

**THERMAL AND POWER DELIVERY NETWORK MODELING FOR
EMERGING MICROELECTRONIC INTEGRATION PLATFORMS**

A Dissertation
Presented to
The Academic Faculty

By

Yang Zhang

In Partial Fulfillment
of the Requirements for the Degree
Doctor of Philosophy in the
School of Electrical and Computer Engineering

Georgia Institute of Technology

December 2017

COPYRIGHT 2017 © YANG ZHANG

**THERMAL AND POWER DELIVERY NETWORK MODELING FOR
EMERGING MICROELECTRONIC INTEGRATION PLATFORMS**

Approved by:

Dr. Muhannad S. Bakir, Advisor
School of Electrical and Computer
Engineering
Georgia Institute of Technology

Dr. Azad J. Naeemi
School of Electrical and Computer
Engineering
Georgia Institute of Technology

Dr. Saibal Mukhopadhyay
School of Electrical and Computer
Engineering
Georgia Institute of Technology

Dr. Arijit Raychowdhury
School of Electrical and Computer
Engineering
Georgia Institute of Technology

Dr. Yogendra Joshi
George W. Woodruff School of
Mechanical Engineering
Georgia Institute of Technology

Date Approved: August 17, 2017

The world is a fine place and worth fighting for.

Ernest M. Hemingway

To my family
for their constant love, caring,
support, and understanding.
I love you all.

ACKNOWLEDGEMENTS

First I wish to thank my advisor, Dr. Muhannad S. Bakir. He has always been considerate, caring, and eager to help since I started working in the Integrated 3-D System (I3DS) group in June, 2013. None of my academic and personal milestones would have been possible without his support, advice, inspiration and encouragements. I cannot be grateful enough for what he has done for me. The research freedom and the strong sense of collaboration he planted, the insight and long-term view on research topics he gave, the willingness he showed to spend tremendous time improving every piece of the work, the patience he paid to revise every part of all the research documents, the sincere compliments he delivered, and the passion he conveyed not only made me a better researcher but also a collaborator and mentor. Moreover, his instant “Yes” and “Go for it” responses to all my requests and the generosity for paid vacations made my life and Ph.D study nothing but easier. I also thank him for referring me without any reservation to scholarships, internships and full-time jobs. Because of him, I can always smile and be happy throughout this rough road to finish my thesis.

I wish to thank Dr. Azad J. Naeemi and Dr. Saibal Mukhopadhyay for being in my thesis reading committee and their numerous time and insightful feedback on my work. Dr. Naeemi helped me understand the fundamental limits for integrated circuits from a wide scope and extend my research views to graphene and spin-based devices. Dr. Mukhopadhyay provided me a lot of significant feedback on CAD and circuit design field. I would also like to acknowledge Dr. Arijit Raychowdhury and Dr. Yogendra Joshi for serving on my defense committee. Dr. Raychowdhury helped me better shape the signal channel modeling work and Dr. Joshi’s comments on the thermal part are very constructive. I also want to express my gratitude to Dr. Sung-kyu Lim for his help in my admission to the graduate school and provided two semesters of training in VLSI and CAD tools. In short, it is fortune to study in Georgia Tech and to be advised by so many top faculty members.

I wish to thank my mentors in the I3DS group. During my early years, I was so lucky to have Dr. Yue Zhang as my mentor, who is warmhearted and unselfish. The collaboration with her in developing thermal isolation technologies for 3-D ICs led to fruitful and significant results. Moreover, the personal help she gave and the encouraging words she said are gratefully acknowledged. I trust great people like her will always be blessed. I wish to thank Dr. Li Zheng for taking me into the field of power delivery and microfluidic cooling. I would also like to thank Dr. Xuchen Zhang deeply for various inspirational discussions, sharp comments at key areas in my research, and his education of signal integrity and 2.5-D integration. I cannot thank him enough for taking care of me in many aspects. I would also thank Dr. Chaoqi Zhang, Dr. Paragkumar Thadesar and Dr. James Yang for their kindness, advice and mentorship.

I wish to extend my thanks to other members in I3DS group: Dr. Hanju Oh for the discussions on TSVs and help in HFSS simulations, Thomas E. Sarvey for the collaboration on microfluidic cooling-related projects, William Wahby for the education of system-level modeling for 3D-ICs, Md Obaidul Hossen for so many days and nights working together on power delivery, your modesty and courtesy to me when we have technical debates and the enthusiasm to cheer me up when we met difficulties, Paul Jo and Joe L. Gonzalez for the discussions on MFIs and geometry optimization, Muneeb Zia, Reza Abbaspour, Congshan Wan, Sreejith K. Rajan and Haopeng Zhang.

A big thank you to my industry colleagues, mentors and managers during my summer internships: Dr. Woosung Choi, Dr. Mindy Lee and Dr. Nuo Xu at *Samsung* and Jimmy Cheng, Ted Hong, Dr. Runjie Zhang, Dr. Fuqiang Zhang, Dr. Xin Huang, Derong Liu and Fangzhou Wang at *Oracle*. I would also thank Hesam F. Moghadam and Michael Dayringer from *Oracle Labs* for their support and guidance on power delivery. I also would like to thank my undergraduate advisor, Prof. Xiaoyan Liu from Peking University who opened the door of research to me.

I am grateful for all my peers in Georgia Tech and the beautiful Atlanta area, espe-

cially Dr. Jianyong Xie for teaching me thermal modeling, Dr. Xin Zhao for power delivery modeling and optimization, Sensen Li for teaching me RF IC and optics, Hang Zhang for convex optimization and machine learning algorithms, Dr. Chenyun Pan for the consulting on architectural modeling tools, Dr. Sandeep K. Samal for monolithic 3-D IC knowledge and device modeling, Dr. Yarui Peng for PDN and signal integrity discussions, and Ke Liu for numerical methods and quadratic programming. I would also like to thank Tao Wang, Kexin Yang, Rui Zhang, Yun Long, Tao Zhu, Dr. Jilai Ding, Dr. Xiaotang Du, Xiaojia Zhang, Heng Chi and Yang Wan. I also appreciate the help and services provided by the staff in ECE department, MIRC and MARCUS building. I cannot go through my Ph.D. journey without your support and unlimited help.

I wish to acknowledge some of my old friends since college who always shared with me their research experiences, interesting stories and thoughts on new ideas and technologies, and all the encouraging words: Dr. Ang Li, Dr. Shaodi Wang, Dr. Min Ye, Dr. Chicheng Zhang, Dr. Lisong Li, Dr. Yuhan Yao, Xiaoji Li, Junkai Jiang, Tao Xiong, Wenyi Liu, Yifei Wang, Mingrui Sun and Yefan Liu.

Finally and most importantly, I want to thank my family: my parents and my girlfriend along with her parents. My parents, Changping Zhang and Yuanen Yang, the greatest in the world, not only grew me up with endless love but endowed me with their determination, perseverance, patience, devotion and beautiful minds. I cannot expect more from them. My girlfriend, Tianyi Lu who is the most special one in the world, sacrificed so much for me in the past five years. She has been the best friend, listener, supporter, and companion. I cannot put it in words how thankful and appreciative I am for what she does. Falling in love with her is the most beautiful thing in my life and I just want her to know I love her very much and trust there will be more good things happening in the future. I also thank her parents, Yong Lu and Yili Zhang, without their understanding, support and blessing, the long-distance relationship with their USC Ph.D daughter will never be possible.

TABLE OF CONTENTS

Acknowledgments	v
List of Tables	xiii
List of Figures	xv
Summary	xxi
Chapter1: Introduction	1
1.1 Motivation	2
1.1.1 Thermal challenges for 3-D and 2.5-D integration	2
1.1.2 Power delivery challenges for 2.5-D integration	3
1.1.3 Interaction between thermal, power delivery and power dissipation	4
1.2 Research Objective and Contribution	5
1.3 Organization of the Thesis	7
Chapter2: Thermal isolation technologies for heterogeneous 3-D ICs	9
2.1 Thermal modeling framework	9
2.1.1 Formulation using finite volume methods	10
2.1.2 Modeling of air- and microfluidic-cooled heat sinks	12
2.1.3 Comparison with existing thermal models	13

2.1.4	Thermal simulation flow and validation	15
2.2	Motivation of thermal isolation	21
2.3	Proposed architecture with thermal isolation technologies	22
2.4	Thermal evaluation of the proposed architecture	24
2.4.1	Thermal specification of 3-D stacks	24
2.4.2	Comparison of different 3-D stacks	26
2.5	Design space exploration of the proposed architecture	27
2.5.1	Impact of Microbump	28
2.5.2	Impact of TSVs	29
2.5.3	Impact of die thickness	31
2.6	Experimental demonstration ¹	33
2.7	Conclusion	38
 Chapter3: Thermal evaluation of 2.5-D integration using bridge-chip technology		39
3.1	Introduction	39
3.2	2.5-D and 3-D benchmark architectures	41
3.2.1	2.5-D integration platforms	41
3.2.2	3-D integration platforms	42
3.3	Thermal specifications for 2.5-D and 3-D integration evaluation	43
3.3.1	Layer thickness and material property	44
3.3.2	Geometry parameter and boundary conditions	44
3.3.3	Power maps of integrated dice	46
3.3.4	Microbumps and TSVs dimensions	47

3.4	Comparison of different 2.5-D integrations	47
3.4.1	Impact of the thickness of interposer and bridge chip	49
3.5	Thermal comparison between 2.5-D and 3-D integration	50
3.6	Thermal study of bridge-chip 2.5-D integration	52
3.6.1	TIM properties and die thickness mismatch	52
3.6.2	Impact of die thickness on heat spreading	55
3.6.3	Impact of microbump and underfill on secondary heat path	56
3.6.4	Impact of die spacing on thermal coupling	57
3.6.5	Transient thermal coupling	58
3.7	Conclusion	60
Chapter4:	Power delivery network evaluation and benchmarking for 2.5-D in-	
	tegration using bridge-chip technology	61
4.1	PDN modeling framework	61
4.1.1	Board-, Package- and on-die PDN models	62
4.1.2	PDN model formulation and simulation flow	66
4.1.3	Comparison with existing PDN models	67
4.1.4	Steady-state and transient analysis validation	67
4.2	PDN challenges of 2.5-D integration	71
4.3	PDN evaluation and benchmarking of 2.5-D integration	73
4.3.1	Study cases	73
4.3.2	PDN design parameters and specification	74
4.3.3	Comparison of different 2.5-D integration	75
4.4	Design space exploration of 2.5-D integration	81

4.4.1	Impact of total current requirement	81
4.4.2	Impact of on-die metal layers	82
4.4.3	Impact of TSV and overlap area	83
4.4.4	Inserting vias in bridge-chip	85
4.5	Conclusion	85
Chapter5: Integrated Thermal and Power Delivery Network Co-Simulation Framework		86
5.1	Introduction of thermal and PDN co-simulation	86
5.2	Simulation flow of integrated thermal and PDN modeling framework	88
5.2.1	Steady-state analysis	88
5.2.2	Transient-state analysis	89
5.2.3	Framework algorithm	91
5.3	Modeling methodologies and implementation	91
5.3.1	Thermal model and formulation	91
5.3.2	PDN model and formulation	93
5.3.3	Power update models	95
5.4	Comparison of models with different number of dependencies included	98
5.4.1	Thermal and PDN specification	98
5.4.2	Modeling Scenarios with different number of dependencies	101
5.4.3	Comparison results	102
5.4.4	Accuracy Improvement Compared to Prior Work	106
5.5	Conclusions	107

Chapter6: Digital signal channel modeling for 2.5-D and 3-D integration	108
6.1 Circuit models of digital signal channels in 2.5-D and 3-D integration	108
6.2 Comparison of integration platform latency, energy efficiency and bandwidth density	112
6.3 Impact of technology parameter scaling	115
6.3.1 Technology process scaling	116
6.3.2 Impact of interconnect wire length in signal channels	118
6.4 Impact of temperature on signaling in 2.5-D and 3-D integration	120
6.4.1 2.5-D and 3-D signaling comparison revisit with the impact of temperature	121
6.4.2 Thermal and electrical tradeoffs on die spacing in 2.5-D integration	123
6.5 Summary	124
Chapter7: Conclusion and future directions	126
7.1 Summary of the presented work	127
7.2 Summary of the future directions	129
References	132
Vita	144

LIST OF TABLES

1	Comparison of different thermal modeling	14
2	Efficiency and accuracy comparison	18
3	The specification of simulated stack	24
4	The comparison of different architectures	26
5	Thermal specification	44
6	Thermal comparison of bridge-chip 2.5-D and 3-D integration	51
7	Comparison of different PDN modeling work	68
8	Validation Results	69
9	Parameters for PDN model	75
10	Transient state analysis results	80
11	Parameters for thermal model	99
12	Parameters for PDN model	100
13	Simulation Model	101
14	Results for different detailed models	103
15	Results after Different Number of Iterations	106
16	Physical dimensions of each parameter of signaling models	109

17	Equation for parasitics estimation	111
18	Comparison of different integration platforms	113
19	Bandwidth density of each integration platform	115
20	Channel length, minimum inverter size, and supply voltage of each process technology using <i>PTM</i> device library	116

LIST OF FIGURES

1	2.5-D chip stack using (a) bridge-chip technology (b) interposer technology. (c) non-embedded bridge-chip using multi-height microbump technology.	1
2	3-D chip stack (a) TSV-based (b) monolithic nanoscale via based.	2
3	PDN structures of three different 2.5-D integration platforms (a) interposer (b) bridge-chip	4
4	The interactions between temperature, PDN noise, and power for 3D-IC. . .	4
5	Illustration of finite volume scheme	10
6	Illustration of non-conformal meshing.	11
7	Illustration of boundary nodes	12
8	Thermal simulation flow.	15
9	Steady-state validation setup.	17
10	Steady-state validation results.	17
11	Impact of mesh size: tradeoff between efficiency and accuracy.	18
12	Transient thermal validation experiments (a) a 2-die 3-D stack (b) transient thermal validation results	19
13	(a) The cross-sectional view of the experimental setup (b) the full chip view of the testbed.	20
14	Comparison between the model and ANSYS simulation results in the power excitation layer along $y = 0.5$ cm	20

15	DRAM retention time reduces exponentially as the temperature goes up . . .	21
16	Proposed architecture with interposer-embedded heat sink, thermal bridge and air gap isolation	23
17	3-D stack (a) with conventional air cooled heat sink (b) with interposer embedded microfluidic heat sink (MFHS)	23
18	(a) Physical structure of the extended heat spreader (b) Lumped resistance modeling for fins of extended heat spreader and TIM	25
19	Power density distribution: (a) Memory die (b) Processor die	26
20	Thermal maps of proposed stack when processor is in active mode (a) without TSVs (b) with TSVs	28
21	The impact of the microbumps. (a) change the number of microbumps (b) change the diameter of microbumps.	29
22	The impact of the TSVs. (a) change the number of TSVs. (b) change the diameter of TSVs.	30
23	Thermal maps of clustered TSVs and uniform TSVs (a) TSVs are clustered in the solid-line rectangle (b) The same amount of TSV are uniformly distributed	31
24	The power map of the processor die. The hotspot (blue square) has power density of $135 W/cm^2$ and the background (grey area) power density is $35 W/cm^2$	32
25	The impact of die thickness. (a) change the thickness of processor. (b) change the thickness of memory.	32
26	(a) Schematic of the designed testbed for evaluation of the proposed thermal isolation technologies. (b) Top tier (low-power) and (c) bottom tier (high-power) layout design	34
27	(a) Microfluidic test setup to evaluate the thermal isolation technologies. (b) Top and (c) bottom view of the stack assembled to a PCB board using wire bonding	36
28	(a) Uniform power density of $10 W/cm^2$ in the bottom tier (Case A), (b) background power of $10 W/cm^2$ plus two hotspots each dissipating $150 W/cm^2$ (Case B), and (c) Junction temperature fluctuation of top and bottom tiers in Case A and Case B	37

29	2.5-D chip stack using (a) bridge-chip technology (b) interposer technology. (c) heterogeneous interconnect stitching technology (HIST).	40
30	Illustration of the envision FPGA-CPU-Memory 2.5-D chip stack using bridge-chip technology (top view)	42
31	3-D chip stack (a) TSV-based (b) monolithic nanoscale via based.	42
32	2.5-D integration using bridge-chip technology with detailed layer information.	43
33	Setup of characterizing effective convection coefficient.	45
34	Power density maps of each die (a) FPGA die, 44.8 W (b) Processor die 74.49 W(c) DRAM die (cell circuit), 5.65 W for cell circuit.	46
35	Top view of thermal profiles of each die in all cases. The bottom die, also the hottest die of DRAM chip stack is plotted. (a) embedded bridge-chip, $T_{max} : 104.92 \text{ }^\circ\text{C}$ (b) interposer, $T_{max} : 102.80 \text{ }^\circ\text{C}$ (c) HIST, $T_{max} : 104.23 \text{ }^\circ\text{C}$	48
36	Illustration of heat spreading effects of (a) the package layer in embedded bridge-chip based 2.5-D, $60.22 \text{ }^\circ\text{C} \sim 104.60 \text{ }^\circ\text{C}$ (b) the interposer layer in interposer based 2.5-D, $61.61 \text{ }^\circ\text{C} \sim 101.14 \text{ }^\circ\text{C}$	49
37	The impact of interposer and bridge thickness.	50
38	Thermal profile of each die in 3-D stack cases (a) Monolithic-3D (b) TSV-3D	51
39	The impact of thermal conductivity of TIM.	53
40	Illustration of die thickness mismatch.	53
41	Impact of die thickness mismatch of (a) processor (b) FPGA (c) DRAM. The solid line in the figures is the case using default TIM filler ($3 \text{ W}/^\circ\text{C}\cdot\text{m}$) and the dashed line is the case for using copper filler ($400 \text{ W}/^\circ\text{C}\cdot\text{m}$). . . .	54
42	The impact of die thickness scaling. The dotted line plots the T_{min} of each die.	55
43	The thermal profile of each die (die thickness is $1 \text{ } \mu\text{m}$). The heat spreading is confined, and the block outlines are clearly observed from the thermal maps.	56

44	The impact of effective thermal conductivity of microbump layer for bridge-chip case.	57
45	Impact of die spacing. As the die spacing increases, the junction temperature decreases.	58
46	(a) Emulated processor power (b) transient analysis results of bridge-chip 2.5-D integration	59
47	The PDN structure hierarchy. From left to right, a lumped model of board-level PDN, a distributed model of package-level PDN and on-chip PDN are shown, respectively.	62
48	The two-layer package PDN model of power/ground planes	63
49	The on-die PDN model. Only one current source and one C4 bump is shown.	64
50	Re-organization of a non-uniform PDN layout	64
51	Map fine-grained power PDN layout to coarse meshing grids (a) vias (b) wires.	65
52	The noise profile of IBM3 (a) Provided results by <i>IBM</i> PG benchmarks (b) modeling results.	70
53	Bump current comparison of IBM3.	70
54	The transient noise the node in IBM2 with maximum error.	71
55	Three different 2.5-D integration platforms (a) interposer (b) bridge-chip (c) HIST	72
56	The current density of each die. (a) die #1 (b) die #2	74
57	Illustration of bridge chip placement: an example with a single bridge chip.	74
58	The IR drop profile of each die for (a) Single die (b) interposer (c) single bridge-chip with a overlap area of $0.5 \times 6 \text{ mm}$	76
59	Illustration of bridge chip placement: an example of 5 bridge chips.	77
60	IR-drop analysis results using 5 bridge chips.	77
61	The IR drop profile of each die for the case with 5 bridge chips.	78

62	(a) Impedance analysis of one on-die PDN node and illustration of the switching current activity (b) waveform #1 1 GHz frequency (c) waveform #2, 4 GHz frequency	79
63	Transient analysis results of waveform #1 (a) Die #1 (b) Die #2 and waveform #2 (c) Die #1 (d) Die #2	80
64	The impact of total current.	82
65	The impact of adding metal layers	82
66	The impact of TSV and overlap area	84
67	The impact of inserting vias in the bridge-chip	84
68	The interactions between temperature distribution, PDN noise and power dissipation	87
69	Steady-state simulation flow.	89
70	Transient-state simulation flow.	90
71	Validation of stable temperature assumption in microsecond scale.	93
72	Block diagram of the PDN structure. The distributed power/ground rail is abstracted for visualization.	94
73	The 2-D piecewise linear model (a) response surface with 8 sampling points (b) the maximum error of models with different number of sampling points.	97
74	The 3D-IC example: processor on memory sack.	98
75	Reference power maps (a) Memory die (2.82W) (b) Processor die (74.49W).	99
76	On-die PDN structure (a) Power/ground pads with wires (b) Interleaved structure of power/ground wires. (c) Dense PDN wires between power/ground pads	100
77	Steady state analysis result of full-model case (a) IR-drop of memory (4.71 mV) (b) Thermal of memory (84.35 °C) (c) IR-drop of processor (33.05 mV) (d) Thermal of processor (85.51 °C).	105

78	Transient power supply noise comparison (a) memory die (b) processor die. PDN-therm is very similar to PDN-power, thus omitted for better visualization.	105
79	Illustration of digital signal channels (a) bridge-chip 2.5-D integration (b) interposer and HIST 2.5-D integration (c) 3-D integration	109
80	Illustration of digital signal channels (a) 2.5-D integration (b) 3-D integration	110
81	Impact of pad scaling on electrical performance of signal channel.	114
82	Delay and energy of the signal channels implemented by different technology nodes.	117
83	Relative difference between bridge-chip/interposer and HIST vs. device process technology (a) delay (b) energy.	118
84	The impact of interconnect scaling on (a) delay (b) energy	119
85	Relative difference between bridge-chip/interposer and HIST vs. wire lengths (a) delay (b) energy	120
86	Impact of temperature on delay and energy of HIST-based 2.5-D integration (a) 45 nm (b) 14 nm	121
87	Comparison of HIST-based 2.5-D and TSV-based 3-D integration (a) delay using 45 nm library (b) energy using 45 nm library (c) delay using 14 nm library (d) energy using 14 nm library	122
88	Thermal and electrical tradeoffs for (a) delay using 45 nm library (b) energy using 14 nm library	124
89	A 3-D chip stack using FOWLP technology.	129
90	Two architectures for nanophotonics-based systems with thermal isolation technologies (a) using air-cooled heat sink (b) using microfluidic-cooled heat sink.	130
91	Demonstrated HIST platforms with active chips (a) two active chips emulating driver and receiver circuitries (b) emulated HIST system	131

SUMMARY

In order to keep pace with rapidly increasing system interconnection requirements, multiple advanced interconnect technologies have been proposed, including 2.5-D and 3-D integration technologies. Despite the benefits of such systems in communication bandwidth, power efficiency, footprint reduction and etc, there are thermal and power delivery challenges, which are potential show stoppers. To enable the design space exploration of these systems from the perspective of temperature and power supply noise, a thermal and a PDN modeling framework based on finite difference methods are developed, implemented and validated, respectively.

To address the thermal coupling issues in heterogeneous 3-D ICs, a stacking structure is proposed using interposer embedded microfluidic cooling, air gap isolation and an extended heat spreader. The proposed architecture is compared to conventional air cooled stack, a significant temperature reduction of the low-power temperature sensitive die is achieved.

Moreover, we explore the design considerations of three approaches of 2.5-D integration from a thermal perspective and compare to 3-D ICs. The impact of several different technology parameters is studied, such as die thickness mismatch and die spacing. Design guidelines are presented for integrating dice with different die thickness and thermal and electrical tradeoffs are discussed for selecting die spacing in such heterogeneous integration systems.

Next, power supply challenges are investigated for interposer and bridge-chip based 2.5-D integration platforms. Interposer based 2.5-D integration may exhibit a worse power supply noise due to TSV parasitics. In bridge-chip based 2.5-D integration, due to the fact that the bridge chips underneath the active dice block access to package power/ground planes, there are higher supply voltage noise in these regions.

Temperature, supply voltage and power dissipation are dependent of each other. The temperature impacts the leakage power and the power grid resistivity. Power dissipation

determines the source current of the chip and is also the excitation of the PDN noise. Reversely, the power supply voltage impacts both leakage and dynamic power. Without considering the interactions between each of the components, the results of the standalone models are inaccurate. On the contrary, the design under worst constraints of temperature and supply noise will be over-optimized. To accurately model temperature, power supply noise and power dissipation, an integrated thermal and PDN modeling framework considering the above dependencies is proposed, developed, implemented and studied.

Lastly, a signaling evaluation framework to benchmark the communication channels of 2.5-D and 3-D integration platforms is presented. The impact of technology scaling, pad size, and interconnect length are shown. HIST platforms show significant latency and power efficiency improvement compared to bridge-chip and interposer based platforms. Moreover, thermal impact on signaling is discussed for 2.5-D and 3-D integration.

CHAPTER 1

INTRODUCTION

Emerging applications such as the Internet of Things (IoT), cloud computing, and machine learning based artificial intelligence have presented performance, communication bandwidth, and functionality challenges to conventional integrated circuits (ICs) and electronic systems [1]. To address these challenges, novel heterogeneous computing fabrics based on processor (CPU), FPGA/GPU/ASIC accelerators, and high density memory [2] have been widely proposed and studied [3, 4] to increase system throughput, computation capability, and efficiency. However, one of the biggest bottlenecks for such systems is the inter-die bandwidth which can cause functional blocks to be idle during data transfer [5], leading to lower system performance.

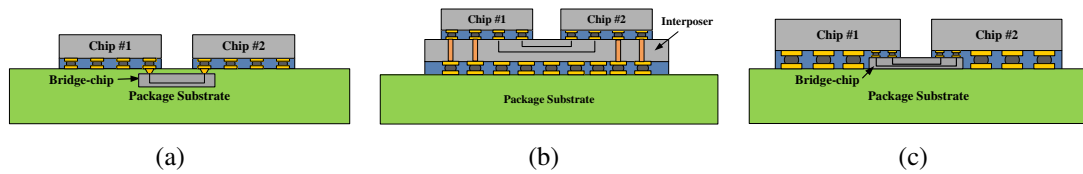


Figure 1: 2.5-D chip stack using (a) bridge-chip technology (b) interposer technology. (c) non-embedded bridge-chip using multi-height microbump technology.

In order to keep up with rapidly evolving off-chip communication requirements, multiple integration platforms using advanced interconnect technologies have been explored and demonstrated. Silicon interposer-based 2.5-D integration of FPGA dice and data converters achieves an aggregate bandwidth in excess of 400 Gb/s [6], as shown in Fig. 1(b). Three-dimensional (3-D) processor-on-memory integration using through silicon vias (TSVs) exhibits a maximum memory bandwidth of 510.4 Gb/s at 277 MHz [7], as shown in Fig. 2(a). Monolithic 3-D integration is another promising option, which achieves even higher bandwidth than TSV-based 3-D integration resulting from the utilization of shorter and

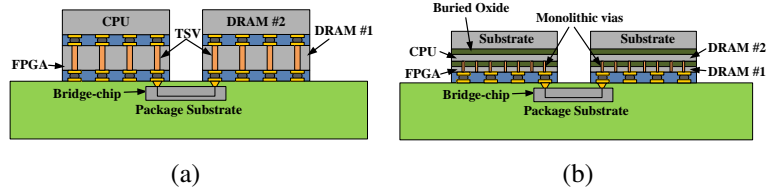


Figure 2: 3-D chip stack (a) TSV-based (b) monolithic nanoscale via based.

denser nanoscale vertical vias [8], as shown in Fig. 2(b). Moreover, there has been recent interest in multi-die packages using bridge-chip technology including embedded multi-interconnect bridge (EMIB) technology [1, 9] and heterogeneous interconnection stitching technology (HIST) [10] to enable 2.5-D microsystems, as shown in Fig. 1(a) and Fig. 1(c), respectively. In its simplest form, bridge-chip technology utilizes a silicon die with high density interconnects for inter-die communication. The performance metrics of these 2.5-D integration technologies are comparable to interposer-based 2.5-D and 3-D solutions but many other benefits are offered, including the elimination of TSVs.

However, as multi-die packaging continues the trend of placing more high-performance (i.e. CPU, GPU, FPGA) chips in a package, the total power density is expected to increase beyond 100 W/cm^2 [11]; the impedance of the power delivery network will be larger, and air cooled heat sinks will become incapable of cooling the whole chip without keeping much of the silicon dark. Therefore, in spite of the bandwidth and power efficiency benefits brought by 2.5-D and 3-D ICs, there are thermal and power delivery challenges that are potential show stoppers.

1.1 Motivation

1.1.1 Thermal challenges for 3-D and 2.5-D integration

The thermal challenge is comprised of two distinct components with each requiring separate optimization and technology solutions: first, stacking dice in 3-D or 2.5-D increases the total power density and the thermal resistance of dice to the atop attached heat sink;

secondly, stacked dice will experience unwanted thermal crosstalk, particularly between high-power die and low-power temperature sensitive components. While significant effort in the literature has addressed the need for improved cooling [12, 13, 14, 15, 16], less effort has gone towards exploring the negative effects of thermal coupling between dice and any solutions to reduce inter-die thermal coupling.

Prior efforts have extensively conducted thermal simulation and analysis for TSV-based 3-D IC [17], monolithic-based 3-D IC [18, 19], and interposer- or glass-based 2.5-D [20, 21, 20] integration system. However, there are no thermal modeling efforts focusing on 2.5-D bridge-chip-based interconnection platforms (as shown in Fig. 1(a)) nor is there any analysis on the impact of technology parameters such as die thickness mismatch and die spacing. Moreover, previous thermal efforts have generally focused on one of the above technologies; there is a need for thermal benchmarking of all these approaches.

1.1.2 Power delivery challenges for 2.5-D integration

The trends of lower supply voltage, higher current demand and increased power density are making power delivery in high-performance digital systems an increasingly difficult challenge [22]. Due to the resonances generated from the interactions of the board-, package-, and die-level parasitics, it is difficult to ensure power integrity over a wide frequency range. Moreover, there is an increasing interest in placing multiple dice into a single package using three-dimensional (3-D) and 2.5-dimensional (2.5-D) integration technologies [9, 10, 17], which exacerbates the power delivery challenges.

Power supply noise (PSN) in traditional single-chip [23, 24, 25] and 3-D ICs [26, 27, 28] have been extensively studied in the literature. However, 2.5-D integrated electronics have not been investigated as thoroughly. Specifically, 2.5-D integrated electronics have several unique attributes that require consideration. For example, for embedded multi-die interconnect bridge (EMIB) technology as shown in Fig. 3(b) [9], signal interconnections and driver circuits are placed, generally, on the edges of the dice and above the bridge-

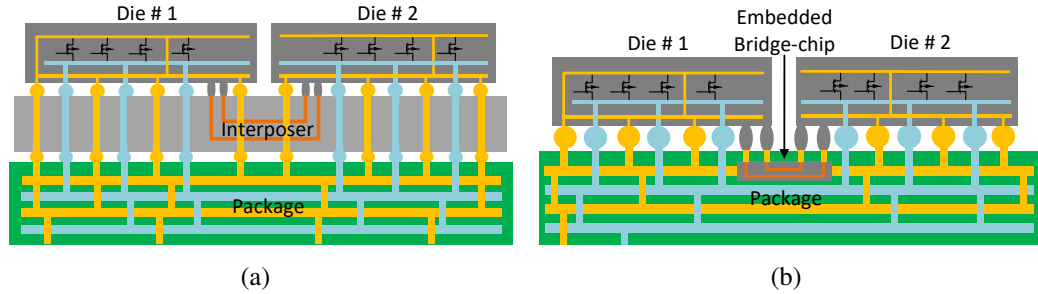


Figure 3: PDN structures of three different 2.5-D integration platforms (a) interposer (b) bridge-chip chips, which may lead to a reduction in the power/ground C4 bumps that have access to the package-level power/ground planes. However, there are no PDN modeling efforts focused on bridge-chip based 2.5-D integration.

1.1.3 Interaction between thermal, power delivery and power dissipation

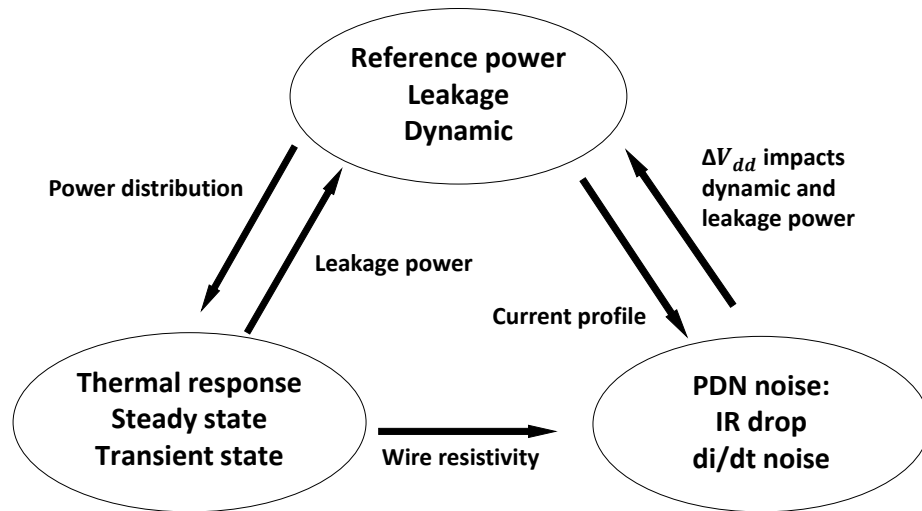


Figure 4: The interactions between temperature, PDN noise, and power for 3D-IC.

Fig. 4 shows the dependencies between power dissipation, temperature, and PDN noise [29]. The temperature impacts the leakage power and the power grid resistivity. Power dissipation determines the source current of the chip, and is also the excitation of the PDN switching noise. Reversely, the power supply voltage impacts both leakage and dynamic power. Without considering the interactions between each component in Fig. 4, the results

of the standalone models are inaccurate. For example, Su et. al pointed out that the leakage power was underestimated by as much as 30% without including the impact of temperature and power supply voltage [30]. Hence, it is essential to build a thermal and PDN noise co-analysis flow to explore the design space of novel integration platforms and answer what-if type questions on the impact of technology parameters.

Prior work focused on either developing the individual thermal [31, 18, 32] and PDN models [24, 33, 27] or studying part of the interactions [34, 28, 30]. There are no frameworks capable of performing steady-state and transient analysis on thermal and PDN noise incorporating the impact of their variation on power dissipation. To meet this need in the literature as well as to close the loop shown in Fig. 4, there is a need to build a framework that is capable of simultaneously analyzing the temperature, PDN noise, power dissipation and the interactions between them for both steady-state and transient analysis. Additionally, there is a need of a comparison of models with different parts of the interactions to understand the accuracy and necessity of each model. With a complete and comprehensive model, the traditional design methodology under worst scenarios will remove some pessimism and reduce the design constraints.

1.2 Research Objective and Contribution

This thesis focuses on understanding the details and challenges discussed above, and developing thermal and power delivery network simulation frameworks to analyze and address them. This work consists of five key topics summarized as follows.

1. **Thermal isolation technologies for heterogeneous 3-D ICs.** First, we build a computationally efficient thermal model to perform a parametric study to answer what-if type questions and enable fast design space exploration. The model has been validated with *ANSYS* for both steady and transient state simulations with a maximum error of less than 7%. Next, we propose a novel stacking structure with microfluidic cooling embedded in the interposer, thermal isolation between the memory and pro-

cessor dice, and a thermal bridge above the isolated die. Last, we use the thermal model to evaluate the proposed architecture and compare with two baseline architectures. The new architecture exhibits thermal benefits over conventional stacks and is of high value in the heterogeneous integration of high-power and low-power dice. In addition, we thermally explore our proposed system as a function of microbump density, TSV density and geometry, die thickness, and other system parameters.

