp. 94–125, 2008.
- [67] K. Etesam-Yazdani, "Continuum and Subcontinuum Thermal Modeling of Electronic Devices and Systems," Carnegie Mellon University, 2006.

- [68] K. Etessam-Yazdani and H. Hamann, "Fast and Accurate Simulation of Heat Transfer in Microarchitectures Using Frequency Domain Techniques," *IPACK*, pp. 1-5, 2007.
- [69] K. Etessam-Yazdani, H. F. Hamann, and M. Asheghi, "Spatial Frequency Domain Analysis of Heat Transfer in Microelectronic Chips with Applications to Temperature Aware Computing," in *Proceedings of IPACK*, 2007.
- [70] P. Rosinger, B. M. Al-Hashimi, and K. Chakrabarty, "Thermal-Safe Test Scheduling for Core-Based System-on-Chip Integrated Circuits," *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, vol. 25, no. 11, pp. 2502-2512, 2006.
- [71] J. T. Hsu and L. Vu-Quoc, "A rational formulation of thermal circuit models for electrothermal simulation. I. Finite element method [power electronic systems]," *IEEE Transactions on Circuits and Systems I: Fundamental Theory and Applications*, vol. 43, no. 9, pp. 721-732, 1996.
- [72] M. R. Stan, K. Skadron, M. Barcella, W. Huang, K. Sankaranarayanan, and S. Velusamy, "HotSpot: a dynamic compact thermal model at the processor-architecture level," *Microelectronics Journal*, vol. 34, no. 12, pp. 1153-1165, Dec. 2003.
- [73] J. T. Hsu and L. Vu-Quoc, "A rational formulation of thermal circuit models for electrothermal simulation. II. Model reduction techniques [power electronic systems]," *IEEE Transactions on Circuits and Systems I: Fundamental Theory and Applications*, vol. 43, no. 9, pp. 733-744, 1996.
- [74] J. M. Wang and T. V. Nguyen, "Extended Krylov subspace method for reduced order analysis of linear circuits with multiple sources," in *Proceedings of the 37th Annual Design Automation Conference (DAC)*, 2000, pp. 247-252.
- [75] J. W. Sofia, "Analysis of thermal transient data with synthesized dynamic models for semiconductor devices," *IEEE Transactions on Components, Packaging, and Manufacturing Technology, Part A*, vol. 18, no. 1, pp. 39-47, 1995.
- [76] F. Christiaens, E. Beyne, and P. Division, "Transient Thermal Modeling and Characterization of a Hybrid Component," in *Proceedings of the 46th Electronic Components and Technology Conference*, 1996, pp. 154-164.
- [77] F. Christiaens et al., "Compact transient thermal models for the polymer stud grid array (PSGATM) package," in *Eurotherm Seminar No. 58*, 1997.

- [78] V. Székely and T. Van Bien, "Fine structure of heat flow path in semiconductor devices: a measurement and identification method," *Solid-State Electronics*, vol. 31, no. 9, pp. 1363-1368, 1988.
- [79] Y. Yang, R. Master, G. Refai-Ahmed, and M. Touzelbaev, "Transient Frequency-Domain Thermal Measurements With Applications to Electronic Packaging," *IEEE Transactions on Components, Packaging and Manufacturing Technology*, vol. 2, no. 3, pp. 448-456, 2012.
- [80] V. Székely, "Identification of RC Networks by Deconvolution: Chances and Limits," *IEEE Transactions on Circuits and Systems-I: Fundamental Theory and Applications*, vol. 45, no. 3, pp. 244-258, 1998.
- [81] Y. C. Gerstenmaier and G. Wachutka, "Calculation of the temperature development in electronic systems by convolution integrals," in *Proceedings of the Sixteenth Annual IEEE Semiconductor Thermal Measurement and Management Symposium*, 2000, pp. 50-59.
- [82] D. Schweitzer, "A Fast Algorithm for Thermal Transient Multisource Simulation Using Interpolated Zth Functions," *Components and Packaging Technologies, IEEE Transactions on*, vol. 32, no. 2, pp. 478-483, 2009.
- [83] V. Székely, "On the Representation of Infinite-Length Distributed RC One-Ports," *IEEE Transactions on Circuits and Systems*, vol. 38, no. 7, pp. 711-719, 1991.
- [84] H. S. Carslaw and J. C. Jaeger, *Conduction of Heat in Solids*, 2nd ed. Oxford, U.K.: Oxford University Press, 1959.