2. **Thermal evaluation and benchmarking of 2.5-D integration using bridge-chip technology.** Thermal benchmarking of a number of 2.5-D integration approaches is performed and compared to 3-D ICs for completeness. Thermal modeling shows that the evaluated 2.5-D integration approaches exhibit similar thermal characteristics, but show significant improvements compared to 3-D IC solutions with the same power consumption. Moreover, bridge-chip-based 2.5-D integrated systems are explored as a function of bridge-chip thickness, thermal interface material properties, microbump properties, die thickness, die thickness mismatch, and die-to-die spacing along with transient analysis to investigate time-domain thermal coupling.
3. **Power delivery network benchmarking and evaluation of heterogeneous 2.5-D integration using bridge-chip technology.** First, a computationally efficient model is developed to answer what-if type questions and flexible enough to modify to explore different configurations of emerging integration architectures. The model has been validated against *IBM* open source power grid benchmarks for both steady and transient state simulations with a maximum error of less than 7.29% and 0.67%, respectively. Second, the power delivery networks of the interposer and bridge-chip based 2.5-D integration platforms are evaluated. The simulation results show that interposer based 2.5-D integration may exhibit a worse power supply noise due to the TSV parasitics. In bridge-chip based 2.5-D integration, under the assumption that the bridge-chips underneath the active dice block access to package power/ground

planes, some power delivery challenges are highlighted. In order to mitigate power supply noise (PSN), it is suggested to minimize the bridge-chip-to-active dice overlap area and to use multiple smaller bridge chips instead of a single large one.

4. **Integrated thermal and power delivery network co-Simulation framework for single-die and multi-die assemblies.** A thermal and power delivery network (PDN) co-simulation framework for single-die and emerging multi-die configurations is presented. The proposed framework incorporates the interactions between temperature, supply voltage, and power dissipation. The temperature dependencies of wire resistivity and leakage power are considered, and the supply voltage dependencies of power dissipation are modeled. Starting with a reference power dissipation, the framework is capable of evaluating the temperature distribution and PDN noise simultaneously and eventually updating the power dissipation based on the thermal and supply voltage conditions.
5. **Digital signal communication channel modeling, benchmarking and comparison of 2.5-D and 3-D integration platforms.** An electrical performance evaluation framework of communication channels of 2.5-D and 3-D integration platforms using compact circuit models is presented. The propagation delay, energy per bit and bandwidth density of each integration platform are benchmarked and compared. By using the modeling framework, the impact of technology scaling and wire length are investigated and discussed. Moreover, thermal impact on 2.5-D and 3-D integration are quantitatively analyzed along with the discussion of the electrical and thermal tradeoff of 2.5-D on die spacing.

1.3 Organization of the Thesis

The rest of the thesis is organized as follows:

1. In chapter 2, a thermal modeling framework is described and validated against *AN-*

SYS. Thermal isolation technologies using air gap, microfluidic cooling and an extended heat spreader are introduced, benchmarked and evaluated.

2. In chapter 3, thermal evaluation and benchmarking for 2.5-D integration using bridge-chip technology is performed. Various technologies are thermally evaluated.
3. In chapter 4, a power delivery network modeling framework is presented and validated against open source *IBM* power grid benchmarks. The framework is used to benchmark and evaluate the power delivery networks for heterogeneous 2.5-D integration using bridge-chip technology.
4. In chapter 5, an integrated thermal and power delivery network co-simulation framework is presented and 3-D memory-on-processor stack is used to qualify the accuracy of different models.
5. In chapter 6, signaling simulation, benchmarking and comparison of different 2.5-D and 3-D interconnected platforms is presented.
6. In chapter 7, the conclusions of this thesis are summarized, along with the discussion on the potential future works.

CHAPTER 2

THERMAL ISOLATION TECHNOLOGIES FOR HETEROGENEOUS 3-D ICs

Significant work has addressed the thermal challenges of 3D integration, e.g. the increasing power density and inter-stack thermal resistance [14, 19, 18, 12], while relatively fewer efforts have been proposed to mitigate the negative effects of thermal coupling between different dice in a 3D heterogeneous stack, and in particular, to minimize inter-die thermal coupling.

In this Chapter, thermal limits and opportunities of heterogeneous 3-D integration architectures are explored and summarized. To solve these challenges, a novel 3-D IC stack architecture using interposer-embedded microfluidic cooling in conjunction with thermal isolation technologies is proposed. In order to evaluate the thermal benefits of the proposed stacks, a thermal modeling framework is developed based on the finite volume method and validated against *ANSYS*. The 3-D stack architecture is compared to other two baseline architectures and then analyzed as a function of TSV/microbump diameter, TSV/microbump number, TSV layout and die thickness.

2.1 Thermal modeling framework

In an IC package stack, there are multiple layers with heterogeneous materials. As the trend to integrate more chips in a single package, the IC stack becomes more complex. Thus, it is time consuming to simulate the whole stack using Finite Element Method software such as *ANSYS* to get high order of accuracy. Modified *hotspot* [32, 35] with an equivalent thermal resistance model decreases the complexity of the geometry, but reduces the solution accuracy. Other Modeling methods based on frequency domain computation, Greens function and cosine or sine transforms [36] are faster but their extension to heterogeneous stacks with non-uniform materials is difficult.

In this Section, we present a thermal modeling framework leveraging the advances of several prior efforts and maintaining a maximum relative error of less than 7% with a maximum speedup of 22X than *ANSYS*. The thermal framework is capable of performing both steady-state and transient-state analysis. Moreover, the model considers different cooling solutions including a conventional air-cooled heatsink as well as a microfluidic cooled heatsink.

2.1.1 Formulation using finite volume methods

The general heat transfer equation is expressed as follows,

$$\nabla \cdot (K(x, y, z) \cdot \nabla T(x, y, z)) + \rho \cdot c_p \cdot \frac{\partial T(x, y, z)}{\partial t} = P(x, y, z) \quad (2.1)$$

where $K(x, y, z)$ and $T(x, y, z)$ denotes the thermal conductivity ($W/m \cdot ^\circ C$) and temperature ($^\circ C$) respectively, ρ is the mass density (Kg/m^3), c_p is the specific heat capacity ($J/Kg \cdot ^\circ C$), and $P(x, y, z)$ is the power excitation density (W/m^3).

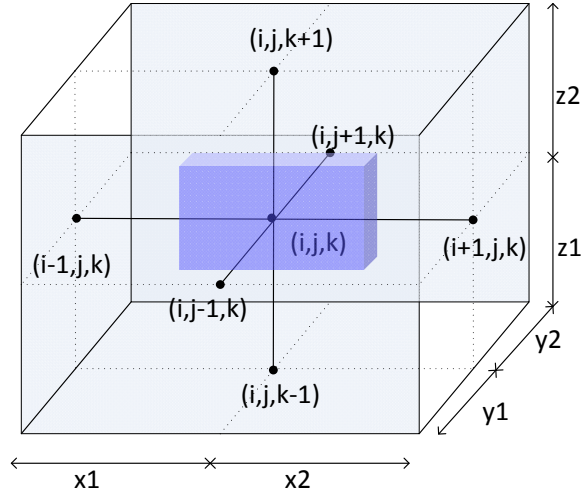


Figure 5: Illustration of finite volume scheme

According to the finite volume methods used in [28, 37], we have the following finite

difference expressions for the nodes inside the stack, which are shown in Fig. 5:

$$\begin{aligned} & \frac{T_{(i,j,k)} - T_{(i-1,j,k)}}{\frac{x_1}{k_x \cdot l_y \cdot l_z}} + \frac{T_{(i,j,k)} - T_{(i+1,j,k)}}{\frac{x_2}{k_x \cdot l_y \cdot l_z}} + \frac{T_{(i,j,k)} - T_{(i,j-1,k)}}{\frac{y_1}{k_x \cdot l_x \cdot l_z}} + \frac{T_{(i,j,k)} - T_{(i,j+1,k)}}{\frac{y_2}{k_x \cdot l_x \cdot l_z}} + \\ & \frac{T_{(i,j,k)} - T_{(i,j,k-1)}}{\frac{z_1}{k_x \cdot l_x \cdot l_y}} + \frac{T_{(i,j,k)} - T_{(i,j,k+1)}}{\frac{z_2}{k_x \cdot l_x \cdot l_y}} + \rho \cdot c_p \cdot V \cdot \frac{\partial T(x, y, z)}{\partial t} = P_{total} \end{aligned} \quad (2.2)$$

where $l_x = (x_1 + x_2)/2$, $l_y = (y_1 + y_2)/2$, $l_z = (z_1 + z_2)/2$; P_{total} is the total power consumption in the shaded cube; $k_{x,y,z}$ is the thermal conductivity along each axis.

In the IC stack, there are multiple materials in each layer. Detailed and homogeneous meshing that guarantees only one material per mesh will increase the problem size and without carefully handling the boundaries between two materials, such case usually results in poor convergence. Therefore, to maintain simulation accuracy as well as increase the efficiency, a hybrid meshing strategy is employed.

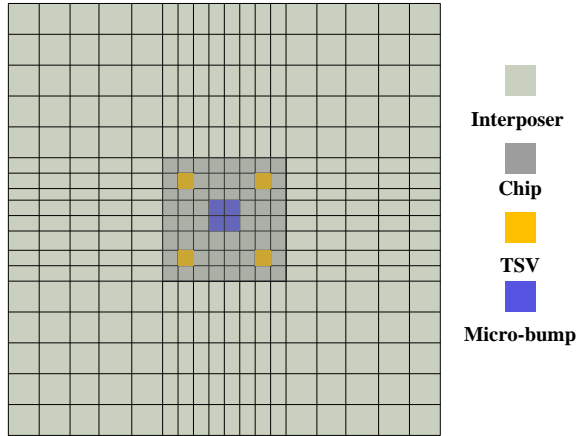


Figure 6: Illustration of non-conformal meshing.

First, we use non-conformal meshes in different domains such as the package, interposer and chips. To avoid divergence issues, we gradually increase the mesh size in the transition area between different domains, as shown in Fig. 6. Second, in each domain, we use effective thermal conductivity modeling methods for the meshes containing more than one material. For example, on-chip and package metal layers are modeled using in-plane

and through-plane thermal conductivity formulated in [38]. Moreover, effective thermal conductivity modeling of the layers with ‘vertical interconnects’ (microbumps, TSVs, etc.) [38] is implemented to further reduce the mesh number.

2.1.2 Modeling of air- and microfluidic-cooled heat sinks

In our thermal modeling, we treat all air-cooled heat sinks as convective boundary conditions, which is applied to an explicitly modeled heat spreader. Therefore in the finite volume method, for the boundary nodes, we apply the convective boundary condition:

$$K \cdot \frac{\partial T}{\partial n} \Big|_{boundary} = -h(T - T_{amb}) \quad (2.3)$$

where, h is convection heat transfer coefficient and n is the direction of the normal vector of the boundary.

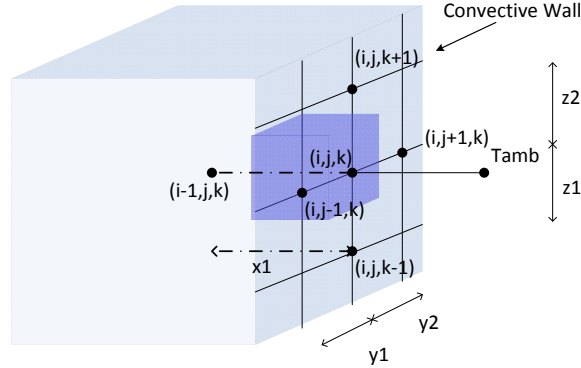


Figure 7: Illustration of boundary nodes

The corresponding finite difference expression is shown below and sketched in Fig. 7.

$$\begin{aligned} & \frac{T_{(i,j,k)} - T_{(i-1,j,k)}}{\frac{x_1}{k_x \cdot l_y \cdot l_z}} + \frac{T_{(i,j,k)} - T_{amb}}{\frac{1}{h \cdot l_y \cdot l_z}} + \frac{T_{(i,j,k)} - T_{(i,j-1,k)}}{\frac{2 \cdot y_1}{k_x \cdot l_x \cdot l_z}} + \frac{T_{(i,j,k)} - T_{(i,j+1,k)}}{\frac{2 \cdot y_2}{k_x \cdot l_x \cdot l_z}} + \\ & \frac{T_{(i,j,k)} - T_{(i,j,k-1)}}{\frac{2 \cdot z_1}{k_x \cdot l_x \cdot l_y}} + \frac{T_{(i,j,k)} - T_{(i,j,k+1)}}{\frac{2 \cdot z_2}{k_x \cdot l_x \cdot l_y}} + \rho \cdot c_p \cdot V \cdot \frac{\partial T(x, y, z)}{\partial t} = P_{total} \end{aligned} \quad (2.4)$$

While in the modeling of microfluidic cooled heat sink, the heating of fluidics cannot be ignored especially when the fluidic velocity is lower than 50 ml/min. In this case, to model

the thermal interactions between the fluidics and the chip, we added the energy balance term described in [39] into our finite difference scheme, as shown below (assuming flow direction in Y-axis):

$$\rho \cdot c_p \cdot v \cdot l_x \cdot l_z \cdot (T_{(i,j,k)} - T_{(i,j-1,k)}) = Q_{heat} \quad (2.5)$$

where, v is the volumetric flow rate, Q_{heat} is the heat carried out in the flow direction.

Finally we use backward Euler scheme to numerically model the time varying terms, as shown below:

$$\rho \cdot c_p \cdot V \cdot \frac{\partial T(x, y, z)}{\partial t} = \rho \cdot c_p \cdot V \cdot \frac{T(x, y, z, t + \delta t) - T(x, y, z, t)}{\delta t} \quad (2.6)$$

2.1.3 Comparison with existing thermal models

Various thermal models have been developed to analyze the thermal profiles of IC packages. We compare our thermal model with existing work using a variety of capabilities including: formulation methods, solving algorithms, steady-state analysis, transient analysis, and capability to model the layers with heterogeneous materials and multi-die packages. The comparison is summarized in Table 1.

Based on the research objectives, the thermal models have different focuses. For the purpose of using the models as an architectural evaluation engine, compact models with ultra-high efficiency are preferred [32, 40]. Such models are usually implemented with a lot of lumped elements, especially for the secondary heat path. For the work focusing on the on-die power-thermal activities [41, 20], the models mainly focus on handling the fine-granularity layout and use an abstracted model to handle the layers with multiple materials. For the work concentrating on the impact of microfluidic cooling from architecture-level [42, 31], they have an emphasis on the interactions between microfluidic cooling and IC packages. Last, there are also thermal efforts focusing on developing co-design frameworks

Table 1: Comparison of different thermal modeling

	Methods	Solving	Steady-state	Transient	Multi-materials layer	secondary heat path	Heatsink	Multi-Die	Research objective
<i>ANSYS</i>	Finite element	Iterative	Yes	Yes	Yes	Yes	Yes	2.5D/ 3D	Detailed thermal centric study
<i>hotspot</i> [32]	Thermal resistance	Iterative	Yes	Yes	Weighted average	Lumped model	Yes	3D	Thermal engine for architecture tools
J. Xie [37]	Domain decomposition	Direct Iterative	Yes	Yes	Yes	Yes	Yes	2.5D/ 3D	Thermal and IR-drop co-analysis
H. Oprins [20]	Data-fitting	Closed form	Yes	No	No	Lumped	Lumped	2.5D/ 3D	Fast thermal simulation for simple package designs
Z. Wan [42]	Thermal resistance	Direct Iterative	Yes	Yes	Not mentioned	Yes	Lumped	3D	Thermal simulator for microfluidic cooled systems
A. Sridhar [31]	Finite difference	Direct Iterative	Yes	Yes	Not mentioned	Lumped	Yes	3D	Compact models for microfluidic cooled systems
A. Ziabari [40]	Power blurred	Convolution	Yes	Yes	No	Lumped	Lumped	3D	Thermal engine for architectural study
Krit. A [41]	Finite element	Iterative	Yes	No	Weighted average	Lumped	Lumped	3D	Thermal engine for floor-planner
This work	Finite volume	Direct Iterative	Yes	Yes	Yes	Lumped	Yes	2.5D/ 3D	Design exploration for 2.5-D and 3-D integration

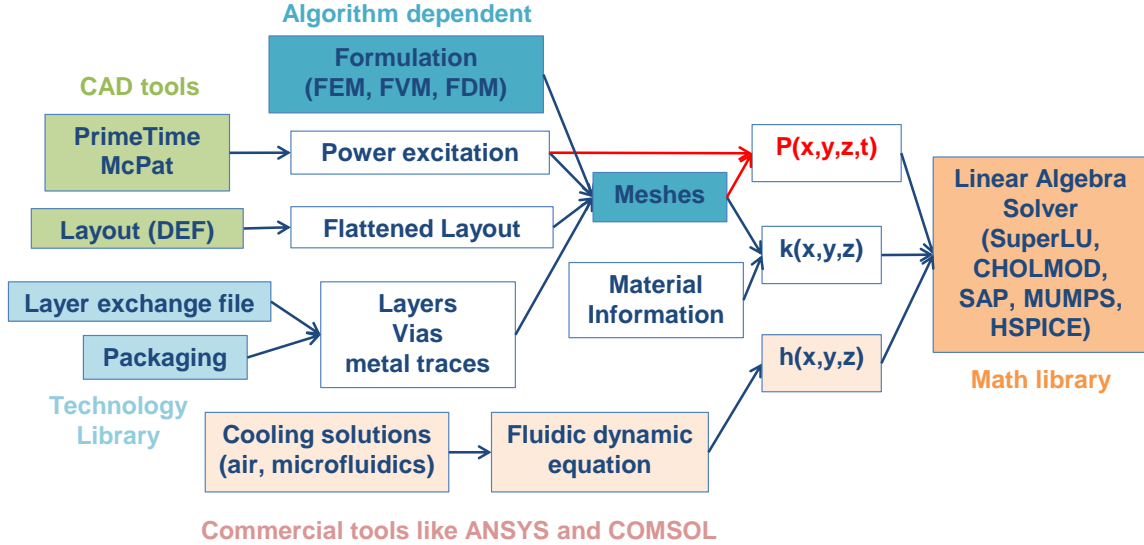


Figure 8: Thermal simulation flow.

to model the interaction between temperature and other system metrics, such as IR-drop [28, 33, 36], such models usually concentrate on steady state analysis to reduce the problem complexity.

The thermal model of this work leverages the advances of several existing models. For each IC package domain, it can easily change the modeling details to meet different simulation requirements in a reasonable computation time compared to other models. Moreover, the model is capable of performing both steady-state and transient analysis for 2.5-D and 3-D integration platforms with either air-cooled or microfluidic-cooled heat sinks.

2.1.4 Thermal simulation flow and validation

The thermal simulation flow is shown in Fig. 8. The analysis tool has three inputs: first, the power consumption of each functional block in the chip along with the layout; second, the geometry and material information in the chip stack, which is usually from a technology library; and third, the boundary conditions characterized from commercial fluidics module of FEM software.

After we obtain $P(x, y, z, t)$, $k(x, y, z)$, $h(x, y, z)$, we can build a linear equation for

steady state analysis, as follows,

$$Y \cdot T = b \quad (2.7)$$

Since Y is symmetric positive definite (SPD), the above equation can be solved using Cholesky factorization or iterative methods such as preconditioned conjugate gradients method. While for the transient analysis, the equation becomes:

$$Y \cdot T + C \cdot \dot{T} = b \quad (2.8)$$

By using backward Euler scheme, the equation is then derived as follows,

$$\begin{aligned} \dot{T} &= \frac{T^{n+1} - T^n}{\Delta t} \\ (Y + \frac{C}{\Delta t}) \cdot T^{n+1} &= b - \frac{C}{\Delta t} \cdot T^n \end{aligned} \quad (2.9)$$

$\frac{C}{\Delta t}$ is a diagonal matrix, thus $Y + \frac{C}{\Delta t}$ is still SPD and the aforementioned algorithms are still applicable. Next, we validate the thermal framework in the following aspects: steady-state, transient-state, and microfluidic modeling.

Steady-state

Fig. 9(a)(b) shows an example 3D stack that was used to validate the steady-state analysis model with ANSYS. The power map of each of the stacked chips is shown in Fig. 9(c)(d). All surfaces are adiabatic except for the top surface, which is defined to have a convection heat transfer coefficient of $40,000 \text{ W}/^\circ\text{C} \cdot \text{m}^2$. The chip size is $1 \text{ cm} \times 1 \text{ cm}$. To reduce the meshing and analysis complexity in ANSYS, we only use 400 uniformly distributed TSVs between the two dice in this validation example. The TSV diameter is $50 \mu\text{m}$, and we assume there is no liner (again, to simplify ANSYS meshing). The thickness of both dice is $50 \mu\text{m}$ and the bonding layer is $5 \mu\text{m}$ thick. The thermal maps of both dice using ANSYS and the thermal model are shown in Fig. 10 and match to within a maximum

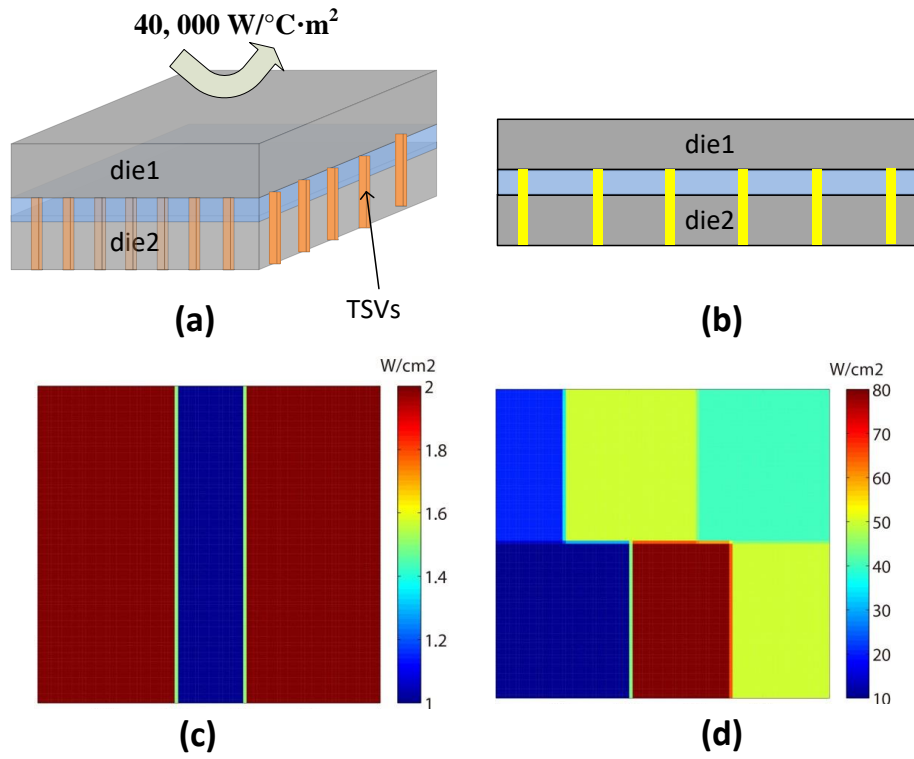


Figure 9: Steady-state validation setup.

relative error of 7% for this example.

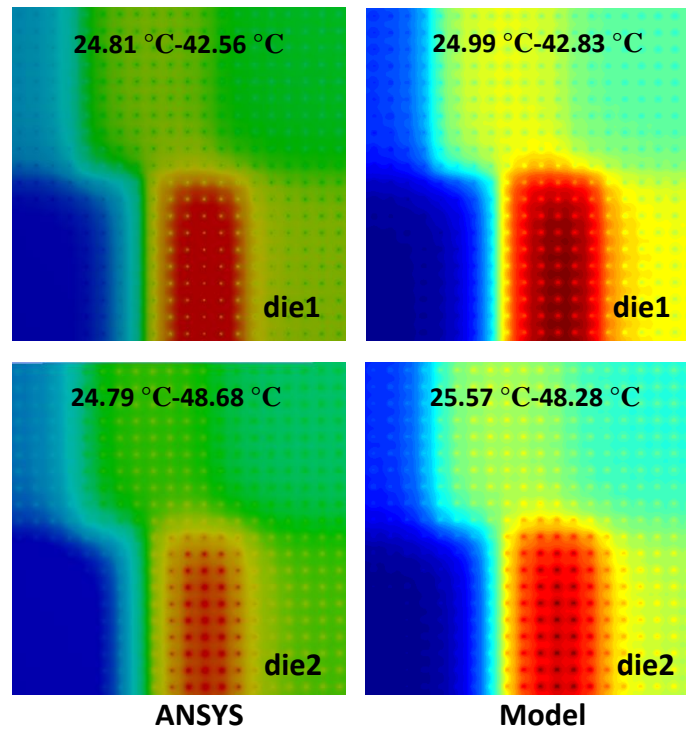


Figure 10: Steady-state validation results.

ANSYS generates meshes very accurately even within the TSVs and the meshes are in a fine granularity in the interfaces between two different materials and layers. Nevertheless, our model still maintains adequate accuracy in maximum temperature and minimum temperature.

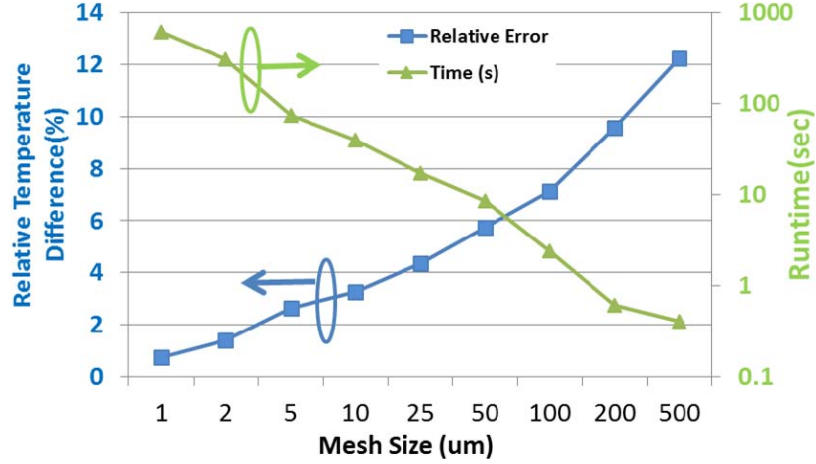


Figure 11: Impact of mesh size: tradeoff between efficiency and accuracy.

Table 2: Efficiency and accuracy comparison

Min Mesh Size (μm)	Difference of <i>ANSYS</i> (%) ¹	Difference of Model	Time of <i>ANSYS</i> (s)	Time of Model (s)
1	0	0.73	1800	615
25	4.81	4.34	321	17
50	7.15	5.72	192	8.6
100	N.A	7.13	N.A	2.4

¹ The error is normalized to the results of *ANSYS* using a 1 μm minimum mesh size

The maximum error occurs in the transitioning area where meshes may not be present in our thermal model. If we polish our meshes to the granularity of *ANSYS*, the accuracy can be further improved at the expense of computational efficiency, as shown in Fig. 11. If we reduce the mesh size to 1 μm , the relative error is less than 1% but the runtime increases by 120 times compared to the default mesh size (100 μm) we use. Nevertheless, our model shows efficiency advantages compared to *ANSYS*, as shown in Table. 2. The model achieves 3X, 18.9X and 22X speedups with comparable accuracy for mesh size of 1 μm , 25 μm and 50 μm , respectively.

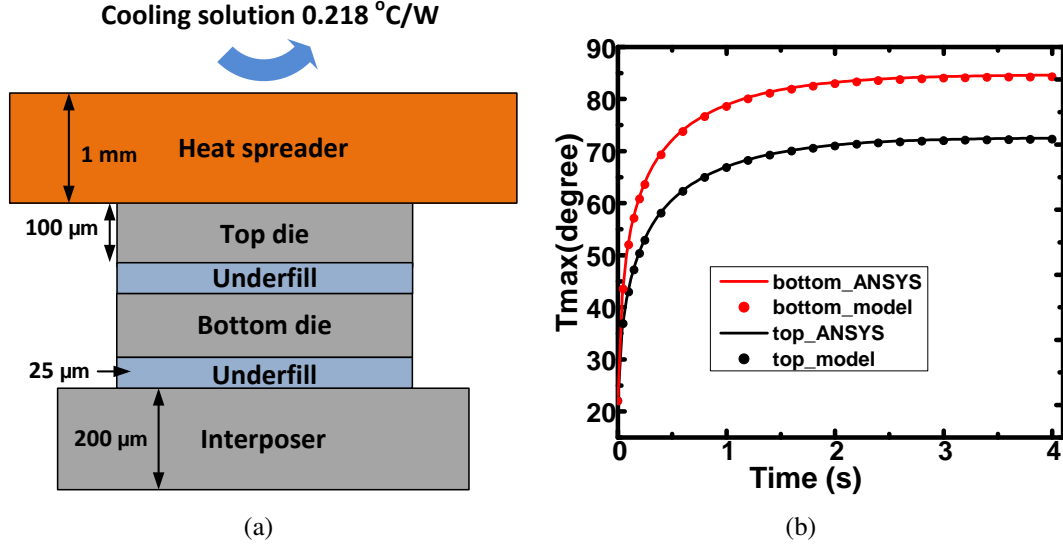


Figure 12: Transient thermal validation experiments (a) a 2-die 3-D stack (b) transient thermal validation results

Transient-state

Fig. 12(a) shows a 3-D stack with two dice used for transient thermal validation against ANSYS. Non-uniform power maps with average power densities of 36 W/cm^2 and 43.8 W/cm^2 are assigned to the top and bottom dice, respectively. The power maps are similar to Fig. 75(b) with appropriate scaling. The heat spreader size is assumed to be $3 \text{ cm} \times 3 \text{ cm}$ with a cooling of 0.218°C/W added to the top surface. Other faces of the stack are assumed to be adiabatic. The interposer size is assumed to be $2 \text{ cm} \times 1.5 \text{ cm}$ and the chip size is set as $1 \text{ cm} \times 1 \text{ cm}$. The thickness is labeled in Fig. 12(a). The starting temperature of the whole stack is assumed to be the ambient temperature, which is 22°C . We add the power excitation from time equal to 0 seconds and perform the transient thermal analysis from 0 to 4 seconds. The results of the maximum temperature of each die are shown in Fig. 12(b). The maximum temperature of both dice in each time point matches ANSYS results with an error of less than 1%. Moreover, the thermal profiles of each time point are compared and the maximum error is also less than 1%.

Microfluidic modeling

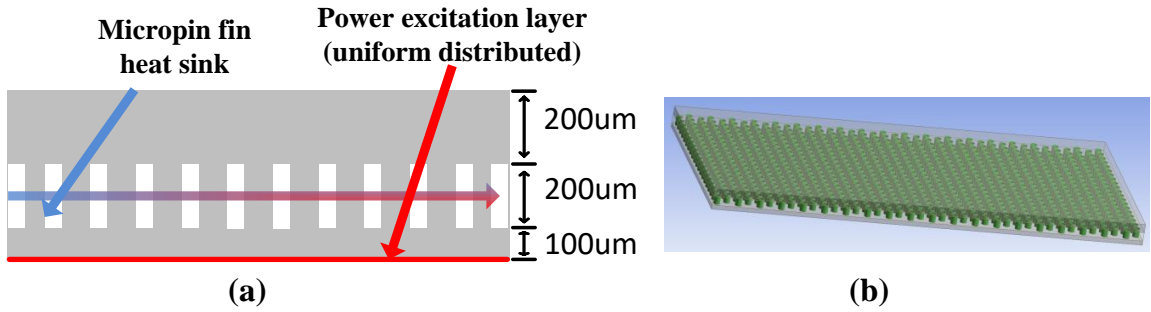


Figure 13: (a) The cross-sectional view of the experimental setup (b) the full chip view of the testbed.

Fig. 13 shows the simulation testbed in both ANSYS and the model. The micropin fins have a diameter of $150 \mu m$ and a pitch (x-, y-axis) of $250 \mu m$ pitch. We vary the power density of the excitation layer from $40 W/cm^2$ to $100 W/cm^2$ and investigate the error between the ANSYS simulation results and the model results. The flow rate is $100 ml/min$.

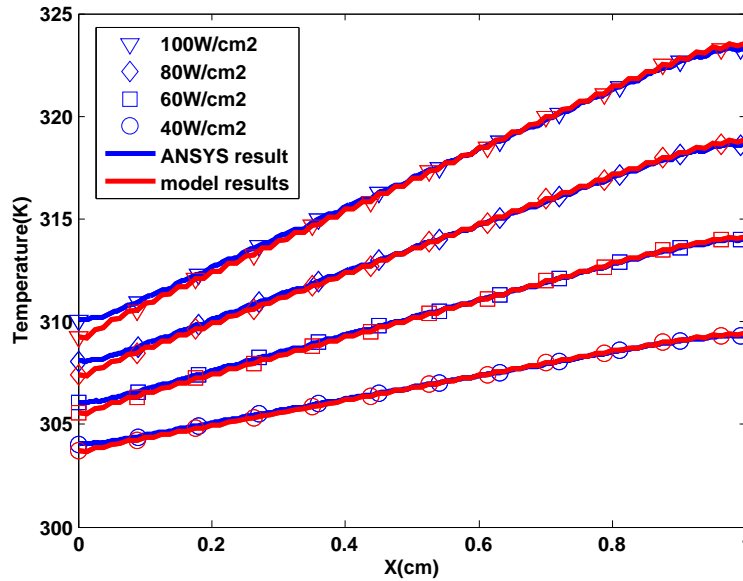


Figure 14: Comparison between the model and ANSYS simulation results in the power excitation layer along $y = 0.5 cm$

Fig. 14 shows the temperature along the center line ($y = 0.5 cm$) in the power excitation layer. The maximum error of the four cases ($40 W/cm^2$ to $100 W/cm^2$) is 9.09%, 9.02%,

8.85%, 8.84%, respectively. However, for the outlet ports, the maximum error is relatively small and less than 1% for all four cases.

2.2 Motivation of thermal isolation

There are two distinct thermal challenges in 3-D ICs with each requiring separate optimization and technology solutions: first, stacking dice in 3-D increases the total power density while simultaneously increases the thermal resistance of dice within the stack to the atop attached heat sink; second, stacked dice will experience unwanted thermal crosstalk, particularly between high-power dice and low-power temperature sensitive components.

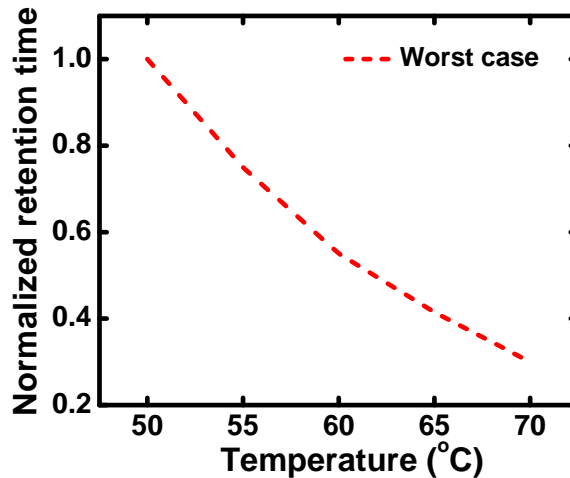


Figure 15: DRAM retention time reduces exponentially as the temperature goes up

For instance, in a DRAM-on-processor 3-D stack, the DRAM usually exhibits a thermal profile similar to the processor die (i.e., a thermal map mirror image) and has a relatively high temperature due to the strong thermal coupling even though DRAM itself dissipates much lower power than the processor [43]. However, higher DRAM junction temperature (say, 90°C) increases the DRAM refresh rate (as shown in Fig. 15 [44]), leading to performance degradation of 8.6% and power consumption overhead of 16.1% [45]. Therefore, the memory die should be thermally decoupled from the processor die. Likewise, in the domain of silicon nanophotonics, a number of components are sensitive to temperature variations. For example, microring resonators are highly temperature sensitive, and thus

complex stabilizer circuits and dense heaters are used to compensate for this thermal variation [46]. With a temperature drift of $8\text{ }^{\circ}\text{C}$, the tuning power is as high as 0.19 nJ/bit (26.7% of total communication power) [47] and will increase under a higher temperature drift.

Tier-specific cooling has been proposed to provide isolation between dice for 2.5-D integration platforms [48, 49]. With a strong cooling solution assigned to each die, the heat spreading between dice is confined. Furthermore, thermally resistive material such as glass interposer enhances the isolation between dice [48]. The simulation results show that the temperature of the low-power die is 40% lower than the high-power die. The prior research focuses on 2.5-D systems, and if the proposed method is used for 3-D, the thermal coupling will not be suppressed since the heat path is in the vertical direction. Porous silicon is proposed to isolate photonics die from the processor die in a 3-D stack [50]. Due to the lower conductivity of porous silicon (100X smaller than silicon), this 3-D integrated system achieves 3.8 ~ 5.4x ring heating power reduction and 11 ~ 23% speedup compared to the baseline without isolation. However, this research ignores the impact of TSVs and the discussion of mechanical reliability is missing. Moreover, the isolated die experiences an elevated temperature, which is not wanted for applications such as memory-on-processor stack.

From the above discussions, it is clear that thermal isolation in a 3-D stack is important for many applications, in particular when decoupling the heating of low-power and temperature sensitive devices from high-power IC.

2.3 Proposed architecture with thermal isolation technologies

To address the thermal crosstalk and cooling needs, the novel 3-D stack architecture shown in Fig. 16 was proposed [51]. The proposed architecture has three novel features. First, a microfluidic heat sink (MFHS) is integrated in the interposer. The purpose of this heat sink is to cool the processor and the extended heat spreader. Second, an air gap is in-

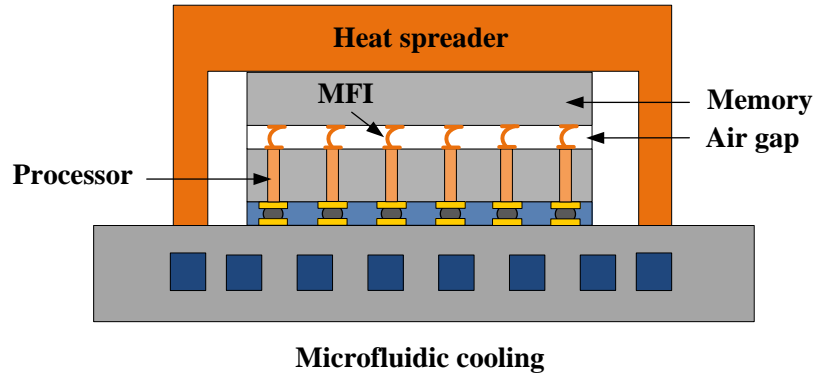


Figure 16: Proposed architecture with interposer-embedded heat sink, thermal bridge and air gap isolation

egrated between the high-power and low-power dice to decouple the thermal crosstalk. Mechanically flexible interconnects (MFIs) are used as chip to chip interconnections [52]. Unlike rigid solder microbumps, MFIs can deform elastically under stress and maintain the electrical connectivity between the memory and the processor tiers. Due to this behavior, MFIs can help eliminate underfill. Finally, a heat spreader is attached on top of the isolated low-power die to provide a cooling path. This architecture is suitable for a wide range of heterogeneous 3-D stacks of high power and low power devices. In this study, we focus on a memory-on-processor stack.

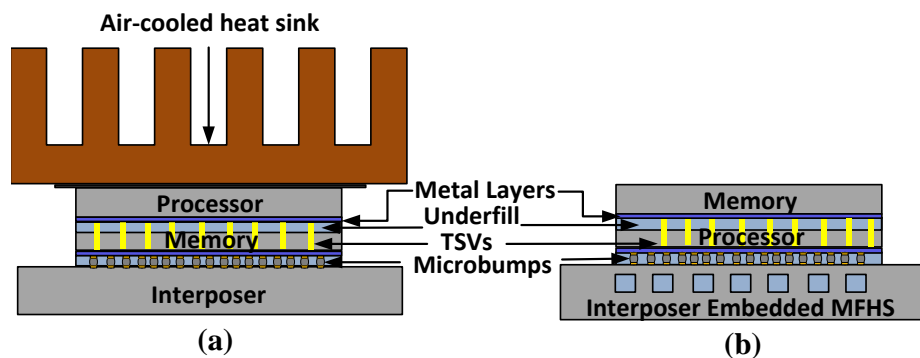


Figure 17: 3-D stack (a) with conventional air cooled heat sink (b) with interposer embedded microfluidic heat sink (MFHS)

As a comparison, Fig. 17 shows two typical 3-D stacks with different cooling solutions. The first 3-D stack is based on current approaches in the literature in which an air-cooled heat sink (with heat spreader) is attached on top of the stack. Since the heat sink is on top,

the thermally optimal architecture for this stack is to place the processor above the memory.

The second 3-D stack is cooled using a microfluidic-cooled interposer [53], as shown in Fig. 17(b). In this case, the processor is on the bottom to minimize the thermal resistance between the processor and the heat sink. Even if the microfluidic-cooled interposer can lower the system temperature compared to the air cooled heat sink, the thermal coupling issue is still unsolved because there are no mechanisms to prevent heat transfer between the two dice.

2.4 Thermal evaluation of the proposed architecture

In this section, we use the thermal framework described above to thermally evaluate the proposed architecture and compare it to other two baseline architectures.

2.4.1 Thermal specification of 3-D stacks

Table 3: The specification of simulated stack

	Conductivity W/mK	Thickness μm
TIM	3	25
Memory die	149	100
Underfill layer	0.9	5
Air gap	0.024	5
Processor die	149	50
Micro-bump	60	40
Interposer	149	200
Copper	400	N/A
SiO_2	1.38	N/A

Table 3 lists the thickness and material of each layer used in the 3-D stack. The chip size is assumed to be $1\text{ cm} \times 1\text{ cm}$. The interposer is set to be $3.5\text{ cm} \times 3.5\text{ cm}$. The interposer embedded microfluidic heat sink is assumed to be the same size as the chip. In our thermal modeling, all heat sinks are treated as convective boundary conditions. The thermal resistances of the air cooled heat sink and MFHS are assumed to $0.5\text{ K} \cdot \text{cm}^2/W$

and $0.2 \text{ K} \cdot \text{cm}^2/\text{W}$, respectively. The ambient temperature is set to be $25 \text{ }^\circ\text{C}$.

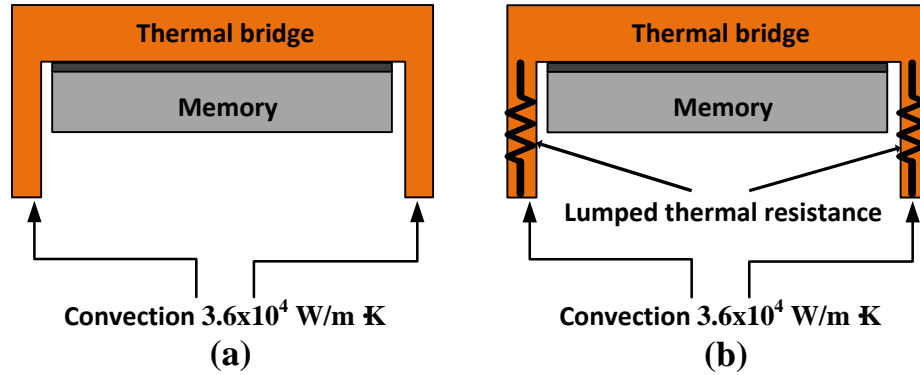


Figure 18: (a) Physical structure of the extended heat spreader (b) Lumped resistance modeling for fins of extended heat spreader and TIM

Without an effective thermal path for the isolated die, the temperature of the isolated die may be relatively large. In Fig. 17, this need is addressed using the extended heat spreader, which can be formed using a modified copper spreader. Fig. 18(a) shows the physical structure of the thermal bridge. The top surface of the copper thermal bridge is $1.5 \text{ cm} \times 1.5 \text{ cm}$ with a thickness of $500 \text{ } \mu\text{m}$ (assuming chip size is $1 \text{ cm} \times 1 \text{ cm}$). A convective boundary condition of $3.6 \times 10^4 \text{ W/m}^2 \cdot \text{K}$ is applied. To simplify the structure, we model the bridge fins and TIM (attaching the bridge to the interposer) as lumped thermal resistors shown in Fig. 18(b); the width of the fin (2 mm) justifies this simplification [54].

The micro-bumps and TSVs are assumed to be uniformly distributed throughout the chip. The default diameter of micro-bumps is $40 \text{ } \mu\text{m}$ and the total number is 1,600. It is assumed there are 10,000 TSVs with a diameter of $5 \text{ } \mu\text{m}$ and a liner thickness of $0.5 \text{ } \mu\text{m}$.

Fig. 19 illustrates the power maps of the memory and processor dice. The memory die layout is based on an 8 Gb 3-D DDR3 DRAM design from *Samsung* [55]. The total power is estimated from the *Micron* DDR3 DRAM datasheet, which gives a value of 2.82 W. The layout of the processor die is based on the *Intel* Core i7 microprocessor [56]. The power is assumed to be 74.63 W. According to estimation from *McPAT* [57], the power distribution is assigned as shown in Fig. 19(b).

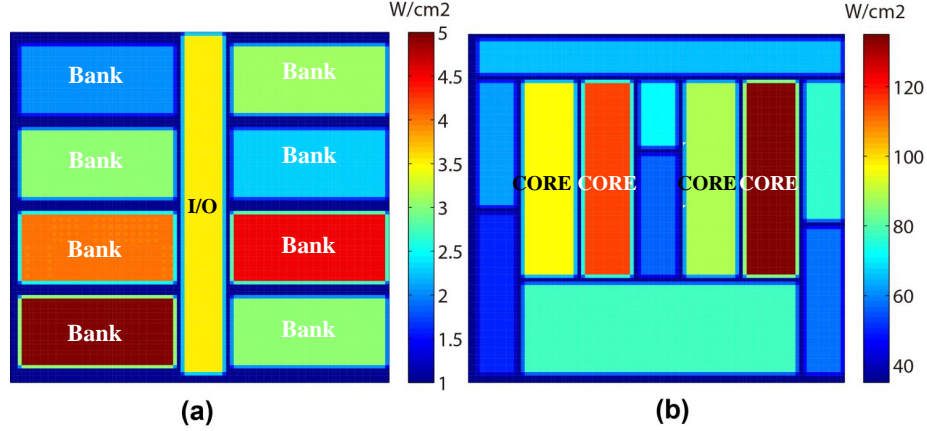


Figure 19: Power density distribution: (a) Memory die (b) Processor die

2.4.2 Comparison of different 3-D stacks

The two baseline stacks are shown in Fig. 17. They are configured with an air-cooled heat sink and an interposer embedded microfluidic heat sink, respectively.

Because the TSVs may influence the thermal decoupling effectiveness of the air-gap, the proposed stack with and without TSVs has been studied to give a worst and best case estimation. For all stacking scenarios, two conditions were studied: processor in standby mode (24.63 W total power) and an active processor (74.63 W total power). When the processor jumps from standby to an active mode, the temperature will increase, and it is important to see how this variation influences the temperature of the memory die.

Table 4: The comparison of different architectures

Unit: $^{\circ}C$	T_{max} (Memory)		T_{max} (Processor)	
	Standby	Active	Standby	Active
Stack with Air cooled heat sink	50.33	75.06	51.63	76.44
Stack with interposer embedded MFHS	46.59	65.38	47.14	66.05
Proposed stack w/o TSVs	31.96	39.63	47.58	64.64
Proposed stack w TSVs	38.88	51.76	44.75	61.44

Table 4 illustrates the maximum temperature of each die under the two different processor modes. From the results, there are three key conclusions. First, the proposed architecture has the lowest temperature in both standby and active modes. In the active state,

the maximum temperature of the three stacks (first three rows) is $76.44\text{ }^{\circ}\text{C}$, $66.05\text{ }^{\circ}\text{C}$ and $64.64\text{ }^{\circ}\text{C}$, respectively. Second, the proposed architecture decouples the heat from the processor to the memory. For the first two scenarios, in both active and standby modes, the memory exhibits a temperature similar to that of the processor. For the proposed stack with thermal isolation, in the standby mode, the memory temperature is $31.96\text{ }^{\circ}\text{C}$, while the processor temperature is $47.58\text{ }^{\circ}\text{C}$; in the active mode, the memory temperature is only $39.63\text{ }^{\circ}\text{C}$ even though the processor temperature is as high as $64.64\text{ }^{\circ}\text{C}$. Third, our proposed stack maintains the memory die temperature fairly independently of the processor die mode. In the first two cases, when the processor transitions from the standby mode to the active mode, the temperature of the memory die increases by $25\text{ }^{\circ}\text{C}$ and $20\text{ }^{\circ}\text{C}$, respectively. For our proposed stack, it increases only by $8\text{ }^{\circ}\text{C}$ when the processor transitions to active mode.

However, when TSVs are inserted in the proposed stack, the thermal coupling increases, as expected, compared to the TSV-free case. In both processor modes, the temperature of the memory die is always higher than the case without TSVs, which indicates coupling. The thermal map of each die is shown in Fig. 20 when the processor die is in active mode. In Fig. 20(b), the temperature distributions of both dice are similar to each other, and the temperature difference of the two dice is only $10\text{ }^{\circ}\text{C}$ compared to $25\text{ }^{\circ}\text{C}$ in Fig. 20(a). Thus, the TSVs clearly impact the thermal isolation, and we will discuss this further later in the Chapter.

2.5 Design space exploration of the proposed architecture

In this section, we thermally study our proposed 3D stack architecture as a function of the cooling capability of the thermal bridge, TSV/microbump diameter, TSV/microbump density, TSV layout and die thickness. Through this analysis, the benefits, limits and challenges of our proposed architecture can be better understood. If not specified, the parameters and power maps are the same as those used in the last section.

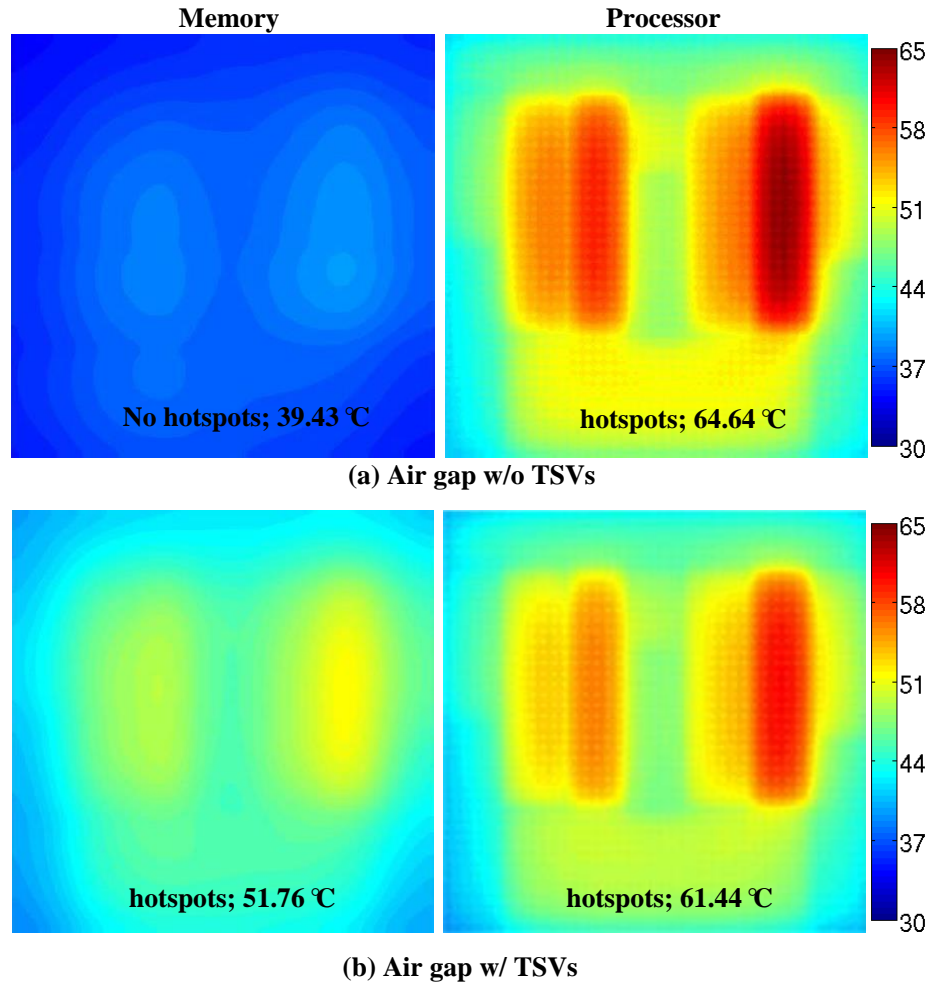


Figure 20: Thermal maps of proposed stack when processor is in active mode (a) without TSVs (b) with TSVs

2.5.1 Impact of Microbump

Micro-bumps affect the primary heat path. Their diameter and number influence the equivalent thermal resistance between the processor and the interposer.

Firstly, we fix the diameter of the microbump to $40 \mu m$ and change the total number of microbumps from 1,600 to 10,000. Secondly, we fix the total number of microbumps to 1,600 and change their diameter from $10 \mu m$ to $50 \mu m$. We also evaluate the cases where there are no TSVs and where there are 2,500 TSVs between the processor and the memory dice.

The results are shown in Fig. 21. As expected, more microbumps or larger diameter

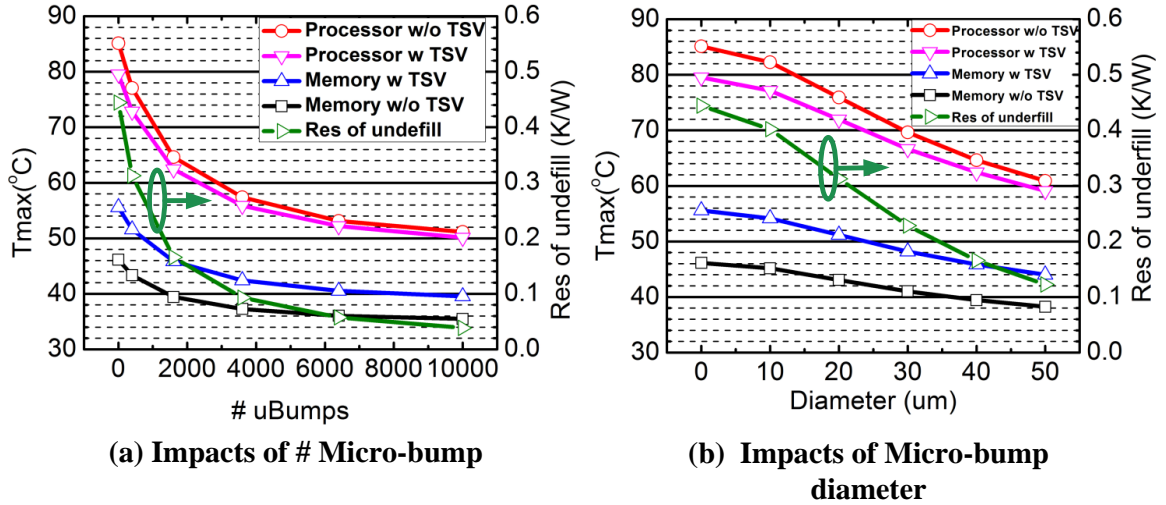


Figure 21: The impact of the microbumps. (a) change the number of microbumps (b) change the diameter of microbumps.

of microbumps lead to a decrease of the equivalent thermal resistance of the layer, which improves the primary heat path. With 3,600 $40 \mu m$ diameter microbumps, the temperature of the processor die is below $60^\circ C$, which is tolerable. In reality, we have very fine pitch ($50 \mu m$) electrical micro-bumps between the chip and the interposer, which leads to a total of 40,000 microbumps; thus the thermal requirement of the micro-bumps can be easily met.

2.5.2 Impact of TSVs

3-D stacks require a large number of TSVs for inter-die signaling and power delivery. These TSVs will bridge the air gap and reduce its equivalent thermal resistance (which is undesirable). The inclusion of TSVs causes the temperature of the memory die to closely track that of the processor die. The diameter and number of TSVs impact the equivalent thermal resistance of the thermal isolation air gap. In order to quantify this impact, two experiments were performed. First, the TSV diameter is fixed to $5 \mu m$ with a $0.5 \mu m$ liner and the TSV number is swept from 1,600 to 10,000. Next, the total number of TSVs is fixed at 10,000 and the TSV diameter is swept from $2 \mu m$ to $10 \mu m$ with a fixed $0.5 \mu m$ liner.

Fig. 22(a) and (b) show the impact of TSV number and diameter, respectively. As the

TSV total volume increases, the air-gap isolation layer becomes more thermally conductive, and the inter-die heat coupling becomes stronger, reducing the temperature difference between the two dice. If $2\ \mu\text{m}$ diameter TSVs are used rather than $10\ \mu\text{m}$, the memory temperature is only $44\ ^\circ\text{C}$, compared to $54.5\ ^\circ\text{C}$ for the $10\ \mu\text{m}$ TSV case. Further scaling of the TSV dimensions will yield additional improvements.

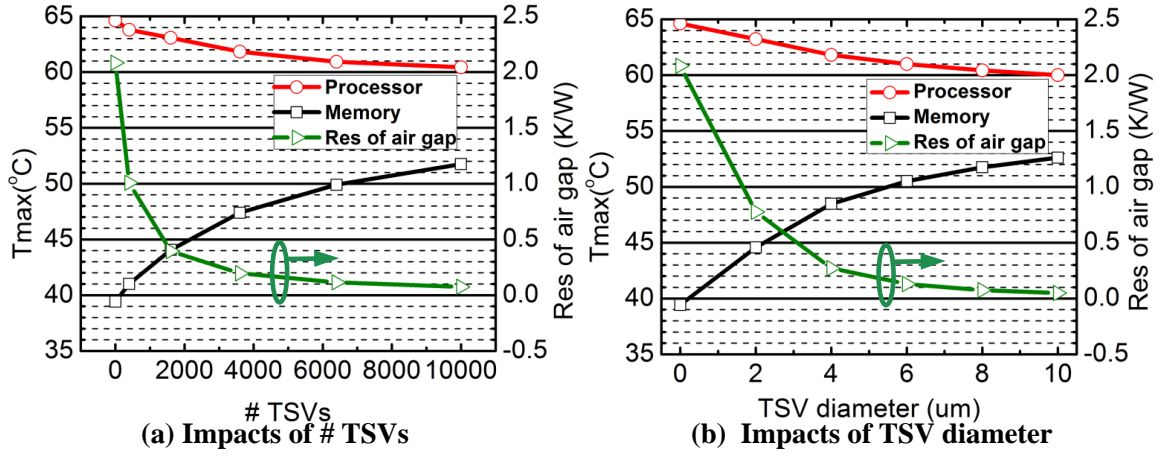


Figure 22: The impact of the TSVs. (a) change the number of TSVs. (b) change the diameter of TSVs.

In 3-D DRAM stacks or wide I/O applications, TSVs are usually clustered towards the center of the die (rather than being uniformly distributed across the die surface). When the TSVs are clustered in a specific area, thermal coupling is expected to occur only in that area. In this way, the heat from the processor die will be localized. For the memory die, the clustered TSVs or “TSV farm” usually acts as I/O pins, and are probably outside of the memory cell circuits (labelled by a dashed-line box in Fig. 23(a)) [55]. Hence, the memory cell circuits will become relatively free from the impact of the processor because there are no TSVs in their area. Inspired by the above analysis, the TSVs are clustered only in the center; an area of $1\ \text{mm} \times 5\ \text{mm}$ is assumed containing 49×100 TSVs, which are labeled by the solid-line rectangle in Fig. 23(a). To make a fair comparison, a uniformly distributed TSV case with 4,900 TSVs is considered. The results are shown in 23(b).

In the clustered TSV case, the maximum temperature of the whole DRAM die drops by $3.55\ ^\circ\text{C}$ compared to the uniform TSV case. However, the maximum temperature of the

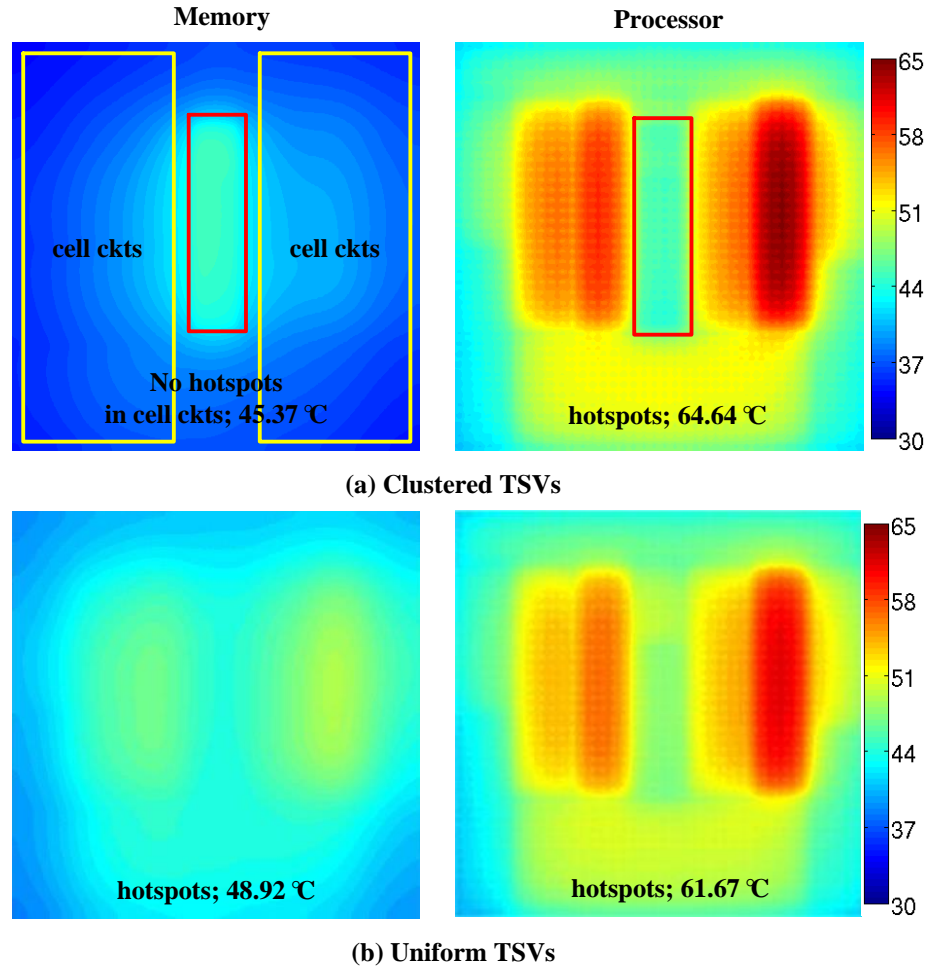


Figure 23: Thermal maps of clustered TSVs and uniform TSVs (a) TSVs are clustered in the solid-line rectangle (b) The same amount of TSV are uniformly distributed

cell array circuits is only $42.27^{\circ}C$, which is a drop of $6.65^{\circ}C$ and is closer to the $39.63^{\circ}C$ junction temperature of the ideal case without TSVs.

2.5.3 Impact of die thickness

The silicon substrate itself serves a useful role as a heat-spreader, further reducing the impact of localized hotspots. Thus, as the die thickness scales down, it becomes very difficult to spread the heat from the hotspot due to increased lateral thermal resistance. In our proposed system, due to the air gap, the temperature of the stack will be more sensitive to die thickness.

In our test case, we assume there is a $2\text{ mm} \times 2\text{ mm}$ 135 W/cm^2 hotspot in the center

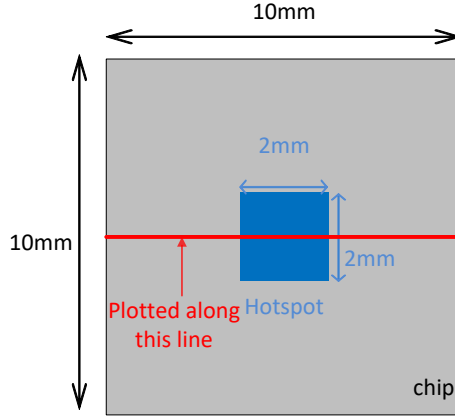


Figure 24: The power map of the processor die. The hotspot (blue square) has power density of $135 W/cm^2$ and the background (grey area) power density is $35 W/cm^2$

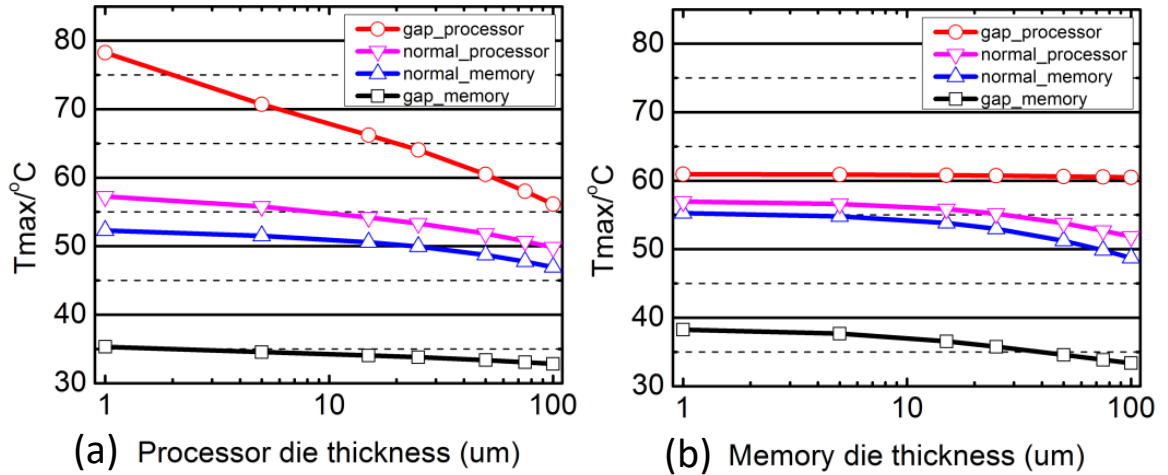


Figure 25: The impact of die thickness. (a) change the thickness of processor. (b) change the thickness of memory.

of the processor die and the background power density is $35 W/cm^2$. The memory die has uniform power density of $1 W/cm^2$. The power map of the processor die is shown in Fig. 24. We separately sweep the die thickness of the memory and the processor from $1 \mu m$ and $100 \mu m$ while fixing the other die at the default thickness. We also compare the results to the normal bonding case (using underfill).

The impact of the processor die thickness is shown in Fig. 25(a). Several observations can be made: firstly, thinning the processor die will increase the temperature of both dice, especially when there is air-gap isolation. With air-gap thermal isolation, the maximum temperature of a $100 \mu m$ thick processor die is only $56 ^\circ C$, while for an ultrathin $1 \mu m$

thick die the maximum temperature is $78\text{ }^{\circ}\text{C}$. Secondly, if there is an air gap, when the processor die is thinned below $50\text{ }\mu\text{m}$, the processor temperature increases very rapidly. Thirdly, for the normal bonding case, the memory die serves as the heat spreader for both dice. In this case, thinning the processor die has a limited impact on system temperature, as shown by the two intermediate lines in Fig. 25(a).

For the impact of the memory die thickness shown in Fig. 25(b), when the air gap exists and provides the thermal isolation, changing the memory die thickness has only a small impact on the processor. For the conventional bonding case, the processor and the memory dice are strongly coupled, thus both dice have similar trends as those of changing the processor die thickness.

2.6 Experimental demonstration ¹

Guided by the above modeling and analysis, a thermal testbed was designed, as shown in Fig. 26(a), to explore thermal coupling and possible solutions [59, 60]. Mechanically flexible interconnects (MFIs) [52] were designed to be clustered in the middle region to further enhance the thermal isolation. Unlike rigid solder microbumps, MFIs can deform elastically under stress and maintain the electrical connectivity between the memory and the processor tiers. Due to this behavior, MFIs can help eliminate underfill and thus reduce the thermal coupling between tiers. The designed and fabricated testbed consists of a low-power and a high-power tier to emulate the heterogeneous 3-D stack shown in Fig. 16; MFIs are used as interconnects between the two tiers (instead of microbumps).

Fig. 26(b) shows the power map and temperature sensor designs for the low-power tier. The low-power tier dissipates a uniform power of less than 5 W. A spiral heater was formed over a $1\text{ cm} \times 1\text{ cm}$ area. Nine resistance temperature detectors (RTDs) were inserted along the middle of the chip in order to measure the temperature profile along the length of the chip. Since the MFIs are clustered in the middle region, the thermal coupling

¹This work is collaborated with Dr. Yue Zhang [58]

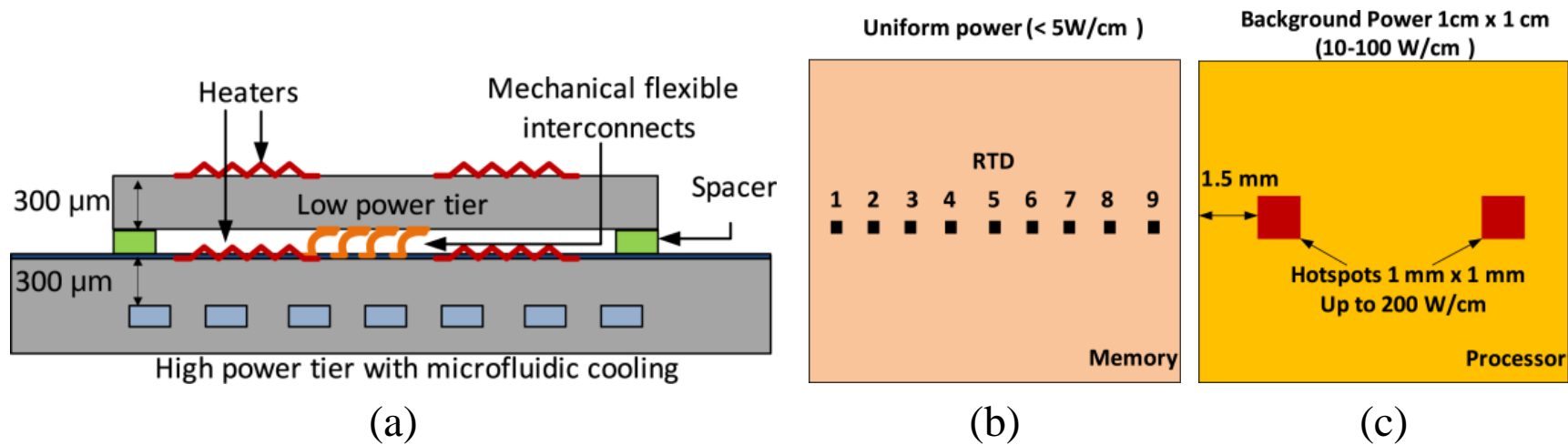


Figure 26: (a) Schematic of the designed testbed for evaluation of the proposed thermal isolation technologies. (b) Top tier (low-power) and (c) bottom tier (high-power) layout design

between the tiers is expected to be non-uniform across the chip; in particular, from the center to the edges of the chip.

Fig. 26(c) shows a schematic illustration of the high-power tier. The chip area is $1\text{ cm} \times 1\text{ cm}$. There are two hotspots on the chip each measuring $1\text{ mm} \times 1\text{ mm}$. The two chips are interconnected with an array of gold-passivated NiW MFIs [52, 61]. The array contains 12×100 MFIs, yielding a total of 1,200 MFIs. This number is chosen based on the Wide I/Os specifications [62]. The MFI design has a pitch of $75\text{ }\mu\text{m} \times 100\text{ }\mu\text{m}$. The total MFI array is $9,940\text{ }\mu\text{m} \times 870\text{ }\mu\text{m}$.

The microfluidic test setup is shown in Fig. 27. Given the micropin-fins are only etched in the high-power tier (bottom tier), the coolant only flows within the high-power tier. The top tier is bonded to the bottom tier through MFIs that are located in the center region. An Agilent DC power analyzer was used to source current into the on-chip heater/RTDs on both tiers. The data logger was used to measure the resistance of the RTDs and extract the junction temperatures [12].

The power maps for the simulated cases are illustrated in Figs. 28 (a) and (b). In Fig. 28 (a), the bottom tier dissipates 10 W/cm^2 across the chip. The junction temperature for each location on both tiers is plotted in Fig. 28 (c) (Case A). Next (Case B), the power density of the two hotspots is increased to 150 W/cm^2 while the background power remains unchanged (Fig. 28 (b)). The corresponding temperature of each chip is plotted in Fig. 28 (c) (Case B). In Case B, the temperature is relatively flat indicating uniform temperature without hotspots. When the power density of the hotspots increases, one obvious observation is that there are two temperature peaks that occur in the bottom die. This is expected because of the large power density of the hotspots. The two temperature peaks are $31.4\text{ }^\circ\text{C}$ and $33.0\text{ }^\circ\text{C}$, respectively. However, also in Case B, there are no obvious hotspots in the upper tier. The temperature of the upper tier gradually increases from $21.1\text{ }^\circ\text{C}$ to $23.1\text{ }^\circ\text{C}$. This demonstrates that the proposed thermal isolation concept effectively minimizes the hotspot coupling between the stacked tiers. In addition, the above described testbed was

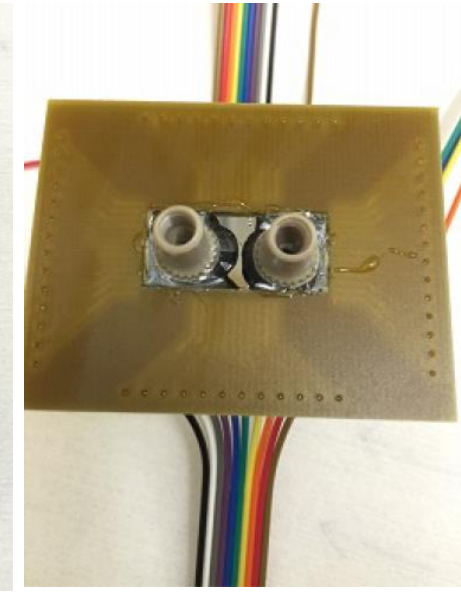
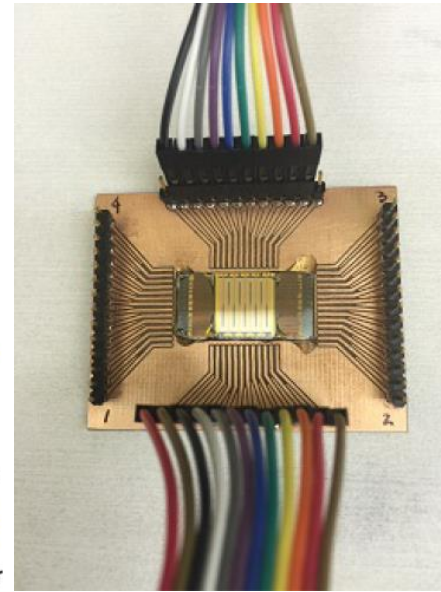
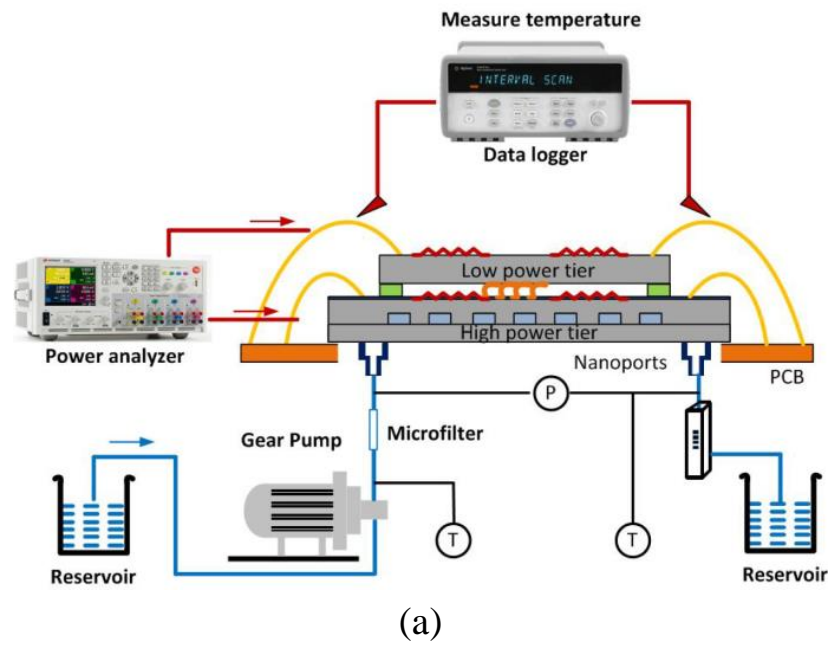


Figure 27: (a) Microfluidic test setup to evaluate the thermal isolation technologies. (b) Top and (c) bottom view of the stack assembled to a PCB board using wire bonding

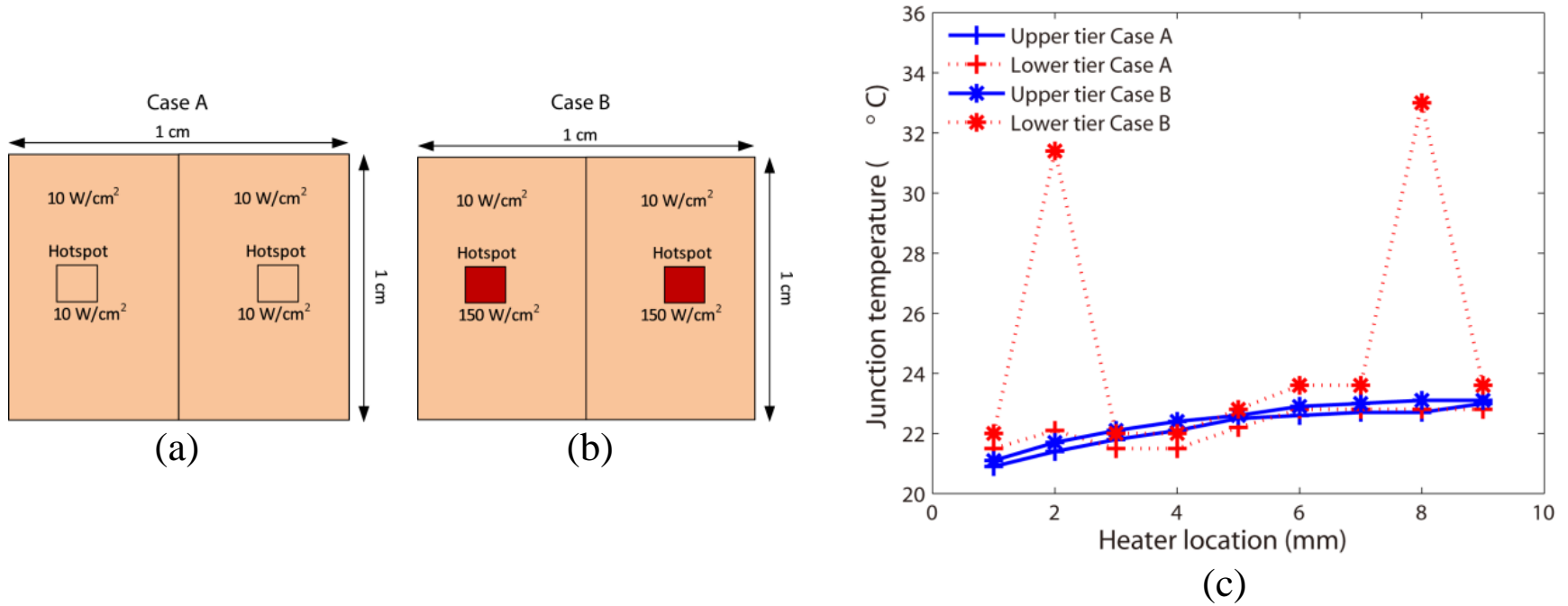


Figure 28: (a) Uniform power density of 10 W/cm^2 in the bottom tier (Case A), (b) background power of 10 W/cm^2 plus two hotspots each dissipating 150 W/cm^2 (Case B), and (c) Junction temperature fluctuation of top and bottom tiers in Case A and Case B

evaluated using the thermal model presented earlier as cross-validation, and the temperature difference is found to be below 7% [59].

2.7 Conclusion

In this chapter, we present the thermal modeling framework and simulation flow. The thermal framework is then validated against *ANSYS* for steady state, transient state and fluidics modeling with a maximum relative error less than 7%, 1%, and 9%, respectively.

A novel stacking structure is proposed with microfluidic cooling embedded in the interposer, thermal isolation between the memory and processor dice and a thermal bridge on top of the memory die. The new architecture exhibits thermal benefits over conventional stacks and is of high value in the heterogeneous integration of high-power and low-power dice. Additionally, we thermally explore our proposed system as a function of micro-bumps, TSVs, die thickness and other system parameters. Specifically, our study benchmarks a memory on processor stack, but the methodologies, analyses and conclusions can be applied to any high-power and low-power stack.

Secondly, tuning of all system parameters is necessary in order to build a thermally-tolerant system: 1) The micro-bumps influence the primary heat path, and must be considered when assessing the thermal performance of the system; 2) The TSVs have an important impact on the thermal isolation layer. With smaller/less/clustered TSVs, the thermal coupling between dice can be minimized; 3) Die thickness also plays an important role. Thinning the processor die below 50 μm will lead to rapid temperature increase of the stack.

Several key conclusions can be drawn from this work. Firstly, our proposed stack can realize lower temperatures for both dice in a memory on processor stack. More importantly, the use of an air gap for thermal isolation causes the memory die to be kept much cooler and to be fairly independent of the fluctuating temperature of the processor because of the thermal isolation.

CHAPTER 3

THERMAL EVALUATION OF 2.5-D INTEGRATION USING BRIDGE-CHIP TECHNOLOGY

In this chapter, 2.5-D integrated circuits using bridge-chip technology are thermally evaluated to investigate thermal challenges and opportunities for such multi-die packages. First, thermal benchmarking of a number of 2.5-D integration approaches is performed and compared to 3-D ICs for completeness. Second, bridge-chip based 2.5-D integrated systems are explored as a function of various technology parameters.

3.1 Introduction

In order to keep pace with rapidly increasing system interconnection requirements [1], multiple advanced interconnect technologies have been proposed, including 2.5-D and three-dimensional (3-D) integration technologies [6, 63, 8]. A key benefit of 2.5-D integration is the ability to assemble heterogeneous dice side-by-side and provide a physical interface comprised of ultra-high density interconnections. Unlike 3-D integration, such 2.5-D integration technologies avoid some of the design, fabrication, and thermal challenges associated with 3-D die-stacking. Nonetheless, there are a number of thermal challenges in 2.5-D integration technologies that we seek to explore and compare to 3-D ICs.

The most widely explored 2.5-D integration technology is based on a silicon interposer approach, as shown in Fig. 29(b). However, while silicon interposer technology has a number of benefits, it may not be a universal solution for 2.5-D integration. To this end, two approaches based on bridge-chip technology have been proposed to explore the formation of 2.5-D systems, as shown in Fig. 29(a) and 29(c).

The first bridge-chip approach uses a silicon chip, called the “bridge” chip, which is embedded into the package. Chip-to-chip interconnects are routed on the bridge-chip and

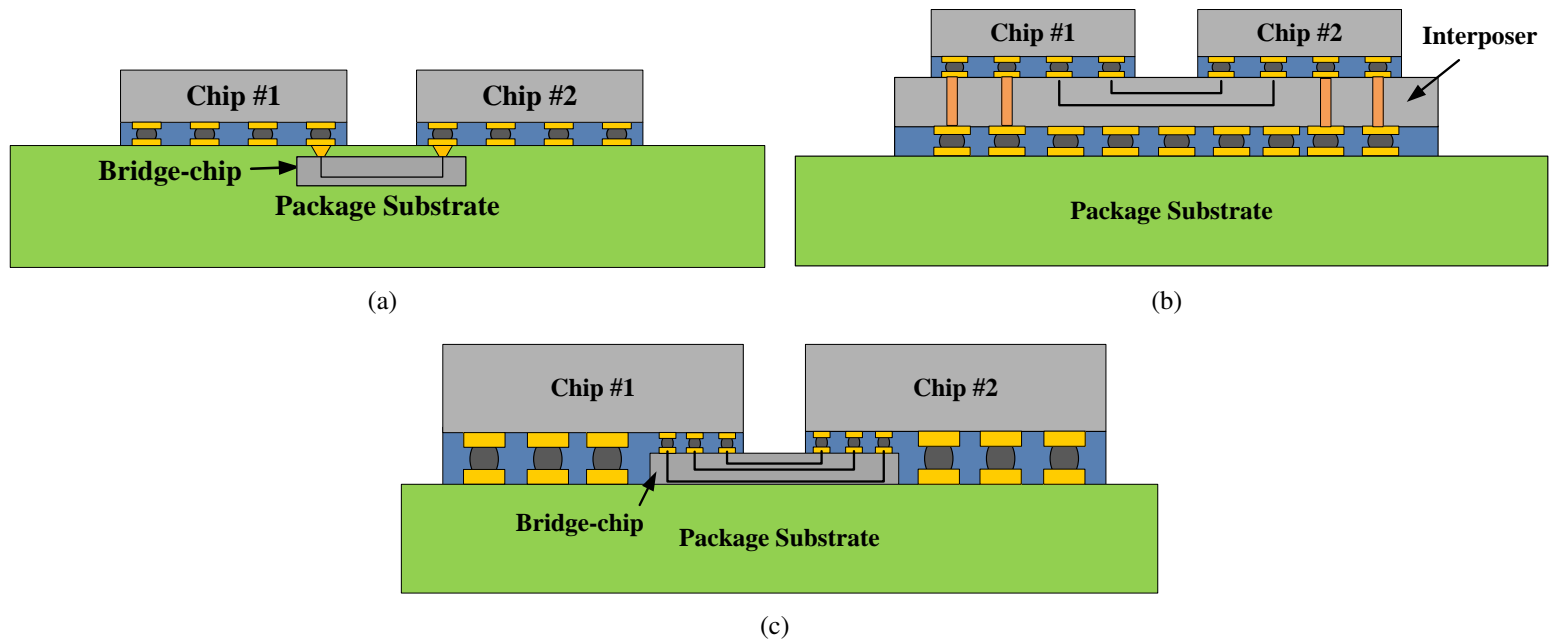


Figure 29: 2.5-D chip stack using (a) bridge-chip technology (b) interposer technology. (c) heterogeneous interconnect stitching technology (HIST).

fine pitch microbumps are used to connect the bridge-chip and the active dice, as shown in Fig. 29(a). The second approach, called heterogeneous interconnect stitching technology (HIST), de-embeds the silicon bridge chip and places it between the active chips and the package, as shown in Fig. 29(c) [10]. Some of the advantages of HIST are as follows: first, the bridge-chip provides similar signal bandwidth density as the silicon interposer technology; second, HIST does not require the use of through-silicon-vias (TSVs); third, HIST is based on die-to-die face-to-face bonding, and thus the interconnections have low parasitics and high density; and lastly, HIST is compatible with any packaging substrate.

Thermal analysis and optimization has been extensively conducted for interposer-based 2.5-D integration [20, 21], as well as TSV-based [17, 64] and monolithic 3-D integration [18, 19]. However, there are no thermal modeling efforts focused on 2.5-D bridge-chip-based interconnection platforms. Moreover, previous thermal efforts have generally focused on one of the above technologies; there is a need for thermal benchmarking of all these approaches.

Therefore, the objectives of this Chapter are twofold. First, we explore the thermal attributes of bridge-chip based 2.5-D integrated ICs and benchmark with other 2.5-D and 3-D solutions. Second, we conduct a deep-dive look at bridge-chip-based 2.5-D integration and evaluate thermal performance as a function of various technology parameters such as bridge chip thickness, TIM properties, microbump properties, die thickness, and die spacing. These studies will help the community understand the thermal limits and challenges facing bridge-chip-based integration technologies.

3.2 2.5-D and 3-D benchmark architectures

3.2.1 2.5-D integration platforms

With the above off-chip technologies, 2.5-D heterogeneous integration of multi-functional chips can be realized. In this Chapter, we focus on high-performance 2.5-D integration of processor, accelerator (GPU, FPGA, or ASIC), and a memory stack. Specifically, we

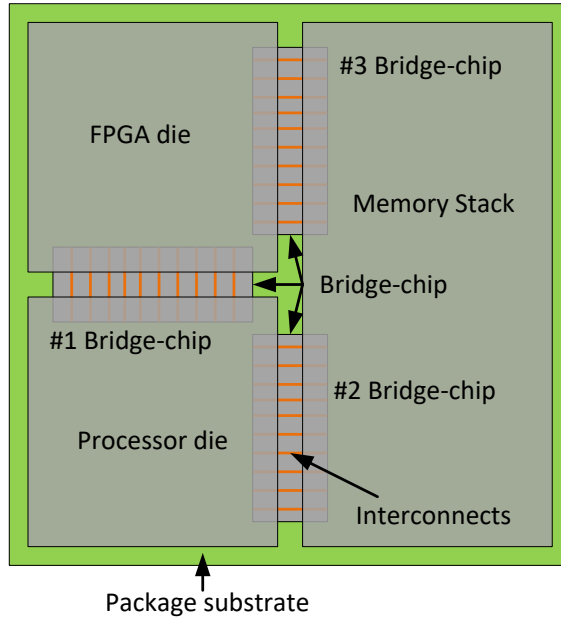


Figure 30: Illustration of the envision FPGA-CPU-Memory 2.5-D chip stack using bridge-chip technology (top view)

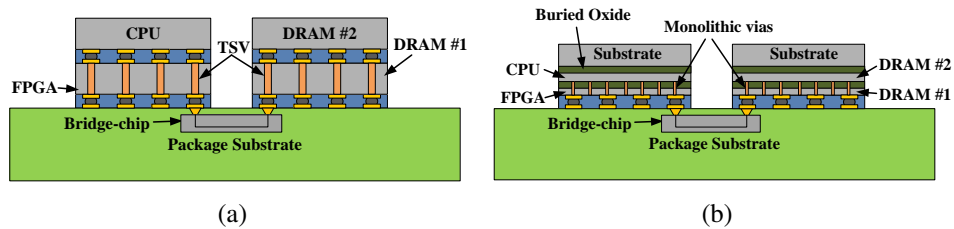


Figure 31: 3-D chip stack (a) TSV-based (b) monolithic nanoscale via based.

envision a CPU-FPGA-Memory 2.5-D microsystem as a test vehicle of the bridge-chip technology and thus, all benchmarks are based on this chip set, as shown in Fig. 30. The FPGA, processor and memory stack are placed side-by-side in a package with bridge-chips underneath the active dice. We assume that the memory stack consists of five tiers, of which the bottom tier is the controller circuit and the other four tiers are memory cells.

3.2.2 3-D integration platforms

Fig. 31 shows two 3-D integration architectures, both of which have been extensively explored in literature. The key enabler for the two stacks is the utilization of vias as chip-

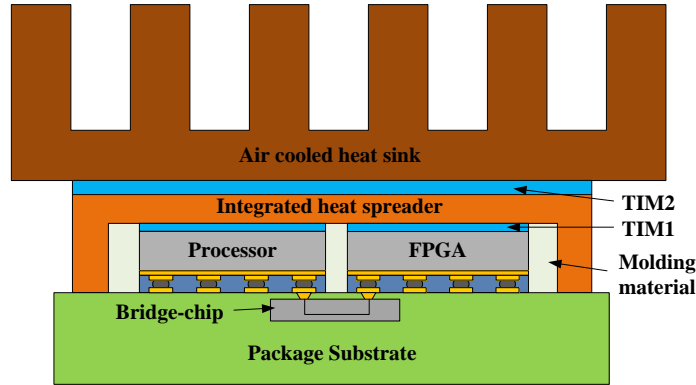


Figure 32: 2.5-D integration using bridge-chip technology with detailed layer information.

to-chip interconnections. The vias in monolithic 3-D are much shorter (\sim a few 100 nm) and smaller ($\sim 100\text{ nm}$) compared to TSVs ($\sim 5\text{ }\mu\text{m}$ diameter and about $40\text{ }\sim 100\text{ }\mu\text{m}$ tall).

In both cases, we still consider a CPU-FPGA-DRAM integrated microsystem. We assume there are two separate 3-D stacks: the first is a CPU-on-FPGA computation stack (CPU is placed on top for thermal consideration), and the second is a DRAM chip stack. To simplify this case study, we assume the FPGA, CPU and memory chips are of the same size. The memory stack has one controller tier and eight memory cell tiers (same storage capacity as 2.5-D cases). The two 3-D chip stacks are placed side-by-side and employ a bridge chip in the package for communication.

3.3 Thermal specifications for 2.5-D and 3-D integration evaluation

Fig. 32 shows a 2.5-D stack with an air-cooled heat sink. The 3-D IC configurations used for evaluation are quite similar except that the chips are stacked. For thermal modeling, we abstract the heat sink and the printed circuit board (PCB) as primary and secondary cooling boundaries, respectively. Both boundaries are modeled using a uniform convection coefficient applied to the top surface of the heat spreader and to the bottom surface of the package substrate, respectively.

3.3.1 Layer thickness and material property

Table 5: Thermal specification

Layer	Conductivity (W/mK)		Thickness (μm)	Heat capacity ($J/^\circ C \cdot Kg$)	Mass Density (Kg/m^3)
	In-plane	Through-plane			
TIM2	3		30	1,000	2,900
Heat spreader	400		1,000	385	8,690
TIM1	3		25	1,000	2,900
CPU/FPGA Die	149		125	705	2,329
Memory Die	149		15	705	2,329
Molding	0.28		N.A.	915	1,790
Metal	61.17	1.62	5	433	7,783
Package	30.4	0.38	1,000	600	1850
Interposer	149		200	705	2,329
Bridge	149		200/25 ¹	705	2,329
Microbump: CPU/FPGA	60		40	227	12,000
Microbump: memory	60		40	227	12,000
Bonding layer: CPU/FPGA	0.9		40	1,000	2,100
Bonding layer: memory	0.9		7.5	1,000	2,100
Copper	400		N.A.	385	8,690
SiO_2	1.38		N.A.	705	2,648
Tungsten	179		N.A.	135	19,250

¹ 200 μm for bridge-chip case and 25 μm height for HIST case.

The layers' information and material properties are summarized in Table 5. On-chip and package metal layers are modeled using in-plane and through-plane thermal conductivity formulated in [38]. Moreover, effective thermal conductivity modeling of the layers with 'vertical interconnects' (microbumps, TSVs, etc.) [38] is implemented to further reduce the mesh number. In the following sections, all the reported simulations are based on these values.

3.3.2 Geometry parameter and boundary conditions

The processor and FPGA dice are assumed to be 1 cm \times 1 cm large. For the 2.5-D cases, the memory die size is assumed to be 1 cm \times 2 cm and for the 3-D cases, the memory die size is assumed to be the same size as the processor and FPGA die, i.e. 1 cm \times 1 cm. As mentioned in Section 3.2, the memory stack has five tiers and nine tiers for the 2.5-D and 3-D cases, respectively. The default spacing between dice is 0.3 mm and the bridge-chip

is set to be $2.3 \text{ mm} \times 7 \text{ mm}$. The package size is $2.23 \text{ cm} \times 2.23 \text{ cm}$. The heat spreader size is assumed to be $4 \text{ cm} \times 3.5 \text{ cm}$.

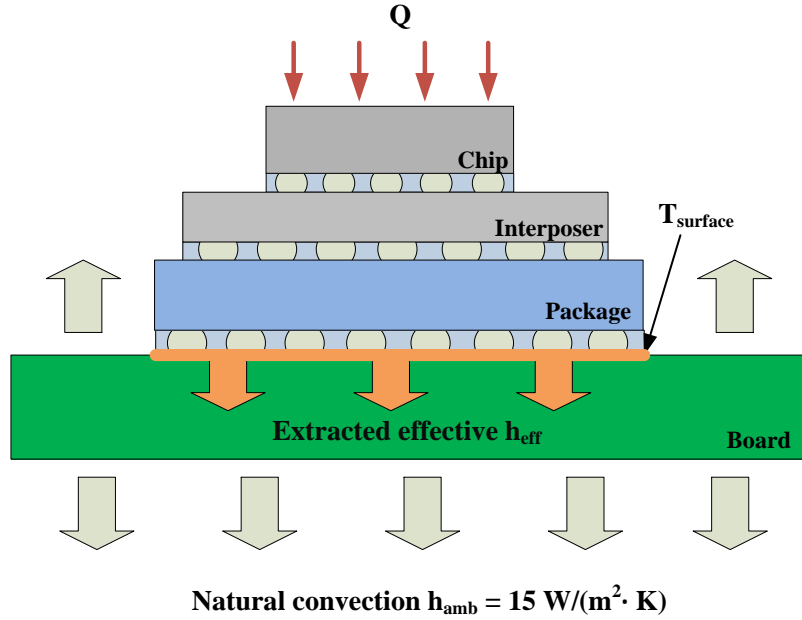


Figure 33: Setup of characterizing effective convection coefficient.

The air-cooled heat sink is assumed to have a case-to-ambient thermal resistance of $0.218 \text{ K}/\text{W}$ [56], and the ambient temperature is assumed to be $38 \text{ }^\circ\text{C}$. For the secondary heat path, we used the method described in [39] to characterize the effective cooling capability of the PCB. We assume the stack is assembled on a $10 \text{ cm} \times 10 \text{ cm}$ PCB and a natural cooling of $15 \text{ W}/^\circ\text{C} \cdot \text{m}^2$ is applied to both surfaces of the PCB, as shown in Fig. 33. In order to extract an effective heat transfer coefficient, so that the whole board does not need to be modeled with the package and dice, a power dissipation of 1 W is applied to the top surface. The effective convection coefficient that the PCB provides can be calculated using the equation shown below:

$$R = \frac{1}{h_{eff} \cdot A} = \frac{T_{surface} - T_{amb}}{Q} \quad (3.1)$$

This value is calculated using weighted average temperature of the bottom surface of the package and is found to be $311 \text{ W}/\text{m}^2 \cdot \text{K}$.

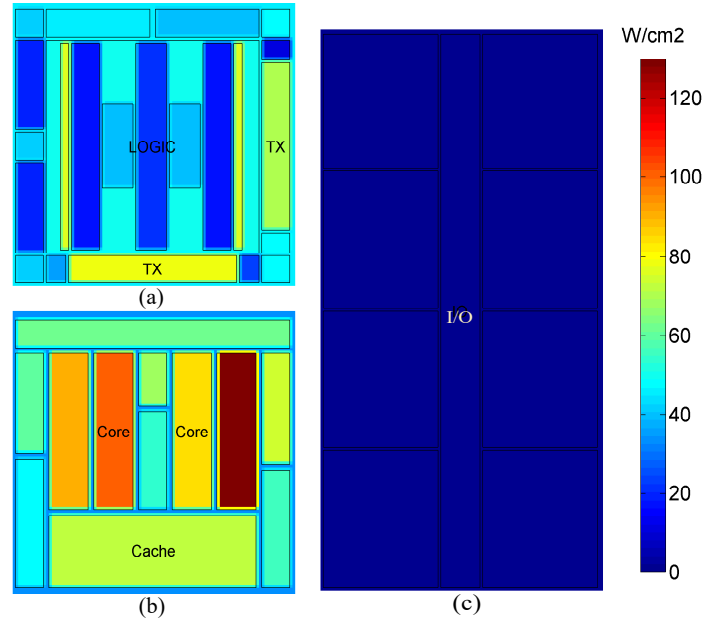


Figure 34: Power density maps of each die (a) FPGA die, 44.8 W (b) Processor die 74.49 W (c) DRAM die (cell circuit), 5.65 W for cell circuit.

3.3.3 Power maps of integrated dice

The layouts of the emulated processor and DRAM dice are shown in [51]. They are based on *Intel Core i7* processor and *Samsung 3-D DRAM*. The processor is assumed to dissipate 74.49 W. For the DRAM chip, the bottom controller is assumed to dissipate a uniform power of 5 W and each DRAM cell tier dissipates 1.46 W. The power profiles of the processor and the DRAM cell tiers are shown in Fig. 34(b) and (c), respectively. The emulated FPGA layout is based on *Altera Stratix V* and *Stratix 10* FPGAs [65]. The FPGA chip power is dependent on application and in our case, we envision the FPGA chip as the server (processor) accelerator. Based on [65], the total power is approximately 44.8 W for the server accelerator and by using the open-source power calculator [66], we can further estimate the power per functional block and emulate the power profile, as shown in Fig. 34(a).

3.3.4 Microbumps and TSVs dimensions

The microbumps we study are categorized into two groups: the first is between the chip and the package or interposer, and the second is between the chip and bridge-chip. For the first group, the diameter and pitch of the microbumps is $40\ \mu m$ and $200\ \mu m$, respectively. For the second group, dense microbumps are used and the diameter and pitch is $20\ \mu m$ and $40\ \mu m$, respectively. Both groups of microbumps are assumed to be uniformly distributed in the corresponding area.

3.4 Comparison of different 2.5-D integrations

All the three 2.5-D IC cases are configured with an air cooled heat sink, as illustrated in Fig. 32. The chip placement and power maps are illustrated in Fig. 30 and Fig. 34, respectively. If not stated, all of the thermal analyses are steady-state results using the maximum power maps shown in Fig. 34 to give a worst case estimation.

Fig. 35 shows the thermal profiles of each die in all cases. The figures are to scale according to die size and spacing listed in Section 3.3. All of the cases exhibit a similar thermal profile because most of the heat is conducted through the attached air cooled heat sink on top (98.18%, 97.17% and 97.19% for the bridge-chip, interposer and HIST cases, respectively). Nevertheless, there is a small difference in the junction temperature resulting from different secondary heat paths in each case. For the interposer-based 2.5-D configurations, heat spreading is enhanced due to the high thermal conductivity of silicon interposer. Therefore, its maximum temperature is the lowest of the three cases (T_{max} is $102.80\ ^\circ C$). Likewise, the bridge-chip is placed closer to the die in the HIST case, and as a consequence, the temperature is lower than the first case ($0.69\ ^\circ C$ cooler). Fig. 36 illustrates the difference in the spreading capability of the bridge-chip and interposer based 2.5-D integration cases. For the interposer 2.5-D case, the thermal profile of the layer beneath the chip (interposer layer) is smoother and exhibits a smaller $T_{max} - T_{min}$ than that

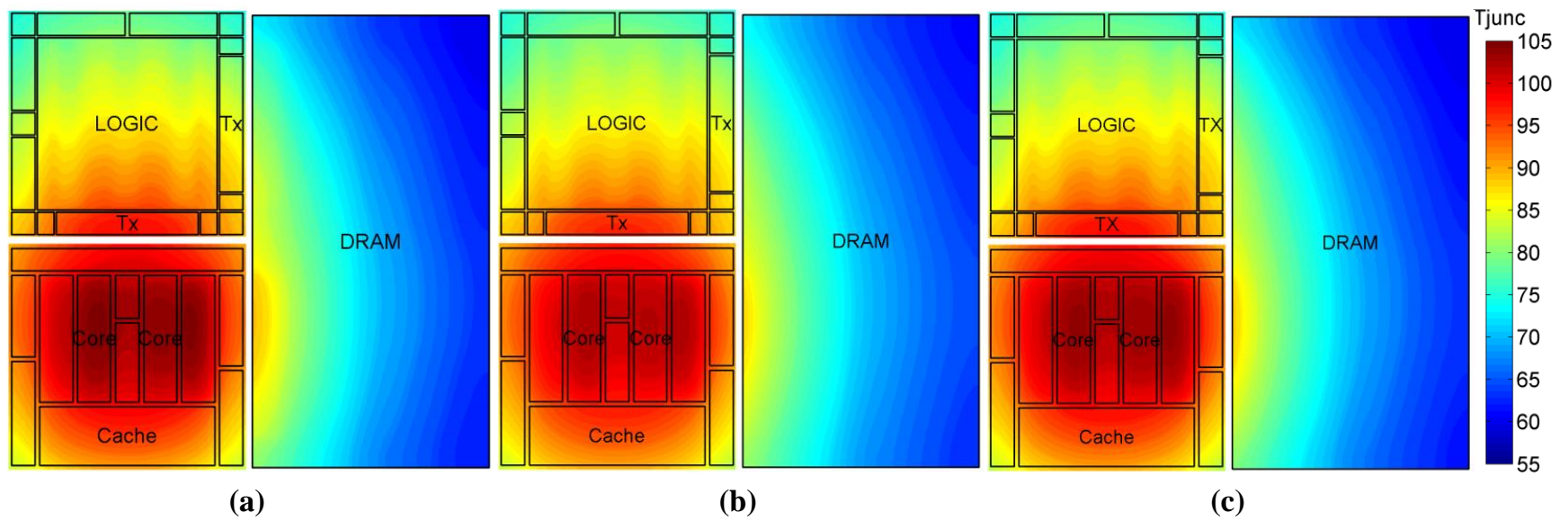


Figure 35: Top view of thermal profiles of each die in all cases. The bottom die, also the hottest die of DRAM chip stack is plotted. (a) embedded bridge-chip, $T_{max} : 104.92 \text{ }^{\circ}\text{C}$ (b) interposer, $T_{max} : 102.80 \text{ }^{\circ}\text{C}$ (c) HIST, $T_{max} : 104.23 \text{ }^{\circ}\text{C}$.

of the bridge-chip case (package layer), which is approximately $4.85\text{ }^{\circ}\text{C}$ lower.

Another observation is the clear lateral thermal coupling between different dice in all cases due to the heat conduction in the heat spreader. The edges of the FPGA and DRAM dice near the processor die are greatly influenced, which creates a relatively larger hotspot area in the FPGA and DRAM dice. To minimize this thermal coupling, it is necessary to apply either independent cooling such as tier-specific microfluidic cooling [14] or thermal isolation technology using an insulator [59] to eliminate thermal coupling.

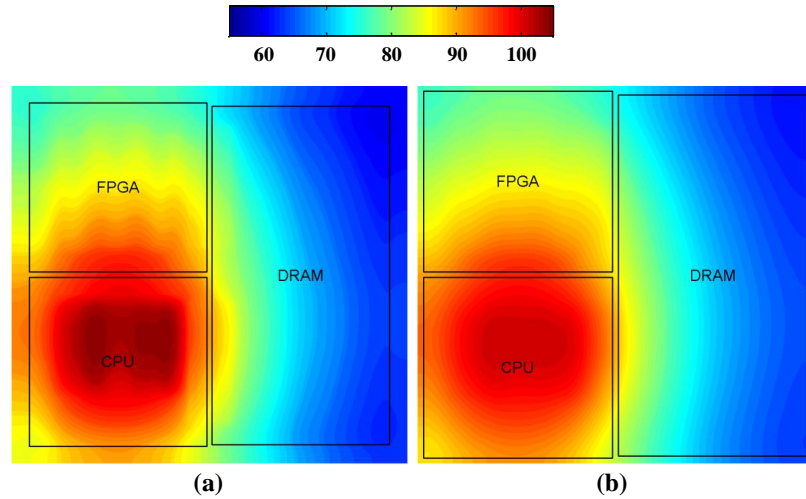


Figure 36: Illustration of heat spreading effects of (a) the package layer in embedded bridge-chip based 2.5-D, $60.22\text{ }^{\circ}\text{C} \sim 104.60\text{ }^{\circ}\text{C}$ (b) the interposer layer in interposer based 2.5-D, $61.61\text{ }^{\circ}\text{C} \sim 101.14\text{ }^{\circ}\text{C}$.

3.4.1 Impact of the thickness of interposer and bridge chip

Based on the above analysis, the thickness of the interposer and the bridge-chip impacts heat spreading. With a thicker silicon layer beneath the dice, the heat spreading is improved and the junction temperature of hottest die is reduced. Therefore, we sweep the thickness of the interposer and bridge chips from 100 nm to $800\text{ }\mu\text{m}$ (for HIST case, the upper bound is $30\text{ }\mu\text{m}$) and plot the maximum junction temperature of each case, as shown in Fig. 37. For the bridge-chip and interposer cases, T_{max} decreases as the thickness becomes larger because of the improved spreading of the thicker interposer and bridge chips. Increasing the

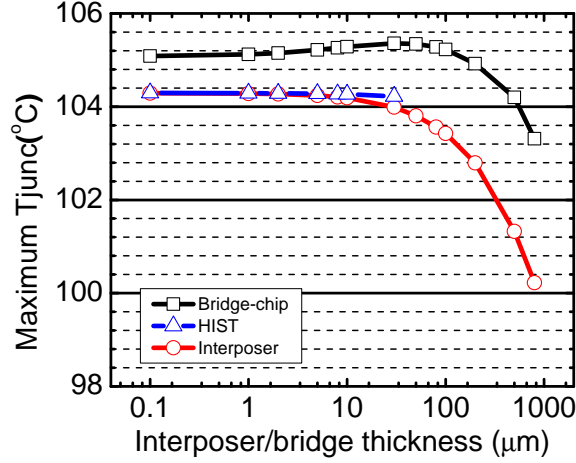


Figure 37: The impact of interposer and bridge thickness.

thickness from 100 nm to 800 μm , T_{max} is reduced by 1.78 $^{\circ}C$ and 4.07 $^{\circ}C$ for the bridge-chip and interposer cases, respectively. The temperature reduction trend is significant when the thickness is larger than 100 μm . For the HIST case, the impact is not significant due to the fact that the thickness of the bridge-chip is limited by the height of the microbumps.

3.5 Thermal comparison between 2.5-D and 3-D integration

3-D integration using TSVs and monolithic nanoscale vias are more advanced integration approaches than 2.5-D integration and have a number of key benefits. However, when multiple chips are stacked vertically, the power density of the resulting 3-D stack will be larger than that of 2.5-D configurations. As a result, the thermal challenges become more difficult to address. In this section, we thermally explore two types of 3-D integration approaches and compare them with bridge-chip based 2.5-D integration.

For the TSV-based 3-D IC case, we assume a die thickness of 125 μm , a TSV diameter of 5 μm , a liner thickness of 0.5 μm , and a pitch of 100 μm . For the monolithic 3-D IC case, the thickness of both the die and buried oxide is 100 nm ; the handle layer is assumed to be 100 μm . The monolithic via diameter is 100 nm and assumed to be tungsten.

Table 6 lists the maximum junction temperature of bridge-chip-based 2.5-D and the two 3-D IC cases. The results show that 2.5-D integration has better thermal attributes than the

Table 6: Thermal comparison of bridge-chip 2.5-D and 3-D integration

Unit: °C	CPU		FPGA		DRAM ¹	
	T_{max}	T_{min}	T_{max}	T_{min}	T_{max}	T_{min}
Bridge-chip	104.92	83.08	98.28	75.02	89.17	60.01
Monolithic	122.29	93.61	124.22	94.25	96.25	63.57
TSV	121.37	94.64	125.62	98.94	98.18	66.57

¹ For DRAM, we show the maximum temperature of the bottom die in the stack (the hottest die).

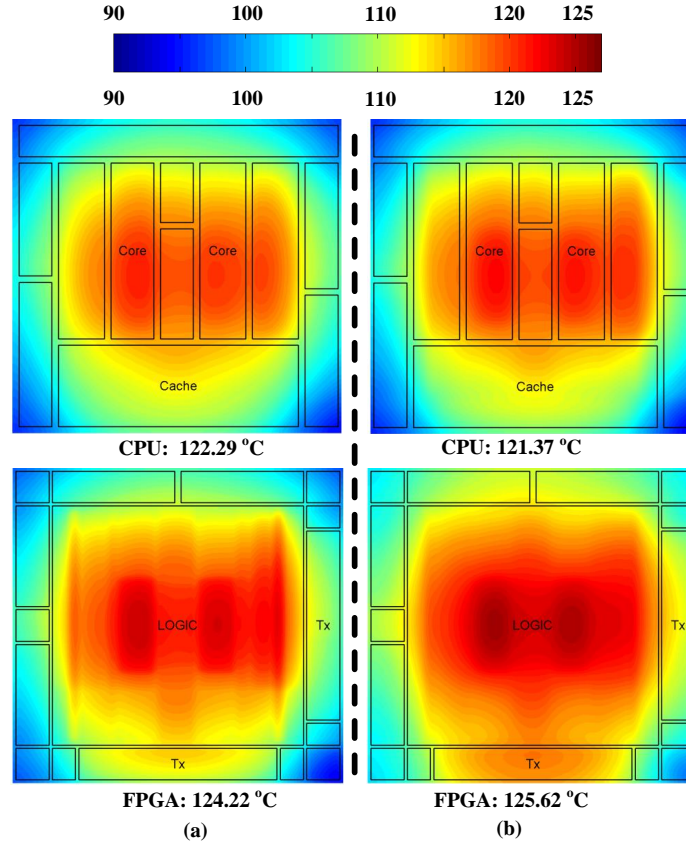


Figure 38: Thermal profile of each die in 3-D stack cases (a) Monolithic-3D (b) TSV-3D

two 3-D IC cases. The maximum junction temperature of the CPU in 2.5-D integration is 17.37 °C and 16.45 °C lower than monolithic and TSV 3-D ICs, respectively. The maximum junction temperature of the FPGA in 2.5-D integration is 25.94 °C and 27.34 °C lower than the monolithic and TSV 3-D ICs, respectively. For the DRAM chips, the maximum junction temperature using 2.5-D integration is 7.08 °C and 9.01 °C lower than TSV and monolithic 3-D ICs, respectively. From a thermal perspective, high power stacks such

as CPU-on-FPGA may not be practical using 3-D technology; while for a low-power stack such as DRAM chips, despite the temperature rise, the maximum temperature does not exceed the thermal limit. Additionally, the DRAM chip can be cooler if less tiers are stacked.

Fig. 38 shows the thermal profiles of the CPU and FPGA dice in the two 3-D IC cases. Compared to the bridge-chip 2.5-D case shown in Fig. 35(a), the 3-D cases have stronger vertical thermal coupling, and the thermal profile of one die exhibits a mirror image of the other. The monolithic case has a smaller active layer thickness thus the spreading is worse than in the 3-D IC TSV case, as shown in Fig. 31. However, the thermal resistance from the FPGA to the heat sink is slightly smaller than the 3-D TSV case, which results in a lower maximum temperature.

3.6 Thermal study of bridge-chip 2.5-D integration

In this Section, we focus on bridge-chip based 2.5-D integration and thermally evaluate as a function of TIM thermal properties, die thickness mismatch, die thickness and die spacing. In addition, transient analysis is performed to further understand die-to-die thermal coupling. Through these analyses, the limits and challenges of bridge-chip based 2.5-D integration are better understood. If not specified, the parameters and power maps are the same as those used in Section 3.3

3.6.1 TIM properties and die thickness mismatch

There are two TIM layers: the first is between the heat spreader and each die (TIM1), and the second is between the heat spreader and the heat sink (TIM2), as shown in Fig. 32. With a good TIM material, the junction temperature will decrease. To evaluate the impact of TIM properties, we sweep the thermal conductivity of TIM1 and TIM2 from $0.9 \text{ W}/^\circ\text{C} \cdot \text{m}$ to $400 \text{ W}/^\circ\text{C} \cdot \text{m}$ (TIM1 and TIM2 are assumed to be the same material). The results are shown in Fig. 39. There is a crossing point in the thermal conductivity at approximately $3 \text{ W}/^\circ\text{C} \cdot \text{m}$, beyond which better TIM material does not yield significant

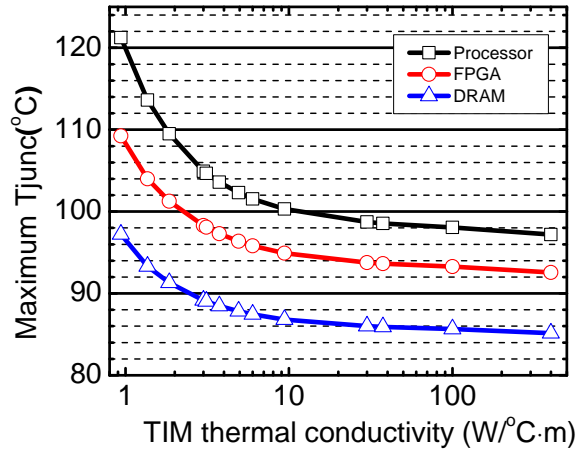


Figure 39: The impact of thermal conductivity of TIM.

benefits. Likewise, changing the TIM thickness leads to similar results.

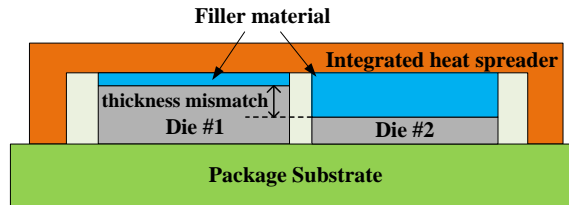


Figure 40: Illustration of die thickness mismatch.

In heterogeneous 2.5-D integration, the chips may be fabricated in different technology nodes and vendors (*Intel 22 nm, 14 nm, TSMC 28 nm, 16 nm*). Therefore, the die thickness may be different and it is necessary to use a material to fill the gap, as illustrated in Fig. 40. The filler candidates are TIM and/or copper (customized heat spreader [67]).

To investigate the impact of die thickness mismatch, we make the following assumptions: first, we assume that the die thickness mismatch is only attributed to the dice; second, we only change the die thickness of one chip and fix the thickness of the other two to the default value; Last, we assume the die with thickness mismatch to be thicker than the other chips.

The results are shown in Fig. 41. There are three observations. First, good fillers are preferred to avoid elevated temperature. If we use copper instead of a TIM (which is not practical), the temperature of each die experiences a nominal change. Second, using the

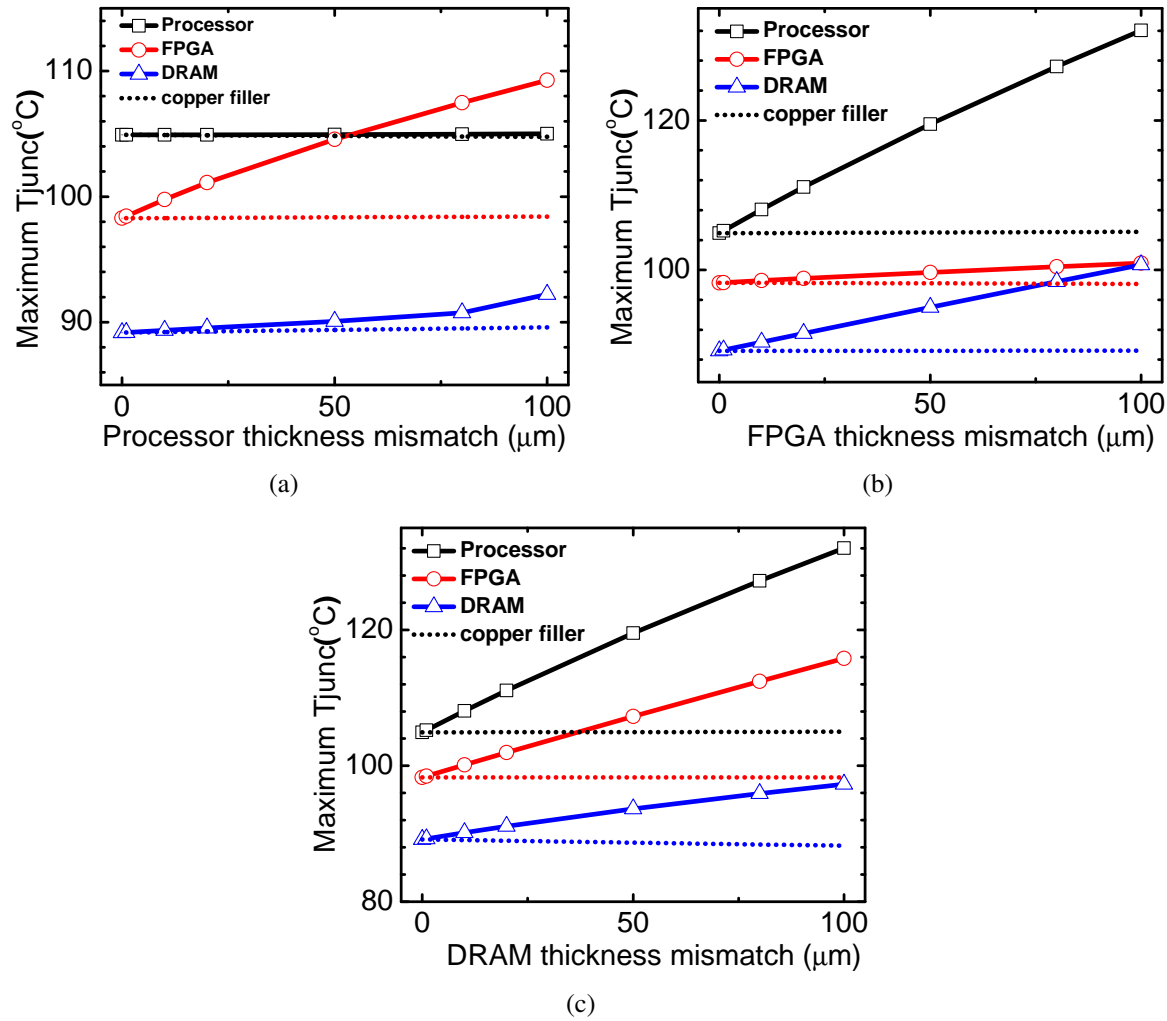


Figure 41: Impact of die thickness mismatch of (a) processor (b) FPGA (c) DRAM. The solid line in the figures is the case using default TIM filler ($3 W/^{\circ}C \cdot m$) and the dashed line is the case for using copper filler ($400 W/^{\circ}C \cdot m$).

TIM filler, the temperature increases as the thickness mismatch increases. Third, when the low-power die is thicker (Fig. 41(b) and Fig. 41(c)), it results in a higher maximum temperature for the whole microsystem. Thus, it is necessary to guarantee that the die with the largest power density has the largest thickness and uses the least filler. In this case, Fig. 41(a) shows when the processor die has a height mismatch of less than $50 \mu m$, the maximum temperature of the whole microsystem does not increase significantly.

3.6.2 Impact of die thickness on heat spreading

The die layer plays an important role in heat spreading, which reduces the localized hotspot temperature. Therefore, as the die thickness scales down, the lateral thermal resistance increases and heat spreading becomes reduced. We sweep the die thickness from $1 \mu m$ to $750 \mu m$ and the results are plotted in Fig. 42. Although the maximum temperature of each die does not change significantly ($2.08^\circ C$, $1.63^\circ C$ and $3.48^\circ C$ for processor, FPGA and DRAM die, respectively), the intra-die variation experiences a relatively larger change. For example, $T_{max} - T_{min}$ for the processor changes from $25.57^\circ C$ to $15.70^\circ C$ when die thickness is changed from $1 \mu m$ to $750 \mu m$.

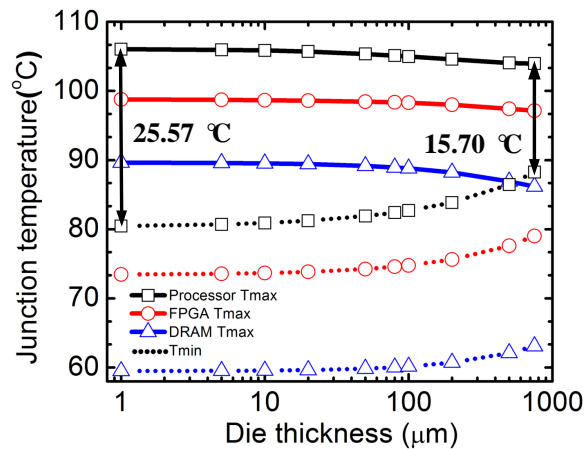


Figure 42: The impact of die thickness scaling. The dotted line plots the T_{min} of each die.

Fig. 43 shows the thermal profile of each die when the die thickness is $1 \mu m$. The block layout is demarcated in the thermal profile as a result of poor heat spreading.

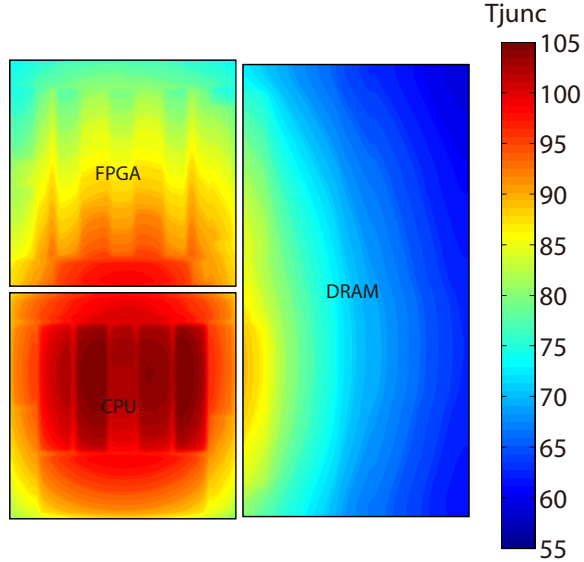


Figure 43: The thermal profile of each die (die thickness is $1 \mu m$). The heat spreading is confined, and the block outlines are clearly observed from the thermal maps.

3.6.3 Impact of microbump and underfill on secondary heat path

The thermal properties of the microbump and underfill impact the secondary heat path. When the effective thermal resistance is reduced, T_{max} of the whole assembly will minimally decrease. On the other hand, due to the fact that most of the heat is conducted through the main heat path, even if the effective thermal resistance of the microbump layer becomes poor, it is expected that T_{max} would not change significantly.

However, if the microbump and underfill become thermally resistive, they form an insulation between the active device layer and the interposer or bridge layer. As a consequence, the interconnection wires will have a lower temperature. Based on the modeling of the thermal impact on the interconnection metric of bandwidth (BWD) over energy per bit (EPB), the BWD/EPB metric can be improved by approximately 7.76% if the temperature is reduced by $30 \text{ }^\circ\text{C}$ [68]. To investigate the impact of microbump and underfill, we fix the number and diameter of microbumps at the default value and only change the underfill and microbump thermal conductivity to change the effective resistance of the microbump layer.

The effective conductivity of the microbump layer is defined as:

$$k_{eff} = k_{bump} \cdot \frac{N \cdot A_{bump}}{A_{chip}} + k_{underfill} \cdot \left(1 - \frac{N \cdot A_{bump}}{A_{chip}}\right) \quad (3.2)$$

$$A_{bump} = \frac{\pi \cdot D^2}{4}$$

where A_{chip} is the chip area, N and D are the number and diameter of the microbumps.

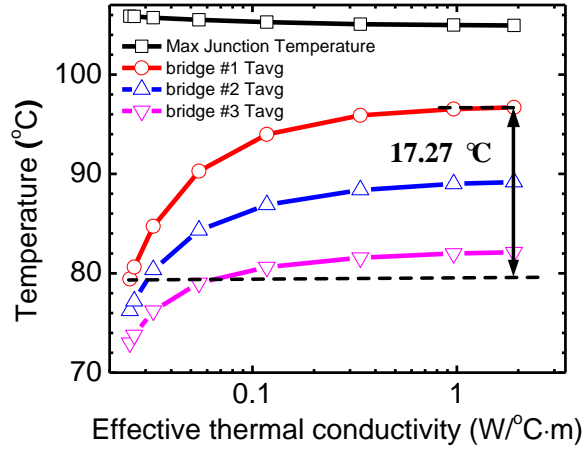


Figure 44: The impact of effective thermal conductivity of microbump layer for bridge-chip case.

Fig. 44 shows the junction and interconnection temperatures as a function of effective thermal conductivity of the microbump layer for the bridge-chip case (similar for interposer and HIST cases). When the effective conductivity of the microbump layer is reduced, there is minimal change in the maximum junction temperature. On the other hand, the average temperature of the bridge-chips (where the chip-to-chip interconnections are located) experiences a temperature change as high as 17.27 °C (bridge-chip #1). This implies that a low conductivity material for the microbump layer helps maintain a lower temperature for chip-to-chip interconnections.

3.6.4 Impact of die spacing on thermal coupling

Fig. 45 shows that as the die spacing increases (heat spreader is kept the same size), the junction temperature of each die decreases. However, the rate of temperature reduction of each die is not the same. For the die with smaller power, the rate is larger and implies

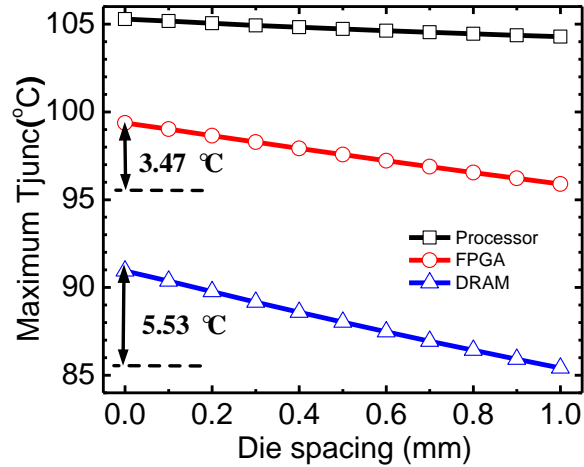


Figure 45: Impact of die spacing. As the die spacing increases, the junction temperature decreases.

that the low-power die is more vulnerable to thermal coupling. The FPGA and DRAM junction temperatures drop by $3.47\text{ }^{\circ}\text{C}$ and $5.53\text{ }^{\circ}\text{C}$, respectively. On the other hand, the processor die temperature has a nominal temperature change when the die spacing increases from 0 mm to 1 mm . Since DRAM performance degrades in the extended temperature range (above $85\text{ }^{\circ}\text{C}$) [45], it is meaningful to take advantage of this effect and carefully design the spacing between the DRAM die and the other high-power chips (of course, there are tradeoffs between thermal considerations, off-chip interconnection metrics such as delay/energy, and integration density as will be discussed in Chapter 6).

3.6.5 Transient thermal coupling

The thermal profiles shown in Fig. 35 represent the final steady-state results but thermal coupling between dice is evolving as the chip activity changes. To investigate the time-domain impact of thermal coupling, we perform a transient analysis to show these time-varying activities. We emulate a processor workload with the activity factor shown in Fig. 46(a). The emulated activity factor has a range of 0.01 to 0.80. To simplify the case, we assume the FPGA and DRAM dice maintain a constant power.

The maximum junction temperature of each die is plotted in Fig. 46(b); time domain thermal coupling can be observed. When the temperature of the processor changes, the

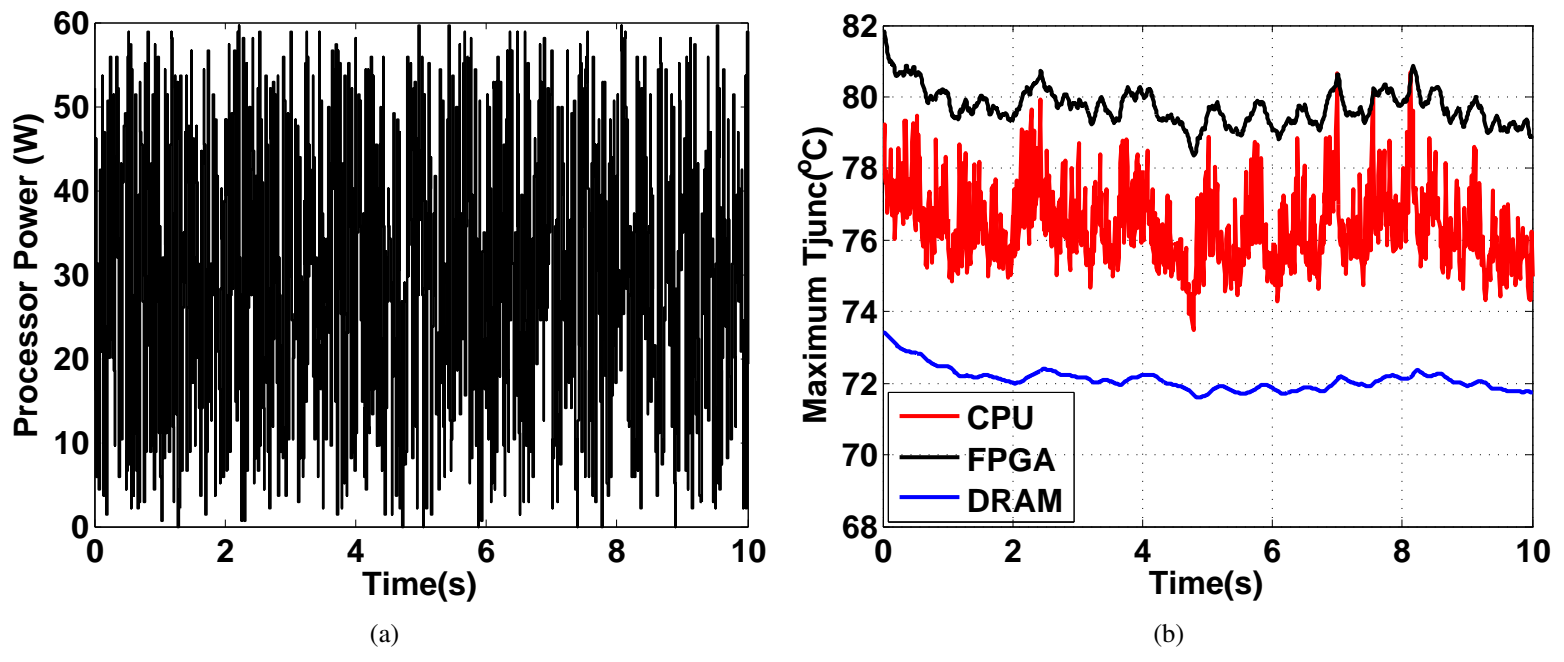


Figure 46: (a) Emulated processor power (b) transient analysis results of bridge-chip 2.5-D integration

other two dice also experience a temperature change, but with a relatively larger response time and smaller variation. The thermal variation of the processor, FPGA and DRAM is $7.23\text{ }^{\circ}\text{C}$, $3.52\text{ }^{\circ}\text{C}$ and $1.81\text{ }^{\circ}\text{C}$, respectively. Due to the lateral distance between the FPGA and DRAM dice to the hotspots on the processor, the two dice respond slower to power changes in the processor.

3.7 Conclusion

This chapter presents a comprehensive thermal study for 2.5-D integration focusing on bridge-chip based technology to identify the thermal limits and challenges in such integration approaches. A CPU-FPGA-DRAM assembly is used as an application example. Bridge-chip 2.5-D integration is compared to interposer and H1ST 2.5-D integration. Compared to bridge-chip 2.5-D integration, interposer 2.5-D integration offers a modest improvements in terms of maximum die junction temperature due to better heat spreading in the interposer layer. Bridge-chip 2.5-D integration is also compared to TSV and monolithic 3-D integration and shows improved thermal response due to smaller power density. We also study the bridge-chip-based 2.5-D integration as a function of bridge chip thickness, microbump properties, TIM thickness, die thickness mismatch, die spacing and die thickness. The simulation results show that the die thickness mismatch should be kept as low as possible and that the hottest die should be the thickest. Moreover, from a thermal perspective, low power dice such as DRAM benefit in maximum temperature by 6.1% when the die spacing is increased from 0 mm to 1 mm. The die thickness plays an important role in heat spreading with thicker die reducing intra-die thermal gradient. Finally, time-domain thermal coupling is investigated. As the temperature of the CPU changes, FPGA and DRAM dice temperatures follow this change.

CHAPTER 4

POWER DELIVERY NETWORK EVALUATION AND BENCHMARKING FOR 2.5-D INTEGRATION USING BRIDGE-CHIP TECHNOLOGY

Power delivery network (PDN) design has been one of the biggest challenges in emerging high-density integration platforms for high-performance computing due to the increased current density and larger parasitics in part from new components such as through silicon vias (TSVs), and reduced voltage levels (which leads to less noise margin). Thus, suppressing power supply noise (PSN) is critical to the success of 2.5-D and 3-D integration platforms. An efficient and accurate PDN modeling framework would help design space exploration and allocate resources more effectively to avoid under- or over- design of PDN.

In this chapter, a PDN modeling framework for emerging integration platforms is presented. The framework is capable of performing both IR-drop and transient analysis. Validation using *IBM* power grid benchmarks shows the IR-drop analysis has a maximum relative error of less than 7.29%, and the transient analysis has a maximum error of less than 0.67% of VDD. The PDN modeling framework is then used to evaluate interposer and bridge-chip based 2.5-D integration platforms. Interposer-based 2.5-D integration may exhibit a worse power supply noise due to the TSV parasitics. In bridge-chip based 2.5-D integration, under the assumption that the bridge-chips underneath the active dice block access to package power/ground planes, there will be some power delivery challenges. In order to mitigate power supply noise (PSN) for bridge-chip 2.5-D integration, several approaches are studied.

4.1 PDN modeling framework

Fig. 47 shows the PDN structure of an IC. Unlike most of the prior work [22, 23, 24] that utilizes a lumped package model, we implement a distributed package-level PDN model to

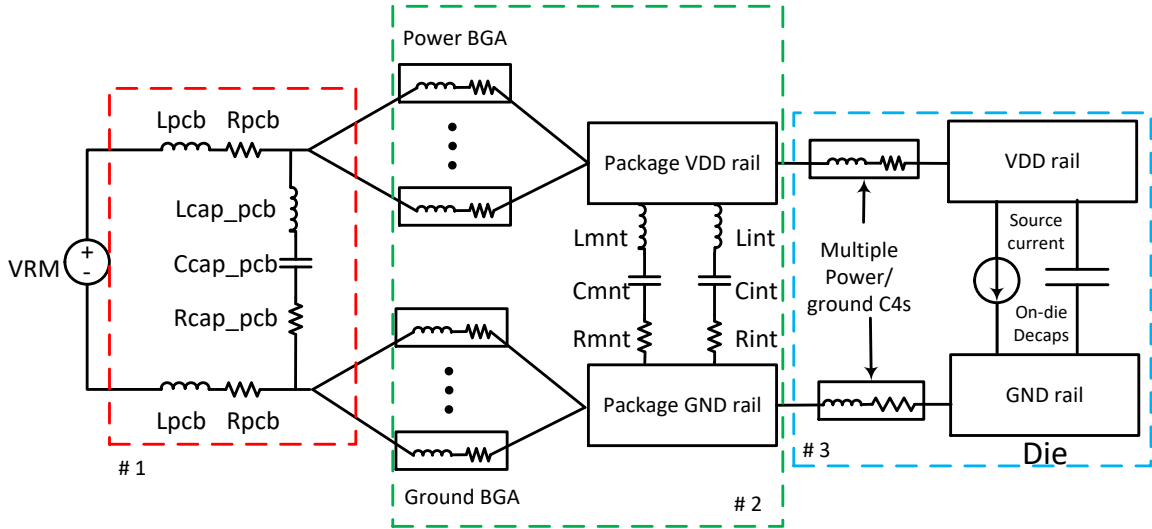


Figure 47: The PDN structure hierarchy. From left to right, a lumped model of board-level PDN, a distributed model of package-level PDN and on-chip PDN are shown, respectively.

reflect the spreading effects of current in the package. This is critical in multi-die packaged systems in which dice may share the package-level PDN.

4.1.1 Board-, Package- and on-die PDN models

In this model, we do not explicitly model the VRM like most of the PDN work [26, 23, 69, 24, 33, 27]; instead, we assume an ideal VRM that is supplying a stable voltage. We use a lumped resistor/inductor network to model the board-level current spreading. Moreover, the equivalent series resistance (ESR) and inductance (ESL) of the board-level decoupling capacitors are included in the model.

Fig. 48 shows the detailed package-level PDN model of power/ground planes. The package power/ground planes are modeled as two layers, where the bottom layer is connected to the motherboard using BGAs, and the top layer is connected to on-die PDN using C4 bumps. Each node in the two layers is connected to six adjacent nodes with a resistor-inductor pair either due to the package traces or inter-layer vias. It is assumed that the surface mounted decaps are only connected to the top layer in the designated areas.

Each R_{sp} and L_{sp} pair in the distributed model represents the current spreading effects.

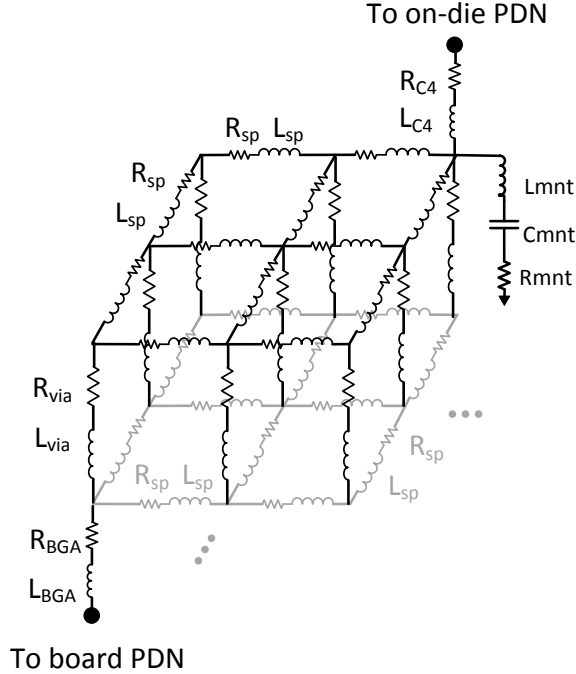


Figure 48: The two-layer package PDN model of power/ground planes

Moreover, each L_{mnt} , C_{mnt} and R_{mnt} pair represents the surface mounted decaps, as shown in Fig. 48. The values of those parameters can be obtained through device characterizations or industry data sheets.

On-die PDN consists of several metal layers, where the power/ground wires are parallel to each other in each layer, but each layer is orthogonal to the layer below/above it (interleaved structure, as shown in the inset of Fig. 49). Prior work has proposed a virtual PDN mesh design using C4 bump granularity with only one metal layer [22, 23, 24]. However, to better reflect the nature of the interleaved PDN design as well as the impact of on-die vias, we model the on-die PDN as a two-layer structure, as shown in Fig. 49. The resistance of R_{top} , R_{bottom} and R_{via} can be extracted from design layout using the process described below.

For each layer of on-die PDN, the metal wires and vias are usually uniformly distributed. If the actual layout is non-uniform, we can calculate the effective wire pitch and via density and re-organize the PDN layout, as shown in Fig. 50.

Next for each layer, we map the fine-granularity PDN layout to coarse mesh grids which

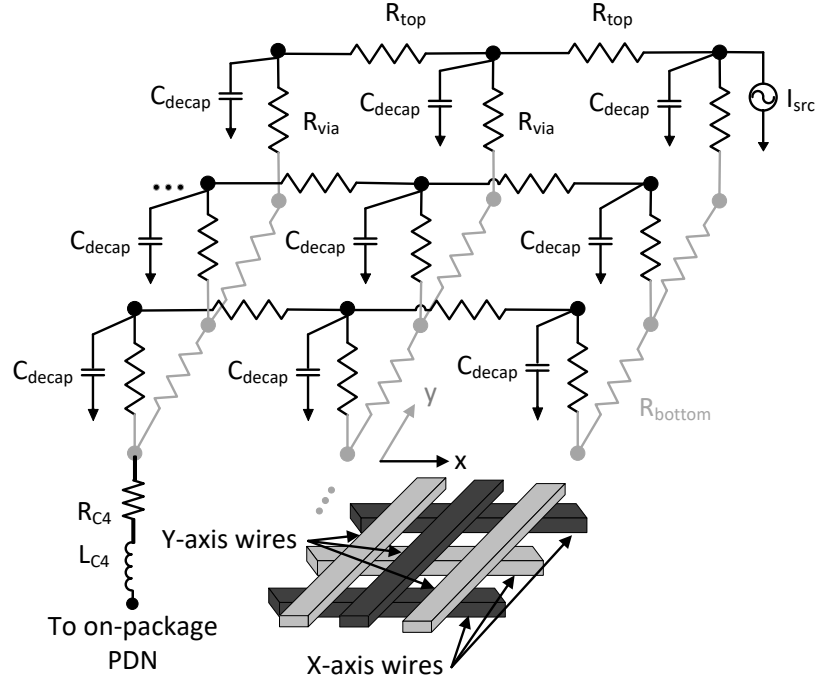


Figure 49: The on-die PDN model. Only one current source and one C4 bump is shown.

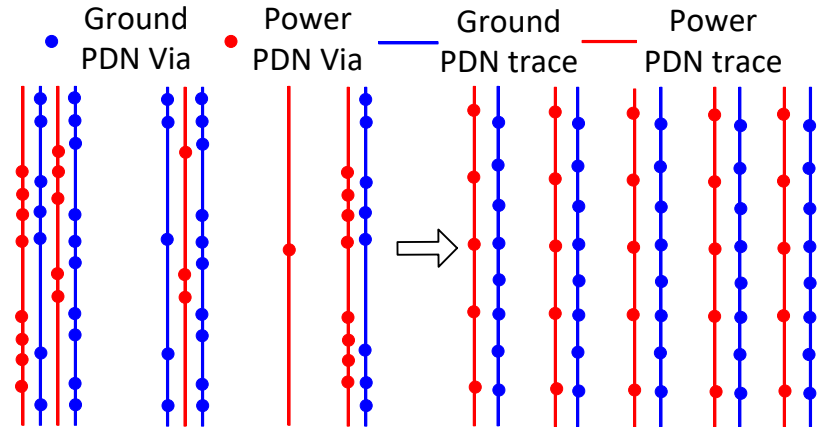


Figure 50: Re-organization of a non-uniform PDN layout

are in C4 bump granularity. Fig. 51(a) and 51(b) illustrates the mapping procedure. For each coarse grid containing multiple vias and metal wires, the equivalent parallel resistance will be calculated and assigned using the equation described in [24], as shown below:

$$R_x = \frac{N_{rows}}{N_{metal.rows}} \cdot \frac{R_{wire}}{N_{cols} - 1} \quad (4.1)$$

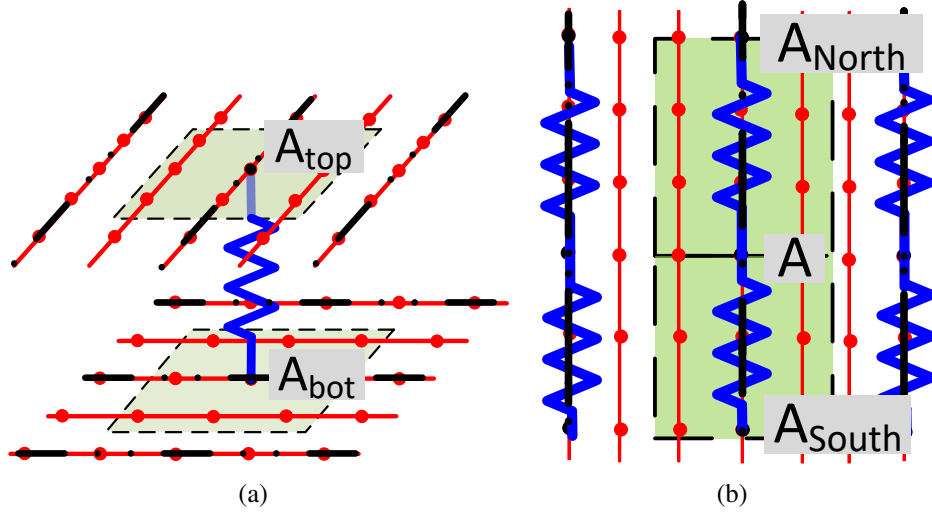


Figure 51: Map fine-grained power PDN layout to coarse meshing grids (a) vias (b) wires.

$$R_z = \frac{A_{chip} \cdot R_{via}}{l_x \cdot l_y} \cdot N_{via} \quad (4.2)$$

where R_x is the resistance of a horizontal branch in the on-die coarse grids; R_{wire} is the resistance of a single wire, N_{rows} , N_{cols} , and N_{metal_rows} are number of rows and columns in the on-die coarse grids, and the total number of metal wires within that layer, respectively. R_y can be calculated using similar equation. R_z is the equivalent via resistance between the neighboring nodes in the adjacent PDN layers, A_{chip} is the total chip area, N_{via} is the total number of vias between the adjacent layers and l_x and l_y is the coarse grid mesh size in horizontal and vertical directions, respectively .

Lastly, all coarse PDN layers with X-axis metal wires are mapped onto the top layer, and with Y-axis metal wires are mapped onto the bottom layer, as shown in Fig. 49. R_{via} in Fig. 49 is the sum of the resistances of vias between adjacent metal layers. Likewise, R_{top} and R_{bottom} are the total parallel resistances between adjacent nodes in all layers with X-axis and Y-axis wires, respectively.

4.1.2 PDN model formulation and simulation flow

The supply voltage noise formulation is shown as follows:

$$\begin{bmatrix} G & A_L \\ -A_L & R \end{bmatrix} \cdot \begin{bmatrix} V(t) \\ I(t) \end{bmatrix} + \begin{bmatrix} C & 0 \\ 0 & L \end{bmatrix} \cdot \begin{bmatrix} \dot{V}(t) \\ \dot{I}(t) \end{bmatrix} = \begin{bmatrix} i_s(t) \\ 0 \end{bmatrix} \quad (4.3)$$

where G is the PDN grid conductance matrix; R and A_L represent the coefficients of nodal voltage $V(t)$ and branch current $I(t)$ in Kirchhoff voltage and current equations, respectively. C and L are the matrices reflecting the capacitive and inductive elements, respectively; $i_s(t)$ is the source current.

For steady state analysis, the time-varying terms are omitted and the branch current $I(t)$ can be expressed in the form of $V(t)$. Eq. 4.3 is then derived in the form of $Y \cdot V(t) = i_s(t)$, where matrix Y is positive symmetric definite. Therefore, the above linear equation can be solved using Cholesky factorization method.

For transient analysis, trapezoid difference scheme can be used to formulate Eq. 4.3, as shown below:

$$\left(\frac{K}{\Delta t} + \frac{U}{2}\right) \cdot X^{n+1} = \left(\frac{K}{\Delta t} - \frac{U}{2}\right) \cdot X^n + \frac{I_s^{n+1} + I_s^n}{2} \quad (4.4)$$

where

$$\begin{aligned} U &= \begin{bmatrix} G & A_L \\ -A_L & R \end{bmatrix} & K &= \begin{bmatrix} C & 0 \\ 0 & L \end{bmatrix} \\ X &= \begin{bmatrix} V \\ I \end{bmatrix} & I_s &= \begin{bmatrix} i_s \\ 0 \end{bmatrix} \end{aligned} \quad (4.5)$$

To accelerate the simulations, we fix Δt which would eventually make $\frac{K}{\Delta t} + \frac{U}{2}$ a constant coefficient matrix. Therefore, we can pre-factorize this matrix before transient simulations using LU factorization. In the solving steps, the triangular factors can be used to solve the

linear equation efficiently.

4.1.3 Comparison with existing PDN models

Various PDN models have been developed to explore IR-drop and switching noise. We compare our PDN model with existing work using a variety of capabilities including IR-drop, transient analysis, impedance analysis, modeling of each PDN level, die stacking and package decap modeling. The comparison is summarized in Table. 7.

Based on the research objectives, the models have different focuses. For thermal and IR-drop co-simulation models of [28] and [33], only resistive elements need to be considered, therefore, they implement a detailed distributed package and PCB models. For the work focusing on the on-die PDN [24, 70, 27] and on-chip voltage regulator design [71], their models have more detailed on-chip PDN models but very abstracted package design. However, for the work investigating the impact of decap placement [72], the distributed package model is implemented. For this work, we leverage the benefits of different models, and focus on on-die and on-package PDN design.

4.1.4 Steady-state and transient analysis validation

To validate the PDN framework, the *IBM* power grid benchmarks [73] have been used. The benchmarks are provided in the *HSPICE* netlist format. There are eight benchmarks for steady-state analysis and six benchmarks for transient analysis. For steady-state results, benchmarks provide the overall noise profile including the noise level at each node. On the other hand, for transient results, the benchmarks provide the waveforms of 20 randomly selected nodes throughout the whole circuit. The benchmark size and the number of metal layers are summarized in the first two columns of Table 8.

We use scripts to extract the layout and RLC information from the provided *HSPICE* netlists, and then we map the PDN layout onto the coarse mesh grids at the granularity of C4 bump pitch. Next, we solve for the supply voltage noise of both states using the above

Table 7: Comparison of different PDN modeling work

	IR-drop	Transient analysis	AC analysis	On-die PDN	Package PDN	Board PDN	VRM model	Multi-Die	Research objective
J. Xie [28]	Yes	No	No	Distributed, single-layer, no vias	Distributed	Distributed	No	2.5-D & 3-D	IR-drop with thermal impact
Y. Shao [33]	Yes	No	No	Distributed, single-layer, no vias	Distributed	Distributed	No	2.5-D & 3-D	IR-drop with thermal impact
C. Pan [72]	Yes	Yes	No	Distributed, single-layer, no vias	Distributed	Lumped	No	No	Impact of decap placement
R. Zhang [24]	Yes	Yes	No	Distributed, multi-layer, no vias	Lumped	Lumped	No	3-D	PDN design and architecture
S. Park [70]	Yes	Yes	Yes	Distributed, single-layer, no vias	Lumped	Lumped	Lumped	No	Co-design with clock tree network
X. Zhang [23]	Yes	Yes	Yes	Lumped	Lumped	Lumped	No	No	Impact of PSN on subthreshold ICs
H. He [27]	Yes	Yes	Yes	Distributed, single-layer, no vias	Lumped	Lumped	Lumped	3-D	PDN challenges for 3-D IC
Z. Zeng [71]	Yes	Yes	Yes	Distributed, multi-layer, with vias	Lumped	Lumped	Circuit modeling	No	On-chip voltage regulator study
This work	Yes	Yes	Yes	Distributed, multi-layer, with vias	Distributed	lumped	No	2.5-D & 3-D	Design exploration for 2.5-D and 3-D integration

mentioned framework. We compare three sets of metrics: current of each C4 bump, IR drop of each node and transient noise of all the selected nodes.

Table 8: Validation Results

Circuits (# of Nodes)	Metal Layers	Bump Current Error (%)	Max IR-Drop Error (%)	Transient Error (%VDD)
Ibm1 (31 K)	2	21.75	20.29	1.84
Ibm2 (127 K)	4	7.14	11.11	0.67
Ibm3 (852 K)	5	3.59	2.21	0.54
Ibm4 (954 K)	6	7.60	0.71	0.12
Ibm5 (1.08 M)	3	6.12	3.03	0.22
Ibm6 (1.67 M)	3	7.29	1.23	0.22
Ibm7 (1.46 M)	6	5.34	5.71	N/A
Ibm8 (1.46 M)	6	5.34	5.71	N/A

Steady state results

The steady-state validation results are summarized in the third and fourth columns of Table 8. Except for the small benchmark cases IBM1 and IBM2, which have highly non-uniform PDN structure, all cases obtain a maximum relative error of less than 7.60% and 5.71% in bump current and IR-drop, respectively. The noise profiles are also compared and the results are well matched. Fig. 52 shows an example of the noise profile comparison of IBM3 as this benchmark has the largest noise gradient. The model accurately captures the distribution of the noise. Fig. 53 shows the bump current comparison of IBM3, in which we sort the current of each bump in an ascending order, and plot both *IBM* provided and our modeling results [24]. Likewise, although the current value spans a wide scale (approximately 5X), the bump current is very well matched.

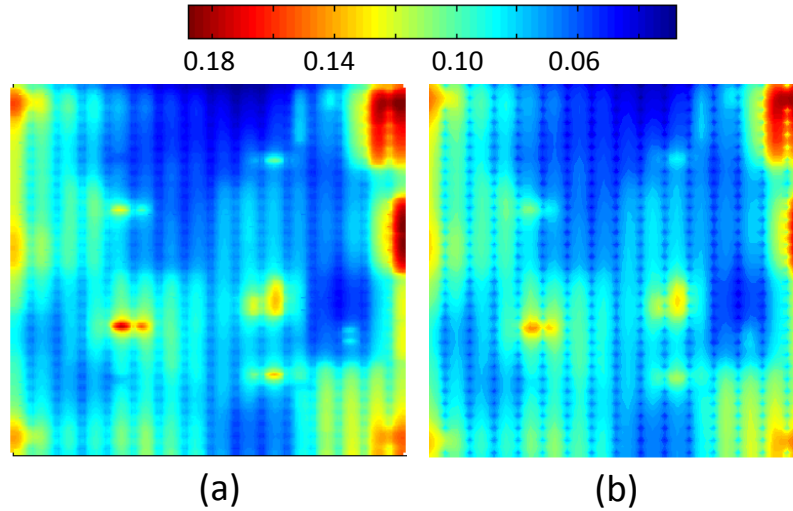


Figure 52: The noise profile of IBM3 (a) Provided results by *IBM* PG benchmarks (b) modeling results.

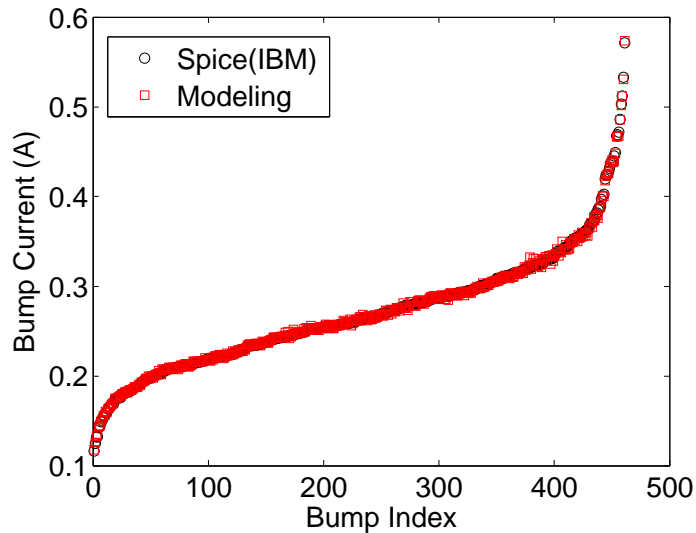


Figure 53: Bump current comparison of IBM3.

Transient state results

Transient validation results are summarized in the last column of Table 8. We normalize the error to supply voltage because some of the benchmark noise values at some time point are small and thus, the relative error can be high. For transient state analysis, except for IBM1, the maximum error for all the cases is less than 0.67% VDD. Due to the fact that the package inductance contributes most of the switching noise, transient state results are

relatively more accurate. Fig. 54 shows the node in IBM2 with the maximum error. Even for this case, the peak noise and waveform are well captured.

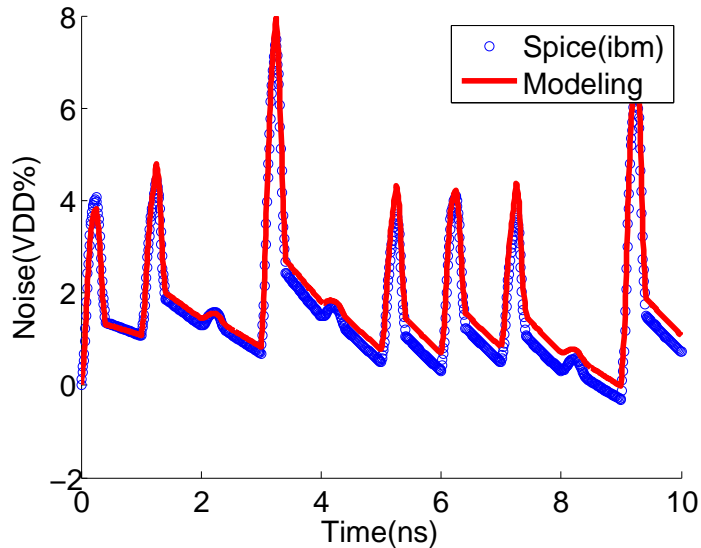


Figure 54: The transient noise the node in IBM2 with maximum error.

4.2 PDN challenges of 2.5-D integration

The trends of lower supply voltage, higher current demand and increased power density are making power delivery in high-performance digital systems an increasingly difficult challenge [22]. Due to the resonances generated from the interactions of the board-, package-, and die-level parasitics, it is difficult to ensure power integrity over a wide frequency range. Moreover, to address the bandwidth and performance limitation of conventional single-die package [74], there is an increasing interest in placing multiple dice into a single package using three-dimensional (3-D) and 2.5-dimensional (2.5-D) integration technologies [9, 10, 17], which exacerbates the power delivery challenges.

Power delivery network (PDN) and power supply noise (PSN) in traditional single-chip [23, 24, 25] and 3-D [26, 27, 28] have been extensively studied in the literature. However, 2.5-D integrated electronics have not been investigated as thoroughly. Specifically, 2.5-D integrated electronics have several unique attributes that require consideration. For exam-

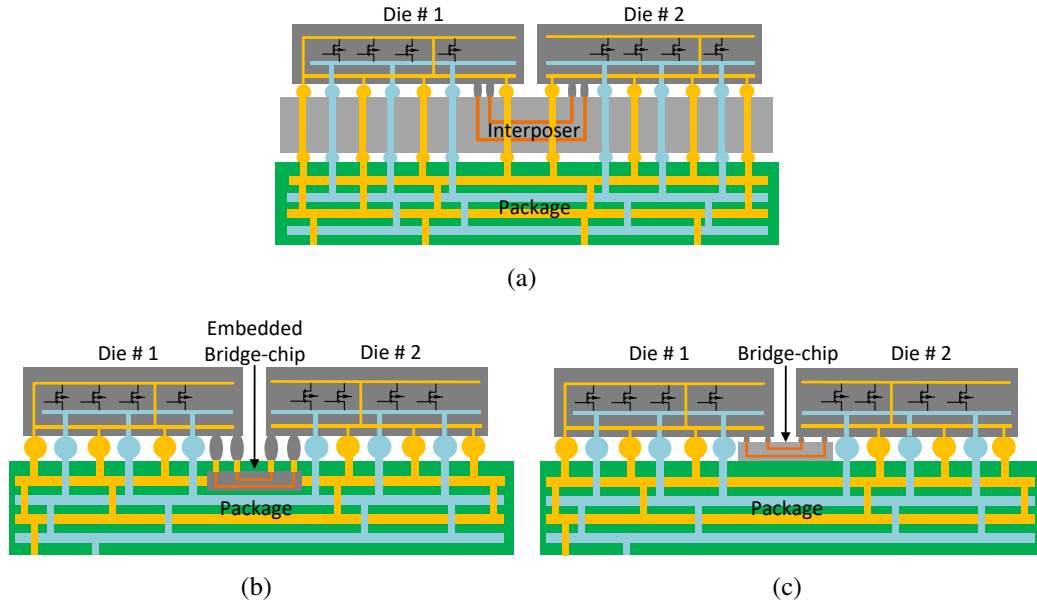


Figure 55: Three different 2.5-D integration platforms (a) interposer (b) bridge-chip (c) HIST

ple, in a silicon interposer based 2.5-D integration as shown in Fig. 55(a) [6, 69], the use of power/ground TSVs increases PDN parasitics. Likewise, for embedded multi-die interconnect bridge (EMIB) as shown in Fig. 55(b) [9] and heterogeneous interconnect stitching technology (HIST) as shown in Fig. 55(c)) [10], signal interconnections and driver circuits are placed, generally, on the edges of the dice and above the bridge-chips, which may lead to a reduction in the power/ground C4 bumps that have access to the package-level power/ground planes. This reduction can lead to increased PSN. Prior work has studied the PDN of interposer-based 2.5-D integration [75, 76], however, there are no PDN modeling efforts focused on bridge-chip based 2.5-D integration. Moreover, there is a need for PDN benchmarking of all these 2.5-D approaches and comparing them with conventional single-die package.

Therefore, in this Chapter, the PDN of 2.5-D integrated electronics (interposer and bridge-chip based interconnections) is evaluated to explore the challenges and opportunities of 2.5-D integrated electronics. Moreover, we explored the impact of several technology parameters such as metal layers, TSVs and bridge chip size to explore the design space of

2.5-D integration. Last, to mitigate PSN in the bridge-chip based assemblies, we propose to insert through bridge-chip vias to mitigate PSN.

4.3 PDN evaluation and benchmarking of 2.5-D integration

4.3.1 Study cases

Fig. 55 shows three 2.5-D integration technologies with different approaches for chip-to-chip interconnection. The first approach utilizes a silicon interposer technology. The second approach is an interconnect-bridge technology described in [9] and utilizes embedded silicon chips within the package to route the chip-to-chip interconnects. The last approach is based on placing a ‘stitch’ chip above the package substrate between the active dice [10].

All three 2.5-D integration technologies impact the PDN and are compared to a single-die package. First, the PDN of interposer-based 2.5-D integration contains TSVs resulting in larger parasitics. Second, for the interconnect-bridge and HIST approaches, since it is unlikely to have vias penetrating the bridge-chip, it is difficult to have power/ground C4 bumps directly interconnected to the package-level power/ground planes in regions of the active dice that overlap with the bridge-chips. As a consequence, the PSN in those regions will be impacted. In this study, we will make the above worst-case assumption for the interconnect-bridge and HIST approaches.

When the silicon die containing the bridge interconnects is embedded within the package, the package-level wiring and routing through these respective regions will be impacted. Despite this effect, our simulation results do not show significant differences between the approaches shown in Fig. 55(b) and 55(c), and thus, we will refer to their results as ‘bridge-chip’ for the rest of the Chapter. Moreover, as a reference to the best achievable results for 2.5-D integration, we also perform simulations for each die in a single-die package, which we will refer to as ‘single-die’ case. In summary, we will present the results of interposer, bridge-chip, and single-die cases in our study.

4.3.2 PDN design parameters and specification

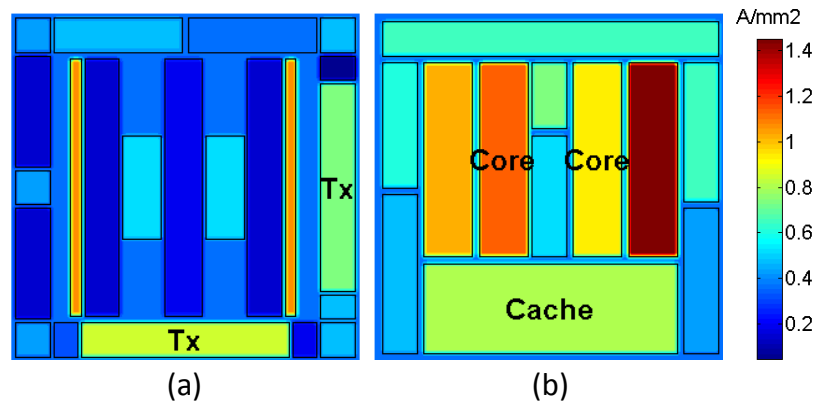


Figure 56: The current density of each die. (a) die #1 (b) die #2

Our framework can model any heterogeneously integrated system, such as processor-memory and processor-accelerator. In this study, we emulate a processor-FPGA 2.5-D integrated package. Die #1 emulates a FPGA die and is assumed to have a total current of 49.78 A [1, 66]. The emulated FPGA layout is based on *Altera Stratix 10* FPGAs [65]. Die #2 emulates a processor with a total current of 82.77 A [51]. The current density maps are shown in Fig. 56. The supply voltage is assumed to be 0.9 V.

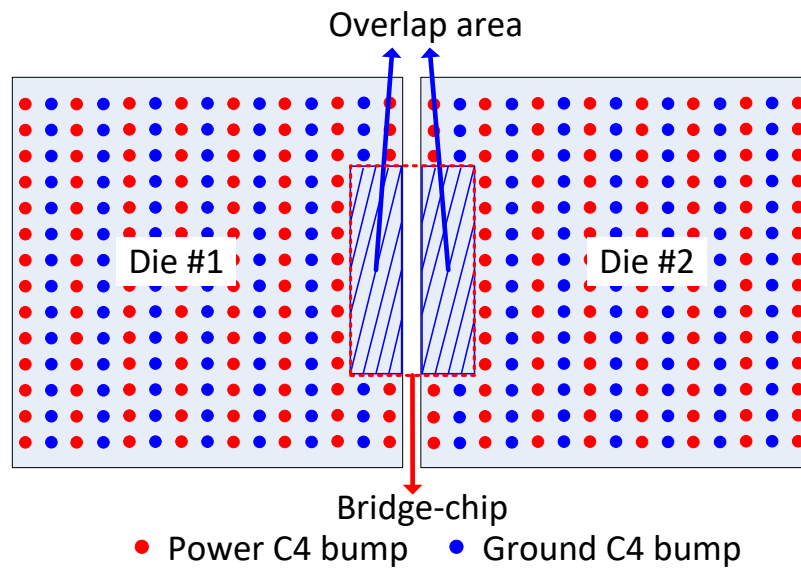


Figure 57: Illustration of bridge chip placement: an example with a single bridge chip.

Both dice are assumed to be $1\text{ cm} \times 1\text{ cm}$, and the package is assumed to be 2.45 cm

Table 9: Parameters for PDN model

Parameter	value
On-die metal resistivity ($\Omega \cdot m$)	1.8e-8
On-die global wire Pitch/Width/Thickness (μm)	39.5/17.5/7
On-die intermediate wire P/W/T (nm)	560/280/506
On-die local wire P/W/T (nm)	160/80/144
on-die decap density (nF/mm^2)	335
C4 bump pitch/R/L ($\mu m/m\Omega/pH$)	200/14.3/11.0
Package effective decap R/L/C ($(m\Omega/pH/\mu F)$)	541.5/220.7/52
Package resistivity/inductance ($m\Omega/mm/pH/mm$)	1.2/24
BGA pitch/R/L ($\mu m/m\Omega/pH$)	500/38/46
TSV R/L ($m\Omega/pH$)	54.2/77.78
PCB R/L ($\mu\Omega/pH$)	166/21
PCB Decap R/L/C ($(\mu\Omega/nH/\mu F)$)	166/19.54/240

$\times 1.8$ cm. The two dice are placed side-by-side with a die spacing of 0.5 mm. The C4 bumps are assumed to be uniformly distributed with P/G/P pattern, as shown in Fig. 57. The bridge chip has a total area of 1.5 mm \times 6 mm and the overlap area with each die is assumed to be 0.5 mm \times 6 mm (I/O area), as shown in the shaded area of Fig. 57.

Table 9 summarizes the parameters used in the PDN simulations. Since the FPGA and processor dice may have different supply voltages, they are assumed to have separate power delivery domains in each package layer and the PDN area in the package is equally assigned for simplification.

4.3.3 Comparison of different 2.5-D integration

IR Drop

IR-drop analysis results are summarized in Fig. 58. For the interposer case, the utilization of TSVs leads to larger noise; the IR-drop is approximately 30.0% (Die #1) and 27.5% (Die #2) larger than the single-chip case.

For the bridge-chip case, compared to the single-die case, the additional noise is mainly due to the absence of C4 bumps that interconnect to the package power/ground planes in the overlapping regions with the bridge chips. The IR-drop is approximately 54.4% (Die

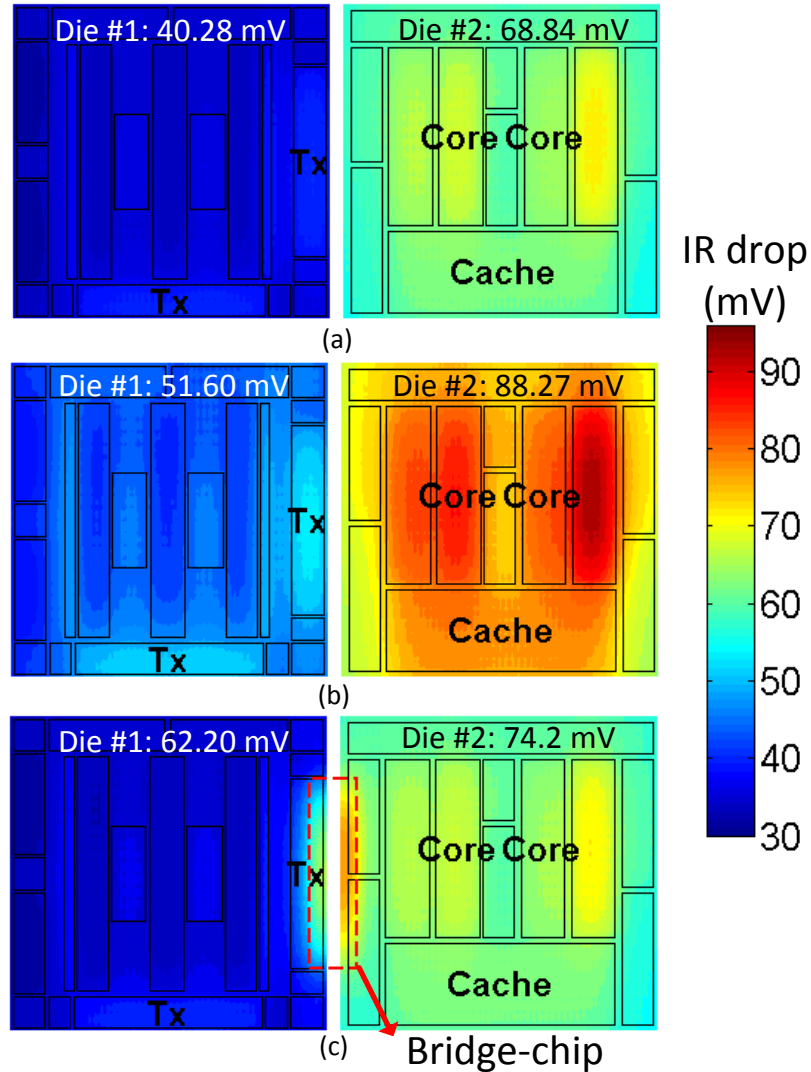


Figure 58: The IR drop profile of each die for (a) Single die (b) interposer (c) single bridge-chip with a overlap area of $0.5 \times 6 \text{ mm}$.

#1) and 7.8% (Die #2) larger than the single-chip case. Therefore, the overlap area between the bridge-chip and the active dice (I/O area) has an important impact on the PSN. Most of the IR drop occurs in the overlap region because of the absence of C4 bumps.

One potential solution to this problem is to break the bridge chip into multiple dice with an aggregated area equating to the original bridge size. With multiple bridge chips, the equivalent power delivery distance from the bumps to the center of the overlap area becomes shorter than the single bridge-chip case and therefore, the IR-drop is expected to be smaller. The bump pattern of a five bridge-chip case is illustrated in Fig. 59.

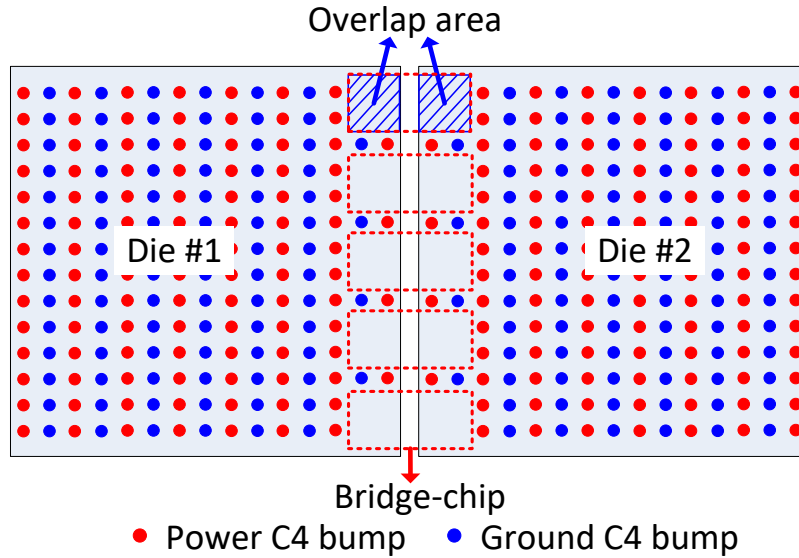


Figure 59: Illustration of bridge chip placement: an example of 5 bridge chips.

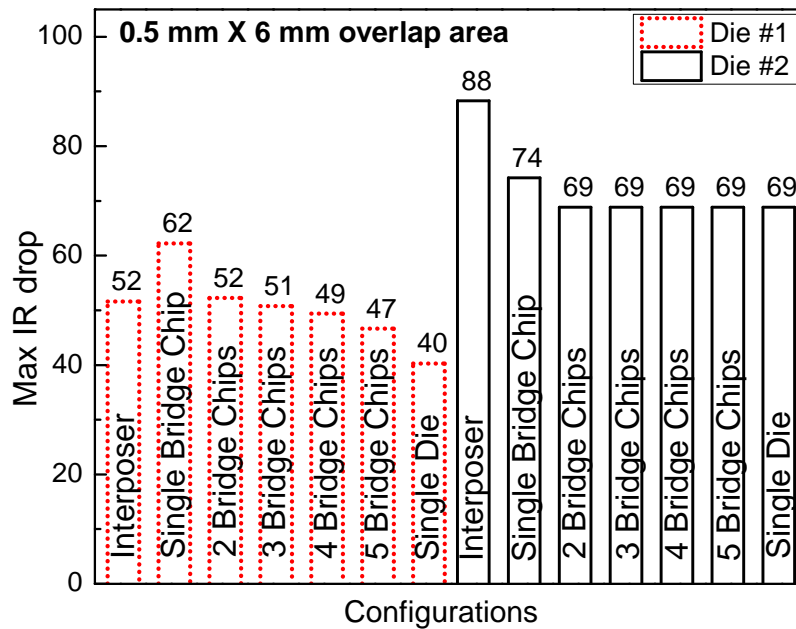


Figure 60: IR-drop analysis results using 5 bridge chips.

To investigate the above effect, we consider multiple bridge chips: from a single large bridge chip to five bridge chips with the same aggregate area. The results are shown in Fig. 60. If we use five bridge-chips, instead of a single large bridge-chip, Die #2 has a nominal IR-drop increase while Die #1 only has a 17.5% IR-drop increase. Compared to the silicon interposer case, the bridge-chip configuration can achieve a smaller IR-drop

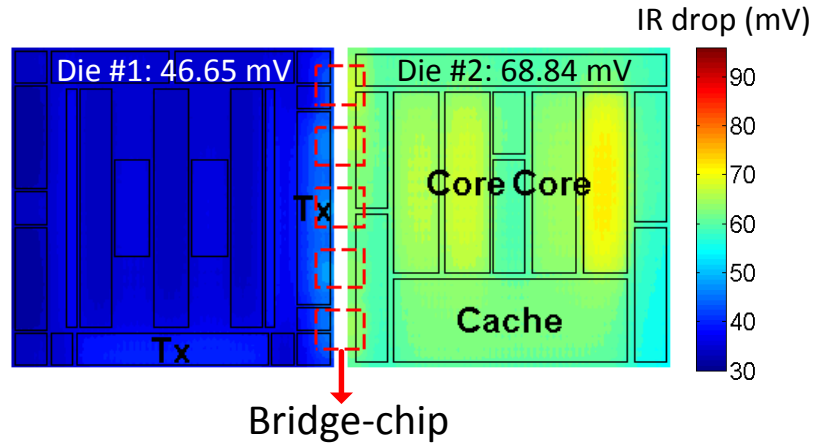


Figure 61: The IR drop profile of each die for the case with 5 bridge chips.

using multiple bridge-chips.

However, there is a tradeoff between manufacturing complexity and power delivery noise mitigation when using multiple bridge chips. For example, as bridge-chip count increases, there is a need for larger number of assembly steps. Therefore, to reduce the manufacturing complexity, we prefer less number of bridge chips. On the other hand, we prefer more bridge chips for smaller PSN. However from Fig. 60, we find that the rate of PSN reduction of both dice almost saturates when the bridge-chip count is five and beyond which indicates there is a diminishing return on PSN as bridge-chip count increases.

Fig. 61 shows the IR-drop noise profile for the case using 5 bridge chips. For the single bridge-chip case, there are two noise hotspots at the edges of the two dice due to insufficient number of C4 power/groups bumps, as shown in Fig. 58(c). However, there are no clear hotspots on the edges if we utilize, for example, five bridge-chips, as shown in Fig. 61.

Transient Droop

For transient analysis, the supply noise results from the switching current. Fig. 62(a) shows the impedance analysis results of an on die node. The chip operating frequency (> 1 GHz) is higher than the resonant frequency (about 150 MHz), therefore we only consider two waveforms with different frequencies (1 GHz and 4 GHz). The two waveforms are

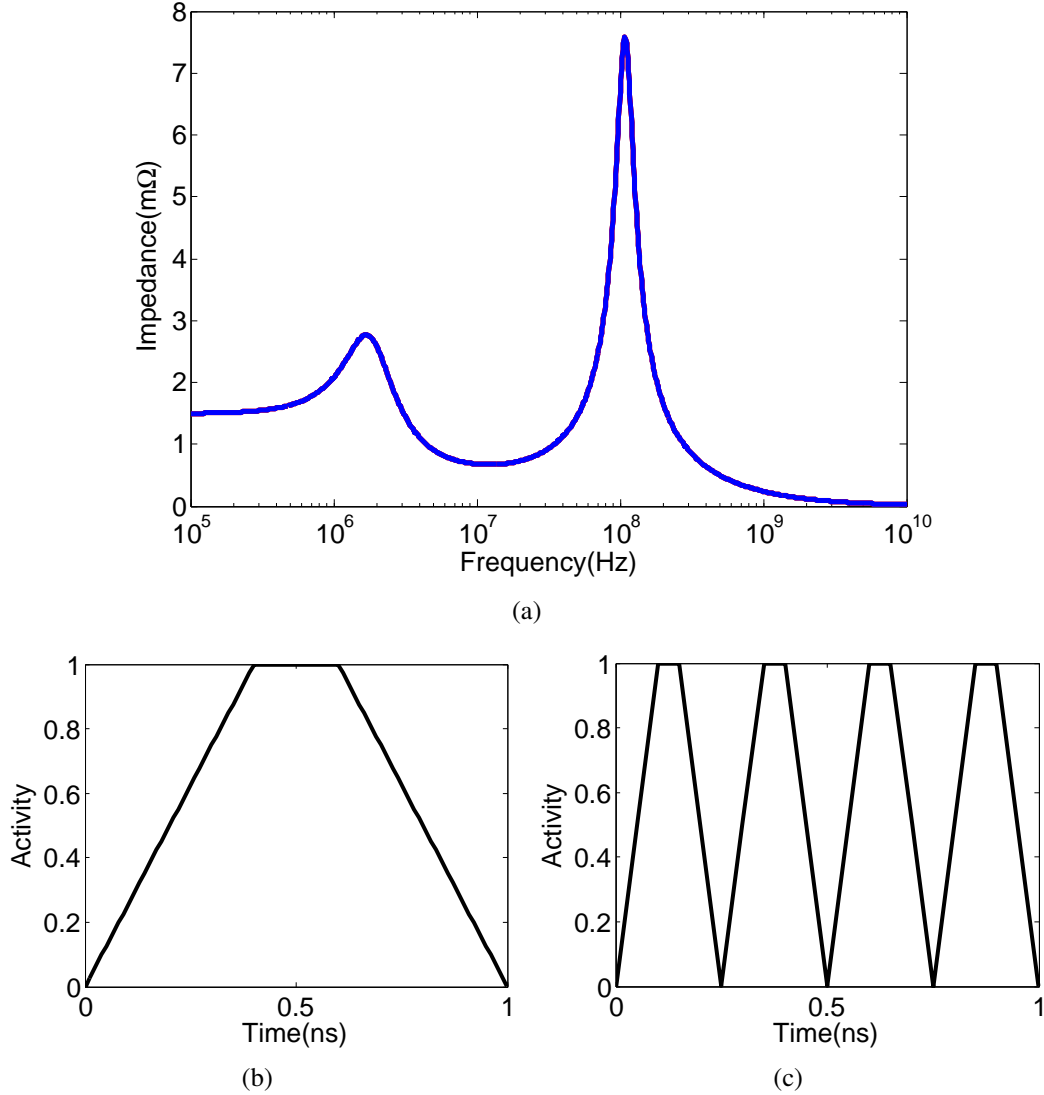


Figure 62: (a) Impedance analysis of one on-die PDN node and illustration of the switching current activity (b) waveform #1 1 GHz frequency (c) waveform #2, 4 GHz frequency

illustrated in Fig. 62(b) and 62(c). Waveform #1 has a rise time, pulse time, fall time and period of 400 ps, 200 ps, 400 ps, and 1000 ps, respectively and waveform #2 is four-time the frequency of waveform #1, as shown in Fig. 62(c). For the bridge-chip case, five bridge-chips are utilized.

The results are summarized in Table. 10. Compared to the single-die case, the interposer case provides the worst PSN due to the inductance of TSVs. However, the difference between the evaluated cases is not as significant as was in the IR-drop analysis since the

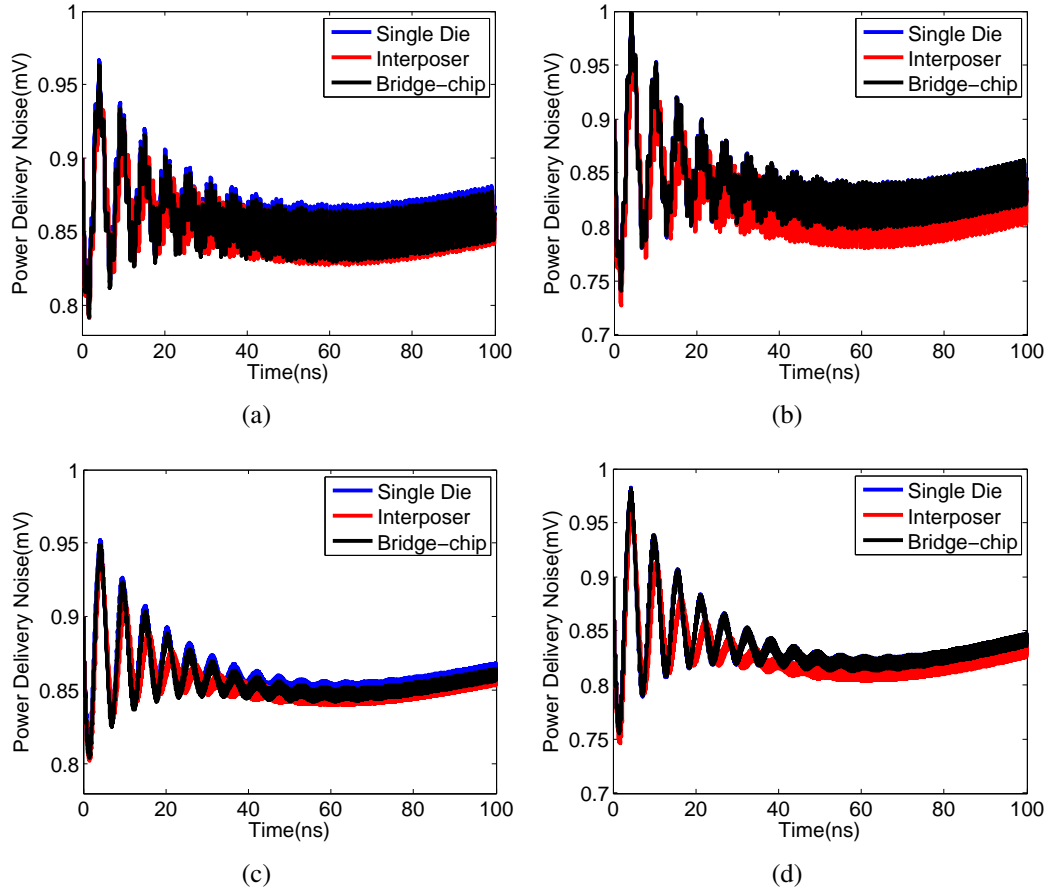


Figure 63: Transient analysis results of waveform #1 (a) Die #1 (b) Die #2 and waveform #2 (c) Die #1 (d) Die #2

Table 10: Transient state analysis results

Unit: mV	Waveform #1		Waveform #2	
	Die #1	Die #2	Die #1	Die #2
Single-die	102.33	159.85	91.11	145.28
Interposer	108.77	172.32	97.83	153.86
Bridge-chip	107.50	159.85	95.91	145.28

inductive parasitics are dominated by the package inductance.

For the bridge-chip case, since there is no hotspot at the edge of Die #2, as shown in Fig. 58(b), the switching noise is the same as in the single-die case. But for Die #1, there are hotspots at the edge (based on our assumed power maps), thus the switching noise is higher than in the single-die case. Nevertheless, the bridge-chip case produces better results than the interposer case. The transient waveforms of the worst node in single-die, interposer

and bridge-chip cases are shown in Fig. 63. The frequency of waveform #1 (1 GHz) is much closer to resonant frequency than that of waveform #2 (4 GHz), therefore, waveform #1 produces a much larger on-die noise swing. Except for that, their mid-frequency and low-frequency responses are similar, which make Fig. 63(a) and 63(c) as well as Fig. 63(b) and 63(d) have similar steady state values.

4.4 Design space exploration of 2.5-D integration

In this Section, we explore the impact of technology parameters such as total current, metal layers, TSV parameters, and overlap area. In addition, we propose to insert power/ground vias into the bridge-chip to mitigate PSN. Through these analyses, the challenges of bridge-chip based 2.5-D integration are better understood.

If not specified, the parameters are the same as those used in Section 4.3. However, for power maps, we assume uniform power distribution to eliminate the impact of on-die power variation. In our study, we focus on one of the two dice and fix the parameters of the other die the same for all the studies. The default total current for Die #1 and Die #2 is 100 A and 50 A, respectively.

4.4.1 Impact of total current requirement

We sweep the total current of Die #1 from 10 A to 100 A and plot the results for single-chip, interposer, single bridge-chip and 5 bridge-chip cases, as shown in Fig. 64. The results show a linear relationship between IR-drop noise and total current. Therefore for high-power systems, using more bridge chips is critical to reduce IR-drop. Under a total current of 100 A, the IR-drop for using 5 bridge chips is approximately 20.9% lower than using a single bridge chip.

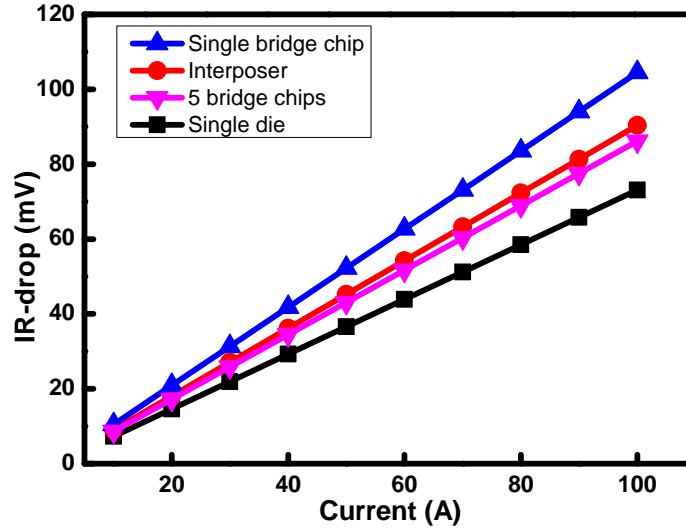


Figure 64: The impact of total current.

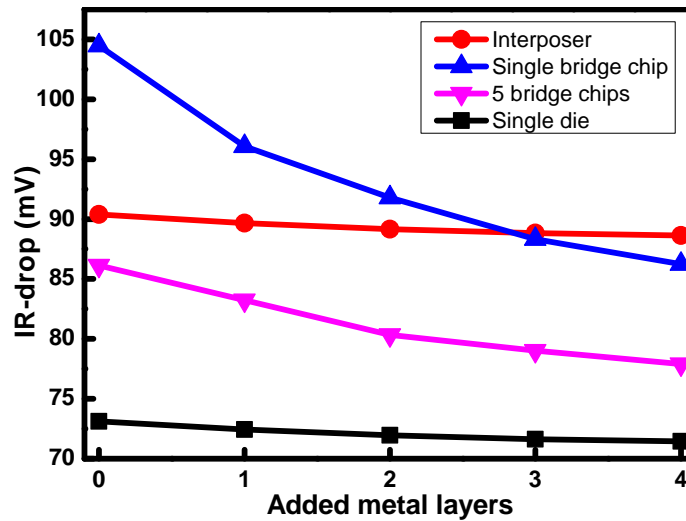


Figure 65: The impact of adding metal layers

4.4.2 Impact of on-die metal layers

On-die Metal layers help laterally spread the current from C4 bumps to functional blocks and therefore play an important role in power delivery. With more global metal layers, the noise is expected to drop. We explore the impact and consider adding 4 metal layers at most to the default configuration, as shown in Fig. 65. In this experiment, we assume a worst case where Die #1 has a uniform total current of 100 A.

From Fig. 65, we observe that adding on-die metal layers will make a large difference

to bridge-chip based integration because the overlap area lacks C4 bumps and the current has to be laterally spread through the on-die PDN. With a more conductive on-die PDN, the current spreading is enhanced. With one additional metal layer added, the IR-drop of the single bridge chip and 5 bridge chips cases reduces by 7.7% and 3.5%, respectively, and will further reduce by 17.3% and 10.5% if four metal layers are added. Similar trend is found for single die and interposer cases, but since the current already has good spreading, the benefits of adding metal layers is not significant.

From a technology point of view, while it is possible to mitigate PSN through adding metal layers, there are cost, manufacturing and performance tradeoffs.

4.4.3 Impact of TSV and overlap area

For interposer and bridge-chip 2.5-D integration, the key parameter is TSV and overlap area, respectively. To investigate their impact, we assume the TSV has a fixed aspect ratio of 15 and sweep the interposer thickness from 30 μm to 300 μm , which results in a TSV diameter of 2 μm to 20 μm . On the other hand, for bridge-chip case, we sweep the overlap area from 0.5 mm \times 6 mm to 2 mm \times 6 mm (width of overlap area changes from 0.5 mm to 2 mm).

The results are shown in Fig. 66. In order to clearly compare interposer and bridge-chip 2.5-D integration, we plot the results of both cases using the same Y-axis. For interposer 2.5-D, as TSV diameter reduces, TSV resistance increases which presents challenges to power delivery. The IR-drop noise of the case using 2 μm diameter TSVs is two times the noise of the case with 20 μm diameter TSVs. For bridge-chip case, as the overlap area increases, the IR-drop inevitably increases since the center of overlap area becomes further to the nearest C4 bumps. However, with multiple bridge chips, IR-drop is less sensitive to the overlap area than the single bridge-chip case and it only incurs an IR-drop increase of 14.0% when the overlap area changes by four times.

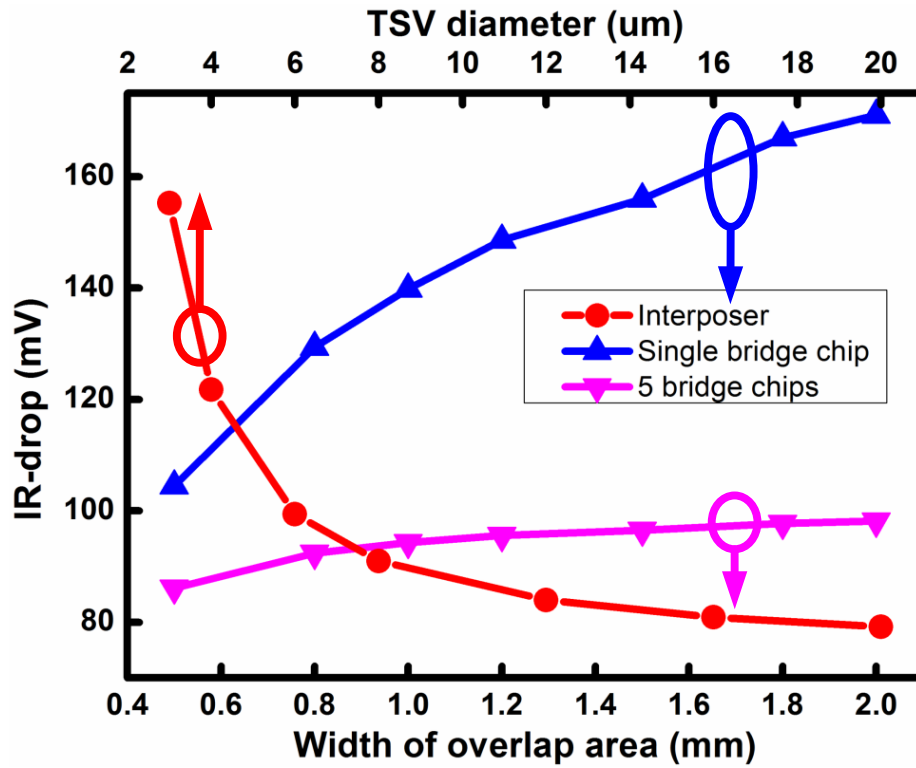


Figure 66: The impact of TSV and overlap area

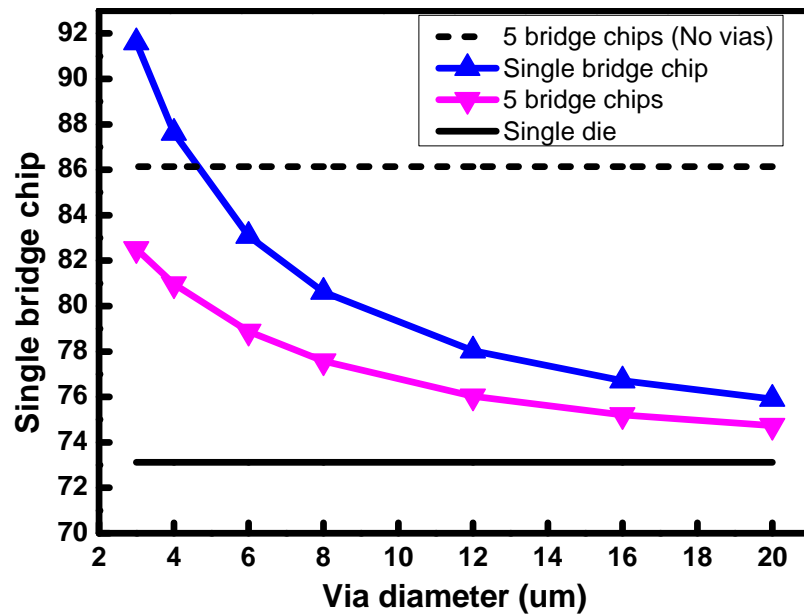


Figure 67: The impact of inserting vias in the bridge-chip

4.4.4 Inserting vias in bridge-chip

Similar to interposer 2.5-D, vias can also be inserted through the bridge chip. With the help of these vias, the power delivery path in the overlap areas is shortened.

Similar to the previous section, we assume that the through bridge-chip vias (TBVs) have a diameter from $2\ \mu\text{m}$ to $20\ \mu\text{m}$. The simulation results of single bridge-chip and 5 bridge-chip cases are shown in Fig. 67. To show the worst and best bounds of PSN, we include the results of 5 bridge chip without vias and the single-die case. With these vias, no matter whether a single or multiple bridge chips are used, the IR-drop can be reduced to the level of a single-die case with appropriate via dimensions. If we use a large via diameter of $20\ \mu\text{m}$, the IR-drop is only 2.2% and 3.8% larger than the single-die case.

4.5 Conclusion

This chapter presents a PDN modeling framework focusing on on-die and on-package PDN, which is used for design space exploration of 2.5-D and 3-D integration platforms. The model is then validated against *IBM* power grid benchmarks, with a maximum relative error of less than 7.29% and 0.67% of VDD for IR-drop and transient analysis, respectively. Interposer and bridge-chip 2.5-D integration technologies are benchmarked and compared. Interposer based 2.5-D integration generally exhibits larger IR-drop and transient droop than bridge-chip based 2.5-D integration due to TSV parasitics. While bridge-chip based interconnection platforms present PDN challenges, especially to the active die regions that overlap with the bridge-chip, results suggest minimizing this overlap region and using multiple bridge chips instead of a single large bridge helps to mitigate PSN. Additionally, adding more metal layers is also useful to reduce PSN, one additional metal layer may achieve 7.7% PSN reduction. Lastly, we propose to insert through bridge chip vias to address PSN issue in bridge-chip 2.5-D, and the results show the IR-drop is only 2.2% larger than single-die case when using a $20\ \mu\text{m}$ via.

CHAPTER 5

INTEGRATED THERMAL AND POWER DELIVERY NETWORK CO-SIMULATION FRAMEWORK

This chapter presents a thermal and power delivery network (PDN) co-simulation framework for single-die and emerging multi-die configurations that incorporates the interactions between temperature, supply voltage, and power dissipation. The temperature dependencies of wire resistivity and leakage power are considered and the supply voltage dependencies of power dissipation are modeled. Starting with a reference power dissipation, the framework is capable of evaluating the temperature distribution and PDN noise simultaneously and eventually updating the power dissipation based on the thermal and supply voltage distributions.

5.1 Introduction of thermal and PDN co-simulation

Fig. 68 shows the dependencies between power dissipation, temperature and PDN noise. The temperature impacts the leakage power and the power grid resistivity. Power dissipation determines the source current of the chip and is also the excitation of the PDN noise. Reversely, the power supply voltage impacts both leakage and dynamic power. Without considering the interactions between each of the components in Fig. 68, the results of the standalone models are inaccurate. For example, Su *et. al* noted that the leakage power was underestimated by as much as 30% without including the impact of temperature and power supply voltage [30]. Hence, it is essential to build a thermal and PDN noise co-simulation framework to answer ‘what-if’ type questions for design space exploration for 2.5-D and 3-D ICs in early design stages.

Prior work focused on either developing the individual thermal [77] or PDN models [24, 78] or studying parts of the interactions [34, 79, 80, 81, 79] shown in Fig. 68. There

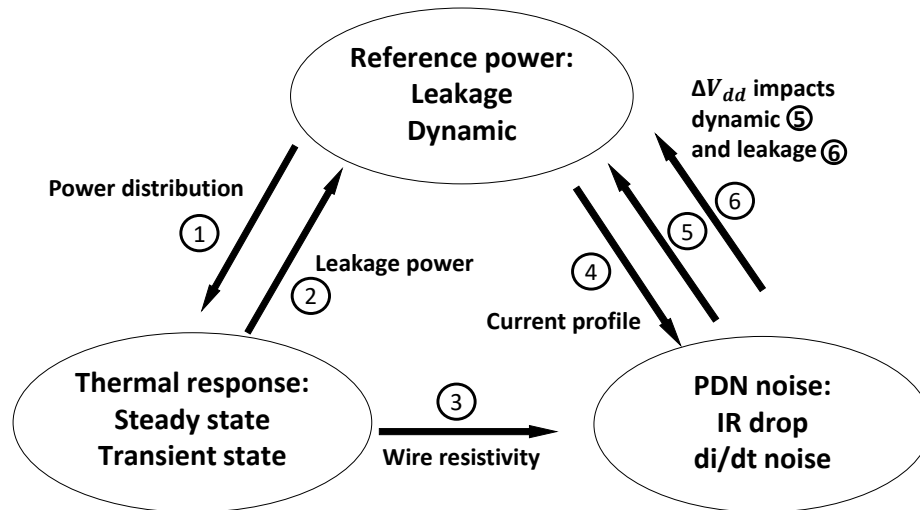


Figure 68: The interactions between temperature distribution, PDN noise and power dissipation

are no co-simulation frameworks capable of performing steady- and transient- state analysis on thermal and PDN noise while incorporating the impact of their variation on power dissipation. Xie *et. al* only studied the interactions between thermal and IR-drop and did not consider the interactions with power dissipation [37]. Although Su *et. al* investigated the temperature and supply voltage dependencies of power dissipation, the thermal impact on wire resistivity and the transient-state analysis were not included [30].

We propose a framework to simultaneously study the temperature, PDN noise, power dissipation and the interactions between them for both steady- and transient-state analysis. The thermal model is based on finite volume method, the PDN model is based on finite difference method and the interaction models are based on the thermal and supply voltage dependencies of power dissipation and the thermal dependencies of wire resistivity. The initial reference power is an input from *McPat* [57] and by using our thermal-power and PDN-power models, we update the power dissipation until the iterations are converged. There are two loops in the framework where the outer loop iterates the thermal-power models and the inner loop iterates the PDN-power models.

5.2 Simulation flow of integrated thermal and PDN modeling framework

A simulation framework for steady-state and transient-state analysis is presented in this section. We assume an architectural tool has already provided the reference power of the chip under uniform temperature and an ideal power supply voltage through the initial power simulations. Since our focus is the impact of supply voltage and temperature, we assume other parameters such as clock frequency remain constant. In the following iterations, the power dissipation is updated by the power models instead of calling the power simulator every iteration. Using the initial power, we start the simulation flow and perform the thermal and power supply noise simulations. At the end of the simulations, the three metrics become consistent with each other within our interaction models. The simulator is implemented using MATLAB because dense matrix operations and calculations are required in the flow.

5.2.1 Steady-state analysis

For steady-state analysis, firstly, the reference power is used to obtain the temperature distribution. With the thermal results, the PDN grid resistance within the chip is updated. The supply voltage profile is then simulated using the updated PDN grid resistance and source current distribution. Next, the leakage and dynamic powers are updated based on the thermal and supply voltage values. For this step, the power and supply voltage form a loop and the Newton method is used to accelerate the convergence rate. After this loop is complete, the thermal profile is updated and checked for whether the convergence has been reached. If not, the simulation restarts. The simulator finally returns the thermal profile, PDN noise distribution and the updated reference power results of the system. Fig. 69 shows the whole simulation flow. Because the thermal conductivity of common materials such as copper, silicon, silicon dioxide etc, is usually constant in the typical temperature range of IC operation, the thermal impact on their material properties is neglected.

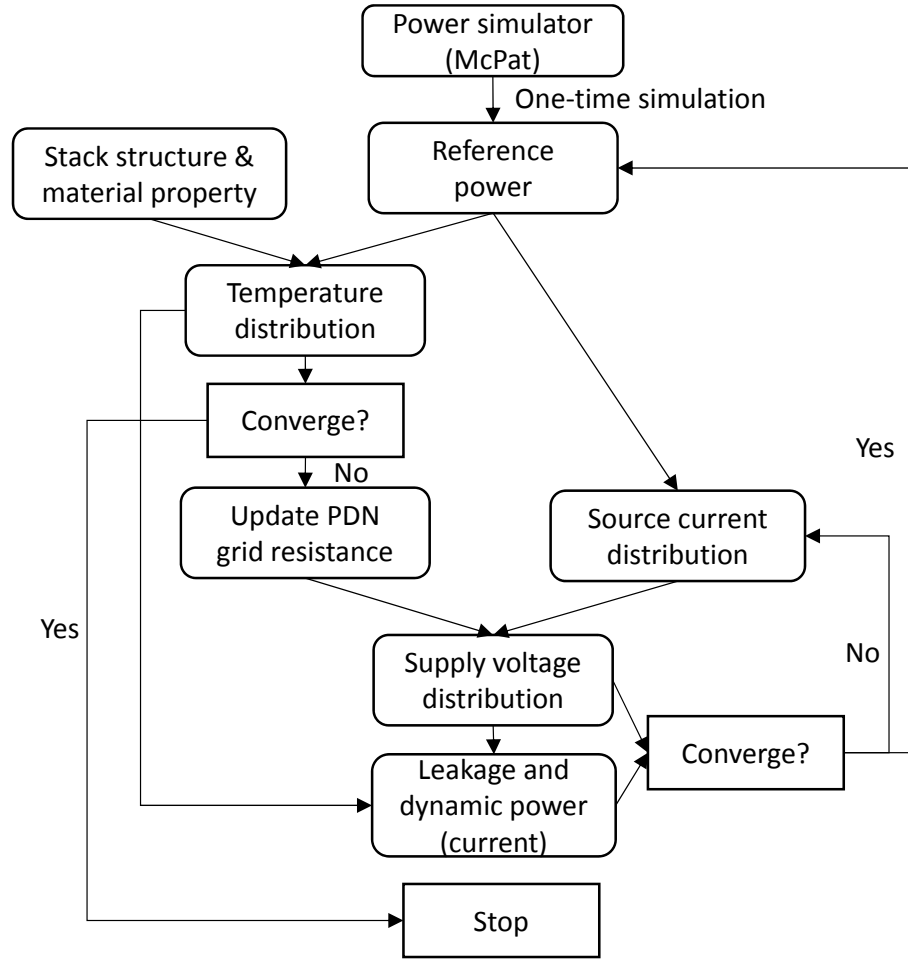


Figure 69: Steady-state simulation flow.

5.2.2 Transient-state analysis

The thermal response time of a typical single or multi-die package with either conventional air cooled heat sink (ACHS) or microfluidic heat sink (MFHS) is much larger than the response time of the PDN. The thermal response time is at least in the milliseconds range [77] and the response time of PDN is within the circuit switching frequency (in the nanoseconds or microseconds range). Therefore, in the PDN simulation time scale (up to microseconds [26]), the thermal profile remains constant. The validation is discussed in Section 3.3.

Based on the above discussion, a one-time steady-state or transient thermal simulation is initially performed to obtain a temperature profile as an input. For this step, one option is to perform a steady-state simulation using a user-defined reference power representing the

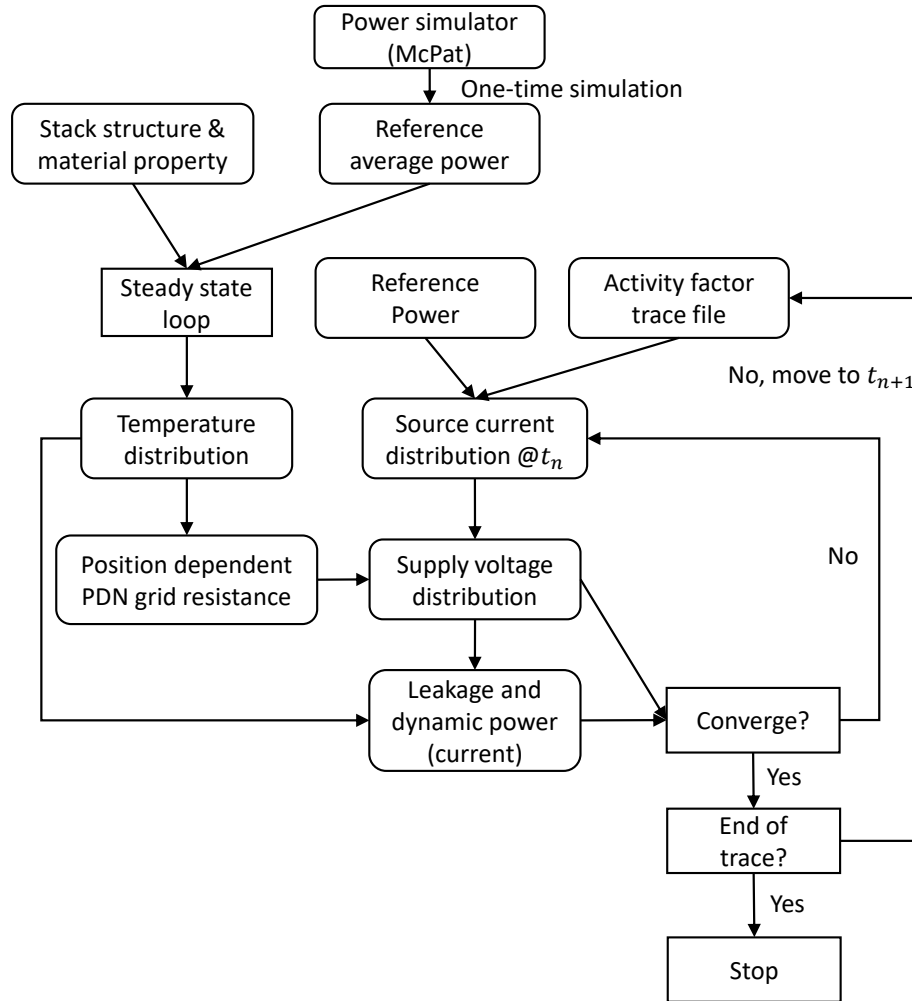


Figure 70: Transient-state simulation flow.

average power or another option is to perform cycle-granularity transient thermal simulation from the very beginning to obtain the final thermal profile. Next, we start the transient simulation of the PDN using the thermal inputs. In each time step, similar power-PDN solving iterations are performed as in the steady-state analysis. Here, we assume leakage and dynamic powers (currents) are changing instantaneously as supply voltage changes [82]. When the loop of the current time step is converged, the simulation of the next time step is started until the end of the trace. The simulator finally returns the transient PDN noise distribution and updated power results. The transient analysis flow is shown in Fig. 70.

5.2.3 Framework algorithm

The pseudo code for the steady-state analysis is shown in Algorithm 1. Transient-state analysis is similar but includes the thermal capacitive and electrical capacitive/inductive elements, therefore we do not show it here. Because the Newton method is used when solving the PDN-power loop, the power update model needs to calculate the partial derivatives over supply voltage (lines 11 to 15). The thermal loop (lines 20 to 28) performs fixed point iterations and for the PDN-power loop (line 22 to 25), the Newton method is used.

5.3 Modeling methodologies and implementation

In this section, the thermal, PDN and power update models are presented. The formulation of the interactions between each component is described. The explicit interactions considered are: temperature-wire resistivity, temperature & supply voltage-leakage power and supply voltage-dynamic power.

5.3.1 Thermal model and formulation

The thermal model has three inputs: the first is the geometry information of the single-chip, 2.5-D module or 3-D stack, the second is the material property of each layer, and the third is the reference power information. The power granularity can be block-level or transistor-level. The formulation of the thermal model is shown below:

$$-G \cdot T_{n+1} + C \cdot \frac{T_{n+1} - T_n}{\delta t} = P_{n+1}(T) + H \cdot T_{amb} \quad (5.1)$$

where T_{n+1} and T_n are the temperatures of the current (to be solved) and the previous (known) time steps, respectively. G is the thermal conductance matrix, C is the heat capacity matrix and $P(T)$ is the power excitation of which the leakage component is temperature dependent. H is a diagonal matrix which represents the convective boundary condition, and T_{amb} is the ambient temperature. The temperature time difference term is only applicable

Algorithm 1 Steady state simulation

```
1: function INITIALIZATION ▷ parse input tables
2:    $G_{elec} \leftarrow$  electrical conductance matrix
3:    $G_{ther} \leftarrow$  thermal conductance matrix
4:    $I_{leak} \leftarrow$  reference leakage current
5:    $I_{dyna} \leftarrow$  reference dynamic current
6:    $V_{node} \leftarrow V_{ref}$  ▷ reference voltage
7:    $T_{node} \leftarrow T_{ref}$  ▷ i.e. ambient temperature
8:    $P_{tot} \leftarrow (I_{leak} + I_{dyna}) \cdot V_{node}$ 
9: end function
10:
11: function PDN_SOLVER( $G, I, V$ ) ▷ solve  $-G \cdot V = I$ 
12:    $G_{new} = G + \partial I_{tot} / \partial V$  ▷ Newton method
13:    $\delta I \leftarrow G \cdot V + I$ 
14:    $output \leftarrow V + \delta I / G_{new}$ 
15: end function
16:
17: procedure STEADY_STATE
18:   Initialization()
19:    $V, T \leftarrow V_{node}, Therm\_solver(G_{ther}, P_{tot})$ 
20:   while  $|T_{node} - T| > \epsilon \cdot T$  do
21:      $G_{elec.new} \leftarrow Resistivity\_Update(T, G_{elec})$ 
22:     while  $|G_{elec.new} \cdot V + I_{tot}| > \epsilon$  do
23:        $I_{tot} \leftarrow Current\_Update(V, T, I_{leak}, I_{dyna})$ 
24:        $V \leftarrow PDN\_solver(G_{elec.new}, I_{tot}, V)$ 
25:     end while
26:      $T_{node} \leftarrow T$ 
27:      $T \leftarrow Therm\_solver(G_{ther}, P_{tot})$ 
28:   end while
29:    $P_{tot} \leftarrow I_{tot} \cdot V$ 
30:    $output \leftarrow T, V, P_{tot}$ 
31: end procedure
```

for the transient analysis and is not applicable for the steady-state analysis.

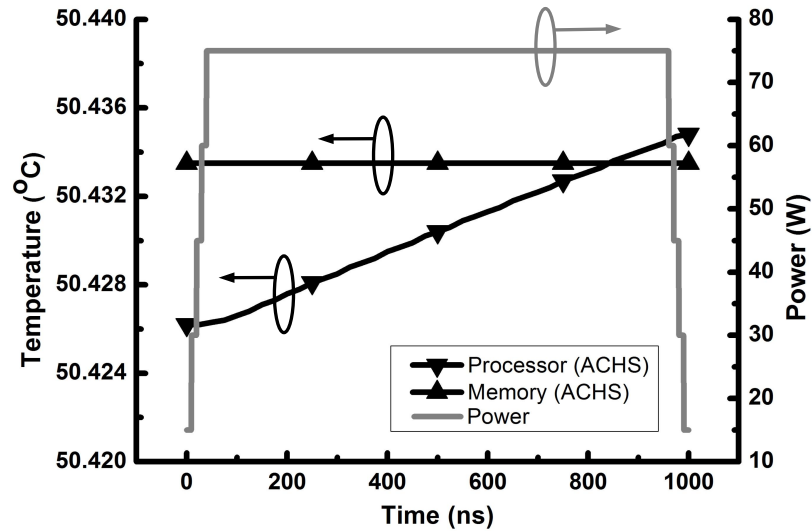


Figure 71: Validation of stable temperature assumption in microsecond scale.

As previously stated in Section 5.2.2, due to the large response time, the temperature profile remains constant in the time scale of microseconds even though the power experiences a sharp change. We simulated a test case consisting of a 3-D stack as shown later in Fig. 74 (Section 5.4) to validate this. The thermal specifications can be found in Section 5.4. Fig. 71 shows there is a nominal change (less than 0.1%) in the maximum temperature of the processor and memory dice when the processor power changes dramatically (memory power remains the same). The thermal profile of the chip is also checked, and it does not undergo any change.

5.3.2 PDN model and formulation

The block diagram of the PDN network is shown in Fig. 72. Here, the distributed power/ground (P/G) resistance network is abstracted for visualization. The detailed structure of the VDD/GND rail is a distributed wire network. Since our study focuses on the on-chip PDN modeling, we use a simplified package model where each P/G port is connected to a lumped resistor and inductor pair. The trapezoid scheme is used to formulate the transient

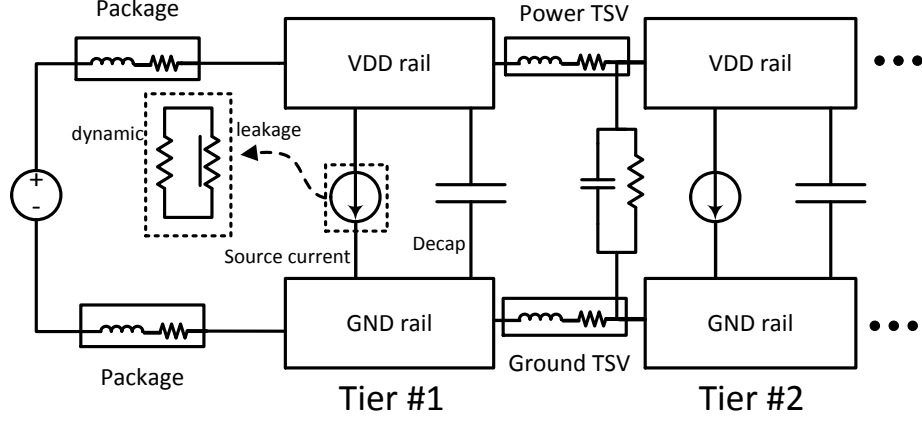


Figure 72: Block diagram of the PDN structure. The distributed power/ground rail is abstracted for visualization.

finite difference equation [25]. The formulation is shown as follows:

$$\left(\frac{K}{\Delta t} + \frac{U}{2}\right) \cdot X^{n+1} = \left(\frac{K}{\Delta t} - \frac{U}{2}\right) \cdot X^n + \frac{I_s^{n+1} + I_s^n}{2} \quad (5.2)$$

where

$$U = \begin{bmatrix} G(T) & A_L \\ -A_L & R(T) \end{bmatrix} \quad K = \begin{bmatrix} C & 0 \\ 0 & L \end{bmatrix} \quad (5.3)$$

$$X = \begin{bmatrix} V \\ I \end{bmatrix} \quad I_s = \begin{bmatrix} i_s(V, T) \\ 0 \end{bmatrix}$$

where $G(T)$ and $R(T)$ are the PDN grid conductance matrix and the PDN grid port resistance matrix (both are temperature dependent), respectively. C AND L are the matrices reflecting the capacitive and inductive elements, respectively. $i_s(V, T)$ is the source current which is dependent on temperature (due to leakage portion) and supply voltage (due to both leakage and dynamic portions). The temperature dependent PDN grid resistivity (G and R) is described below:

$$\rho = \rho_0(1 + \alpha(T - T_0)) \quad (5.4)$$

where ρ_0 is the resistivity under reference temperature T_0 and α is the temperature coefficient of resistivity.

5.3.3 Power update models

In this work, we do not aim to develop a detailed power analysis model such as *McPat* [57] but instead pursue a power update model based on distributed temperature and supply voltage. Therefore, we focus on the impact of supply voltage and temperature, while other parameters such as clock frequency are assumed to be constant. We begin with the power results from *McPat* and by considering the supply voltage and thermal variation through the whole chip, we update the value of power dissipation.

It is assumed that the power of each functional block consists of leakage and dynamic powers. In the PDN model, the power is converted to source current, as shown in Fig. 72. Prior work assumed that the source current (I) can be calculated by $I = P/V_{dd}$, where P is power dissipation and V_{dd} is the value of ideal supply voltage [24]. This simple method does not consider the power dependency of supply voltage and temperature. Instead, modeling the source current as two resistors (the dashed inlet box in Fig. 72) of which values are a function of supply voltage and temperature will capture the dependencies [23].

Leakage power

The relationship between reference leakage power $P_{leak.ref}$ and leakage current source $I_{leak.ref}$, is expressed as follows:

$$P_{leak.ref} = I_{leak.ref} \cdot V_{dd} \quad (5.5)$$

where V_{dd} is the ideal supply voltage of each power grid. Based on the fitting method [30], the actual leakage current of a node, $I_{leak.act}$, can be generalized as:

$$I_{leak.act} = I_{leak.ref} \cdot f(V, T) \quad (5.6)$$

where $f(V, T)$ is the fitted function of supply voltage and temperature of the node in the chip. In this work, we propose to use a 2-D piecewise linear model for $f(V, T)$. The first advantage of a 2-D piecewise linear model is to cover a wide range of voltage and temperature values ([30] only covers a small range around the reference point). Second, the partial derivative of leakage current over temperature and voltage can be easily calculated so that the Newton method can be implemented to accelerate the whole simulation. For a target circuit, the voltage-temperature plane is uniformly meshed based on the number of sampling points. Next, for each sampling point, i.e. (V_i, T_i) , we run HSPICE simulations and collect the data. The leakage of an arbitrary point (V, T) can then be calculated using $f(V, T)$ as shown below:

$$V_i \leq V \leq V_{i+1}; T_i \leq T \leq T_{i+1} \quad (5.7)$$

$$\xi = \frac{V - V_i}{V_{i+1} - V_i}; \eta = \frac{T - T_i}{T_{i+1} - T_i}; \quad (5.8)$$

$$wt = \begin{bmatrix} \xi * \eta \\ (1 - \xi) * \eta \\ (1 - \eta) * \xi \\ (1 - \xi) * (1 - \eta) \end{bmatrix} I_{node} = \begin{bmatrix} I_{leak}(V_i, T_i) \\ I_{leak}(V_{i+1}, T_i) \\ I_{leak}(V_i, T_{i+1}) \\ I_{leak}(V_{i+1}, T_{i+1}) \end{bmatrix} \quad (5.9)$$

$$f(V, T) = \frac{wt^T \cdot I_{node}}{I_{leak.ref}} \quad (5.10)$$

Although it is difficult to use just one circuit and apply the characterization results to other circuits, the fitted $f(V, T)$ of an inverter array is quite accurate to use based on the results of [23, 83]. In this work, we use 50 stage inverter pairs (for each inverter pair, the input of the first inverter is connected to V_{dd} and the input of the second is connected to ground, and the output of both inverters is floated). PTM-MG 20nm (HP) model [84] is

used for *HSPICE* simulation and the parameter range for supply voltage and temperature are $(0.7V \sim 1.0V)$ and $(25^{\circ}C \sim 110^{\circ}C)$, respectively, the reference voltage is $0.9V$, and the reference temperature is $100^{\circ}C$, a pessimistic temperature as most of the IC design tools assume.

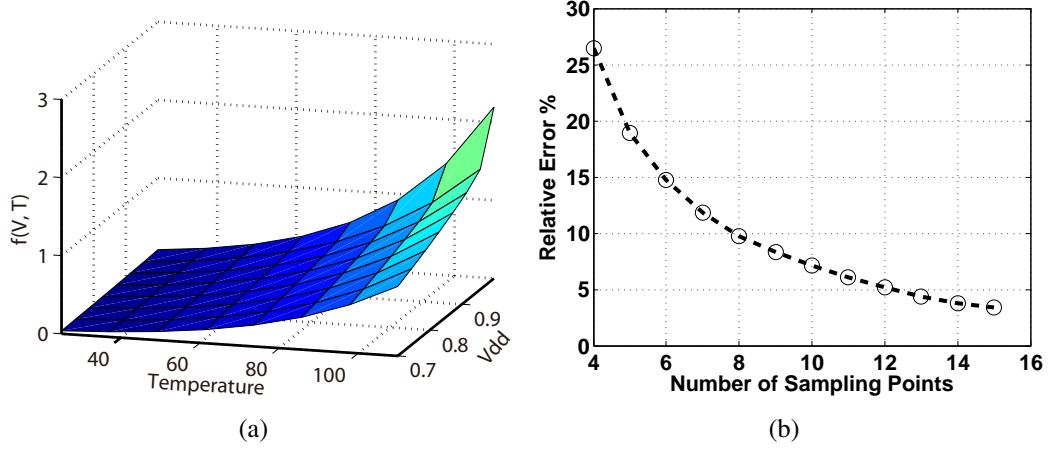


Figure 73: The 2-D piecewise linear model (a) response surface with 8 sampling points (b) the maximum error of models with different number of sampling points.

Fig. 73(a) shows the surface response of the model with 8 sampling points. We generate 1600 random data points in the (V, T) plane for validation. Fig. 73(b) shows that the model accuracy increases as the number of sampling points increases and the error drops below 5% when the number of sampling points is larger than 13.

Dynamic power

The reference dynamic power $P_{dyna.ref}$ is expressed below:

$$P_{dyna.ref} = \alpha \cdot C \cdot f \cdot V_{dd}^2 \propto V_{dd}^2 \quad (5.11)$$

Where α is the activity factor, f is the frequency and C is the total capacitance. The power grid actually gets a supply voltage of $V_{dd.act}$ instead of V_{dd} when calculating reference

power, thus the dynamic power becomes:

$$P_{dyna.act} = P_{dyna.ref} \cdot \frac{V_{dd.act}^2}{V_{dd}^2} \quad (5.12)$$

By converting the dynamic power to a current source, we have the dynamic current update model [30] [23] :

$$I_{dyna.act} = P_{dyna.ref} \cdot \frac{V_{dd.act}}{V_{dd}^2} \quad (5.13)$$

5.4 Comparison of models with different number of dependencies included

In this section, we use a 3-D processor-on-memory stack as an example to demonstrate the capability and accuracy of the modeling work.

5.4.1 Thermal and PDN specification

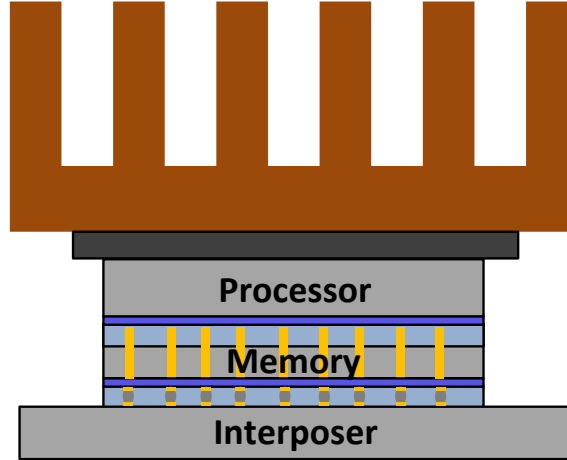


Figure 74: The 3D-IC example: processor on memory sack.

The stack we evaluate is shown in Fig. 74. The processor is placed on top of the memory for thermal considerations. The thickness of each layer and thermal conductivity of each material are shown in Table 11. The reference power maps of the memory and processor dice are shown in Fig. 75(a) and (b) [51]. The reference temperature is $100\text{ }^{\circ}\text{C}$, and the supply voltage is 0.9 V for 22 nm multi-gate ICs [84]. Based on the simulations from

Table 11: Parameters for thermal model

	Conductivity W/mK	Thickness μm
TIM	3	25
Memory die	149	100
Underfill layer	0.9	25
Processor die	149	100
Micro-bump	60	25
Interposer	149	200
Copper	400	N/A
SiO_2	1.38	N/A

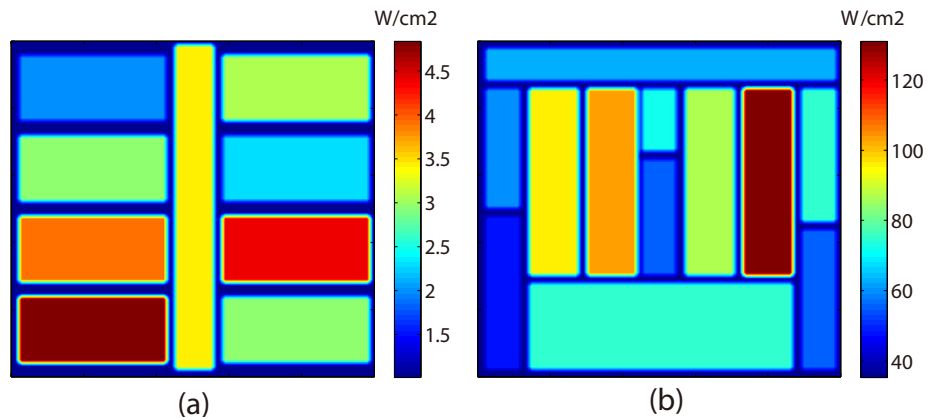


Figure 75: Reference power maps (a) Memory die (2.82W) (b) Processor die (74.49W).

McPat [57], under the reference temperature, 20% of the total power is leakage. The chip size is assumed to be $1\text{ cm} \times 1\text{ cm}$ and the interposer size is assumed to be $2\text{ cm} \times 1.5\text{ cm}$. The heat spreader is assumed to be $4.5\text{ cm} \times 3.5\text{ cm}$ and the air cooled heat sink is converted to a boundary condition of 0.24 W/K [56]. All the other faces are adiabatic. The ambient temperature is assumed to be $38\text{ }^\circ\text{C}$.

The parameters for the PDN model are shown in Table 12. In this paper, we focus on the global on-die PDN and it is assumed to consist of the top two metal layers. Each metal layer is assumed to be $5\text{ }\mu\text{m}$ thick. Fig 76 shows the detailed geometry and configuration of the interleaved global PDN [26]. Besides P/G TSVs, we assume there are 10,000 signal TSVs with $5\text{ }\mu\text{m}$ diameter and $0.5\text{ }\mu\text{m}$ thick liner. The TSVs are assumed to be uniformly distributed because of the high-power processor. The reference source current is calculated

Table 12: Parameters for PDN model

	value
P/G TSV diameter	$5 \mu m$
on-die decap %	10% of the area
# of P/G TSV	676/625
Package inductance	0.1 nH
Package resistance	0.0107Ω
P/G wire thickness	$5 \mu m$
P/G wire width	$3.33 \mu m$
P/G wire pitch	$30 \mu m$

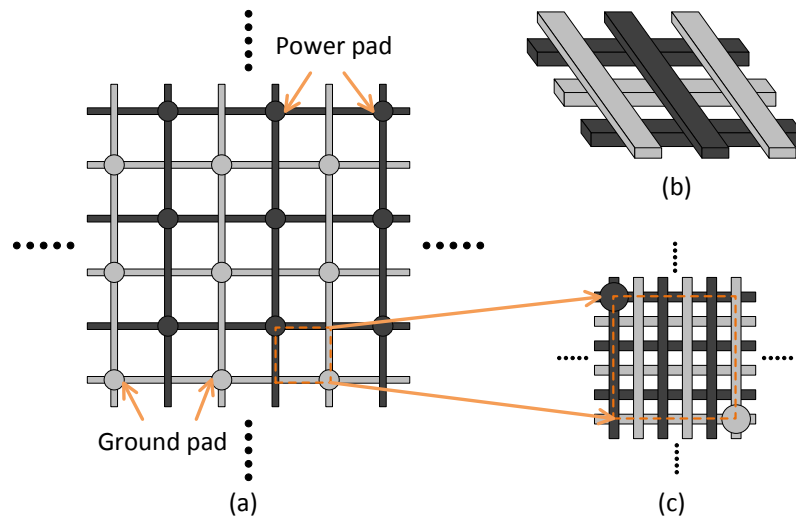


Figure 76: On-die PDN structure (a) Power/ground pads with wires (b) Interleaved structure of power/ground wires. (c) Dense PDN wires between power/ground pads

using Eq. 5.6 and 5.13. For the transient PDN analysis, we consider a worst case scenario where both dice switch from full sleep mode to peak power dissipation (power map shown in Fig. 75) in a rise time of 100 ps and remain in peak power mode for 50 ns. Based on Section 5.3, when performing transient simulation, the thermal maps do not change over the time scale of the PDN analysis. For this case, we assume both dice have been in the peak power state long enough that the thermal profiles with maximum temperatures are reached.

Table 13: Simulation Model

Model	Description	Arrows included in Fig. 68
standalone	Power results from <i>McPat</i> , individual thermal and PDN simulation	① ④
PDN-power	Interactions between power dissipation and PDN are added to standalone models	① ④ ⑤ ⑥
PDN-therm	Thermal impact on wire resistivity is added to PDN-power case	① ③ ④ ⑤ ⑥
PDN-therm-leak	Impact of PDN and thermal on leakage power and the thermal impact on wire resistivity are added to standalone models	① ② ③ ④ ⑥
partial-therm	Interactions between temperature and leakage power are added to PDN-power case	① ② ④ ⑤ ⑥
full-model	Thermal impact on wire resistivity is added to partial-thermal case	① ② ③ ④ ⑤ ⑥

5.4.2 Modeling Scenarios with different number of dependencies

In this section, we compare a number of models to establish the benefits of the proposed work. The first model (denoted as standalone model) provides fixed input power maps for the thermal and PDN models. Standalone thermal and PDN analysis are performed (constant reference temperature is used throughout the chip for PDN analysis). In the second model, we consider the interactions between power dissipation and PDN only (PDN-power model). The thermal effects of power dissipation and PDN wires are not included (same as in the standalone model) but in this case, the final updated power distribution is used to perform thermal analysis.

In the third model, we add the thermal impact on PDN wires to the PDN-power case (denoted as PDN-therm model) [28]. In the fourth model, we consider the thermal impact on both wire resistivity and leakage power but the interactions between PDN and power dissipation are not included (PDN-therm-leak model) [34].

In the fifth model, we add the interactions between thermal and leakage power to

PDN-power model but the thermal impact on grid resistivity is not included (partial-therm model). Lastly, we include all the interactions shown in Fig. 68 (full-model). The six models are summarized in Table 13. In practice, there are usually thermal, power and PDN constraints for a system. If under a configuration any of the metrics are out of bound, the configuration should be changed to meet the constraints such as lowering the target frequency. To model this, it is necessary to include an integrated power module in the framework. The author is aware of these constraints, but in this work, our focus is to present a way to co-simulate these metrics.

5.4.3 Comparison results

The simulation results are shown in Table 14. To easily compare them, the metrics of each model are normalized to those of the standalone model and are shown in the parentheses.

Firstly we analyze all the steady state analysis results (IR-drop, temperature and power). Comparing standalone and PDN-power model, we observe that the standalone model overestimates all the metrics by about 3% ~ 6%. This is because when the IC is operating, not a single power grid receives the ideal power supply voltage due to IR-drop. Based on Eq. 5.6 and Eq. 5.13, with a lower supply voltage, the actual source current becomes smaller than the reference one, resulting in lower power and as a consequence, the simulated temperature is smaller.

The PDN-therm model is similar to the PDN-power model with the only difference being that the thermal impact on wire resistivity is included. For the PDN-power model, the wire temperature is the reference temperature (100 °C) while the temperature for the PDN-therm model is about 90 °C. Based on Eq. 5.4, the wire resistivity changes 3.93% for a 10 °C temperature change (temperature coefficient of copper wire is $3.9 \cdot 10^{-3}/^{\circ}C$). The difference in IR-drop of the processor die between PDN-Power and PDN-thermal models has a good agreement to this number.

PDN-therm-leak model adds the thermal impact on both wire resistivity and leakage

Table 14: Results for different detailed models

Model	die	noise(mV)		temperature (°C)	power(W)	
		IR-drop	Transient		dynamic	leakage
Standalone	Processor	38.79	118.07	92.74	59.59	14.90
	Memory	5.32	83.85	91.31	2.26	0.56
PDN-power	Processor	37.11 (4.33%)	104.35 (11.62%)	89.70 (3.28%)	56.63 (4.97%)	13.91 (6.64%)
	Memory	5.13 (3.57%)	75.19 (10.33%)	88.51 (3.07%)	2.24 (0.88%)	0.56 (0.00%)
PDN-therm	Processor	36.02 (7.14%)	103.36 (12.46%)	89.81 (3.16%)	56.74 (4.78%)	13.95 (6.38%)
	Memory	5.10 (4.14%)	75.47 (11.19%)	88.61 (2.96%)	2.24 (0.88%)	0.56 (0.00%)
PDN-therm-leak	Processor	34.71 (10.52%)	107.57 (8.89%)	87.01 (6.19%)	59.59 (0.00%)	8.31 (44.23%)
	Memory	4.90 (7.89%)	76.50 (8.77%)	85.71 (6.13%)	2.26 (0.00%)	0.31 (44.64%)
partial-therm	Processor	34.41 (11.28%)	97.25 (17.63%)	85.37 (7.95%)	56.89 (4.53%)	7.42 (50.20%)
	Memory	4.74 (10.90%)	69.80 (16.76%)	84.22 (7.76%)	2.24 (0.88%)	0.30 (46.43%)
full-model	Processor	33.05 (14.80%)	96.02 (18.68%)	85.51 (7.80%)	57.02 (4.31%)	7.47 (49.87%)
	Memory	4.71 (11.47%)	70.18 (16.30%)	84.35 (7.62%)	2.24 (0.88%)	0.30 (46.43%)

The relative percentage change in the parentheses is normalized to the results of the standalone model.

power to the standalone model. Due to the significant impact of temperature on leakage power, the leakage estimation becomes more accurate. Thus, the IR-drop and temperature are also closer to the full-model results.

Partial-thermal model includes the thermal impact on leakage power as well as the PDN-power interactions. With these effects adding up, the results become smaller than the first four models: the IR-drop decreases by 11.28%, the temperature decreases by 7.95%, and the dynamic power decreases by 4.53% compared to the standalone model. However, for leakage, it almost drops by half because the actual temperature is lower than the reference temperature (100 °C) and in this temperature range, the leakage has an exponential relationship with temperature, resulting in severe errors. For example, the leakage current at 80 °C is 53.23% of that at 100 °C based on the leakage power model described in Section 5.3.3.

For the full-model, the temperature impact on wire resistivity is included (compared to partial-therm model). As a result, the IR-drop of the full-model becomes a little lower because of the lower PDN impedance at the simulated temperature (versus the reference temperature). Fig. 77 shows the thermal and IR-drop profiles of both dice using the full-model simulation. There is strong thermal and IR-drop coupling between the two dice due to the uniformly distributed TSVs.

A similar trend is found for the maximum transient power supply noise: the standalone model is 11.62%, 12.46%, 8.89%, 17.63%, and 18.68% higher than PDN-power, PDN-therm, PDN-therm-leak, partial-therm and full-model models, respectively. To understand the transient PDN noise, we plot the maximum noise over time, as shown in Fig. 78. Not only is the difference of maximum noise is large, so is the noise profile. The models that include the interaction between power dissipation and supply voltage predict a relatively faster damping profile (PDN-power, PDN-therm, partial-therm and full models). We define the damping rate as the amplitude of the second noise valley divided by that of the first noise valley. The relative difference of the damping rate between standalone and PDN-

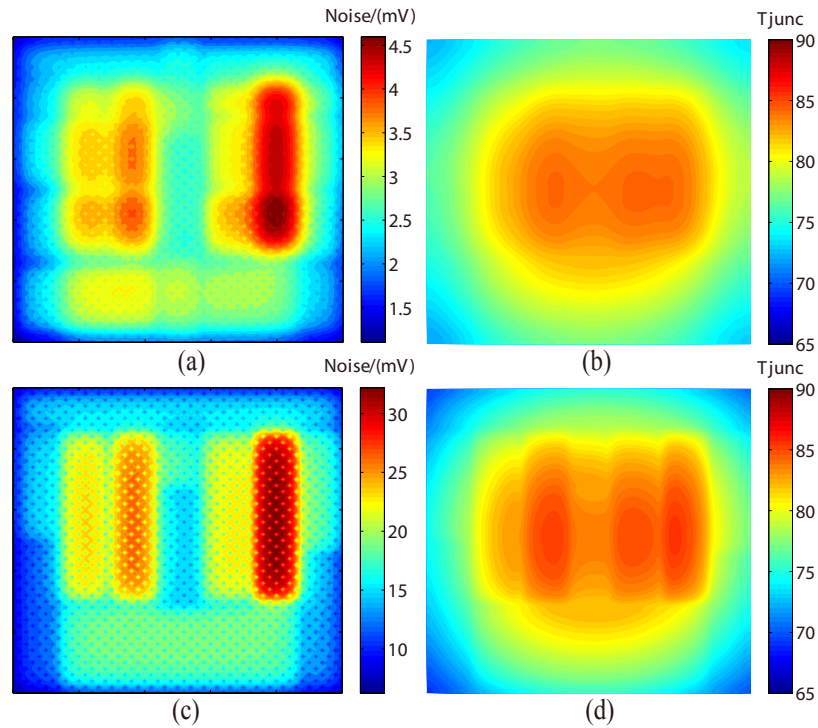


Figure 77: Steady state analysis result of full-model case (a) IR-drop of memory (4.71 mV) (b) Thermal of memory (84.35 °C) (c) IR-drop of processor (33.05 mV) (d) Thermal of processor (85.51 °C).

power, PDN-therm, partial-therm and full models is 20.03%, 20.10%, 18.86% and 18.89%, respectively. This difference results from the power-PDN negative loop that makes the noise damp faster.

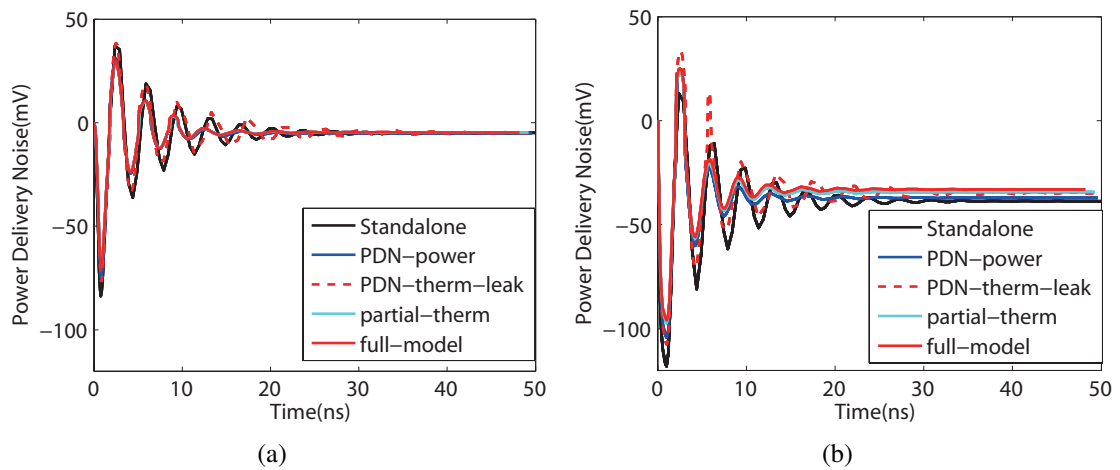


Figure 78: Transient power supply noise comparison (a) memory die (b) processor die. PDN-therm is very similar to PDN-power, thus omitted for better visualization.

In summary, thermal-leakage and PDN-power interactions have a significant impact on steady-state results, and the PDN-power interaction affects the transient PDN noise greatly. However, thermal-PDN interaction has a relatively smaller impact.

5.4.4 Accuracy Improvement Compared to Prior Work

Compared to the PDN-therm model, which only considers PDN-thermal interaction, the full-model achieves an accuracy improvement of 7.66%, 6.22% and 4.64% for IR-drop, transient PDN noise and maximum temperature, respectively; compared to PDN-therm-leak model, which includes both PDN-thermal and thermal-leakage interactions, the full-model achieves an accuracy improvement of 4.26%, 9.79% and 1.61% for IR drop, transient PDN noise and maximum temperature, respectively.

Table 15: Results after Different Number of Iterations

# iteration	IR drop (mV)	Temperature ($^{\circ}C$)	dynamic	leakage
1	34.79	92.86	56.90	9.21
2	33.33	86.78	57.00	7.74
3	33.10	85.71	57.01	7.51
4	33.06	85.54	57.02	7.48
5	33.06	85.52	57.02	7.47
6	33.05	85.51	57.02	7.47
7	33.05	85.51	57.02	7.47

Several leakage power estimation efforts have proposed a DC analysis framework similar to the partial-therm model [30][34]. Nevertheless, these two efforts focused on the leakage power. Moreover, a coarse thermal model was used by [34] to reduce the integration complexity, and as a result, the thermal map was not full-chip scale; only one iteration of the integrated analysis was performed in [30], since more iterations did not increase the estimation accuracy of leakage power significantly. However, we find one iteration is not adequate for obtaining accurate results due to the large change of leakage power in the temperature range between $85 \sim 100^{\circ}C$. Table. 15 shows the full-model results of processor die after several iterations. It is observed that at least 3 iterations are necessary to achieve

a relative error less than 1%.

5.5 Conclusions

In this chapter, we present a thermal and PDN co-simulation framework incorporating the interactions between temperature, PDN noise, and power dissipation. First, compared to prior work, the proposed models show that when we do not consider the interactions, there is a maximum error of 7.66%, 9.79% 4.64 % in IR-drop, transient noise, and temperature, respectively. Second, the integrated simulator is capable of performing fast simulations to answer what-if type questions in early design stages as well as being able to conduct detailed studies such as the impact of a wide range of technology parameters and different power delivery architectures. The modeling framework will benefit the architecture and packaging research communities.

CHAPTER 6

DIGITAL SIGNAL CHANNEL MODELING FOR 2.5-D AND 3-D INTEGRATION

In this chapter, we use repeater-based driver and receiver designs to model the digital signal channels for 2.5-D and 3-D integration platforms. The signaling latency, energy efficiency, and maximum bandwidth density of each integration platform are simulated and compared. In addition, the impact of process technology and I/O dimension scaling is studied. Last, we focus on the impact of temperature and investigate the thermal and electrical tradeoffs of die spacing in 2.5-D integration.

6.1 Circuit models of digital signal channels in 2.5-D and 3-D integration

The 2.5-D bridge-chip, interposer/HIST, and 3-D digital signal channels are illustrated in Fig. 79(a), Fig. 79(b) and Fig. 79(c), respectively. The signal channels consist of input/output (I/O) drivers and receivers, I/O pads, microbumps and chip-to-chip wires [85, 86, 87]. For 2.5-D integration, the chip-to-chip wires are horizontally routed on the interconnect carriers, such as bridge chip and interposer [10, 6, 9]. Note that for bridge-chip 2.5-D, there are additional vias and pads in the package to connect to the embedded bridge chip, as shown in Fig. 79(a). For 3-D integration, the chip-to-chip wires are vertical vias such as monolithic nanoscale vias or TSVs [8, 88, 89].

The equivalent circuit models for the 2.5-D and 3-D signal links are shown in Fig. 80(a) and 80(b), respectively. The parasitics of the pads, microbumps and wires are included. For 2.5-D integration, the chip-to-chip wires are modeled using three segment π -model per millimeter [85, 86]. For 3-D integration, the vias are modeled based on the compact models described in [90, 91, 87]. An optimal Signal-to-Ground (SG) TSV/microbump coupling case is considered. An ESD capacitor of 50 fF is added to both driver and receiver sides, a pre-driver of 102Ω is added prior to the driver and an output resistor of

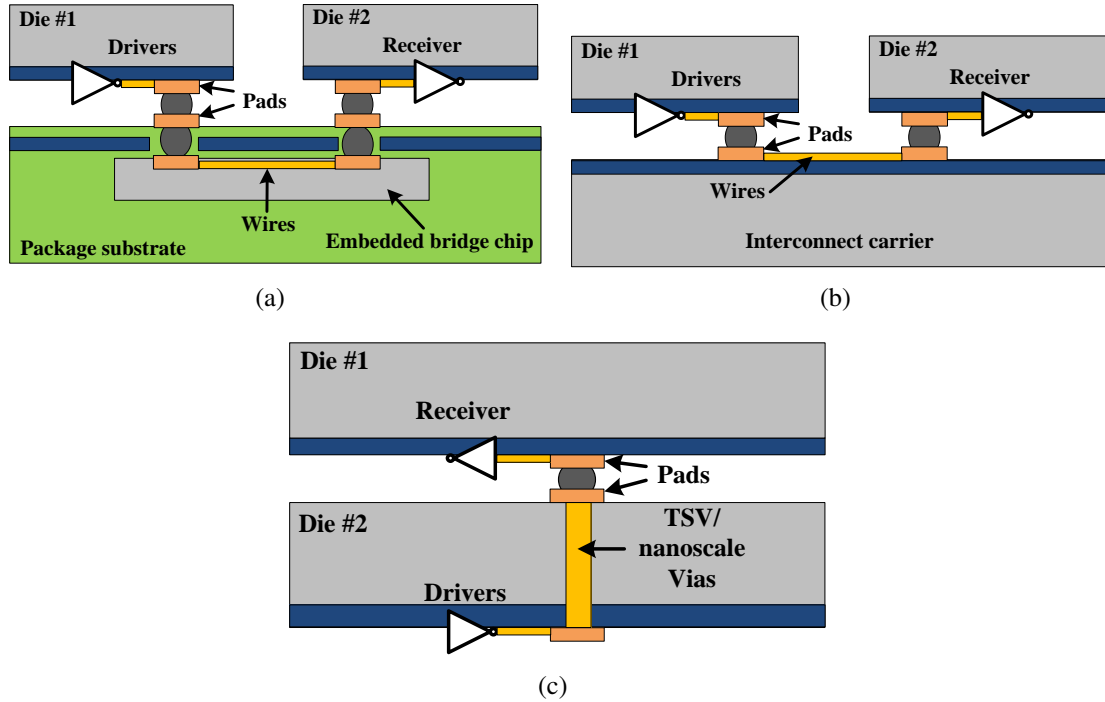


Figure 79: Illustration of digital signal channels (a) bridge-chip 2.5-D integration (b) interposer and HIST 2.5-D integration (c) 3-D integration

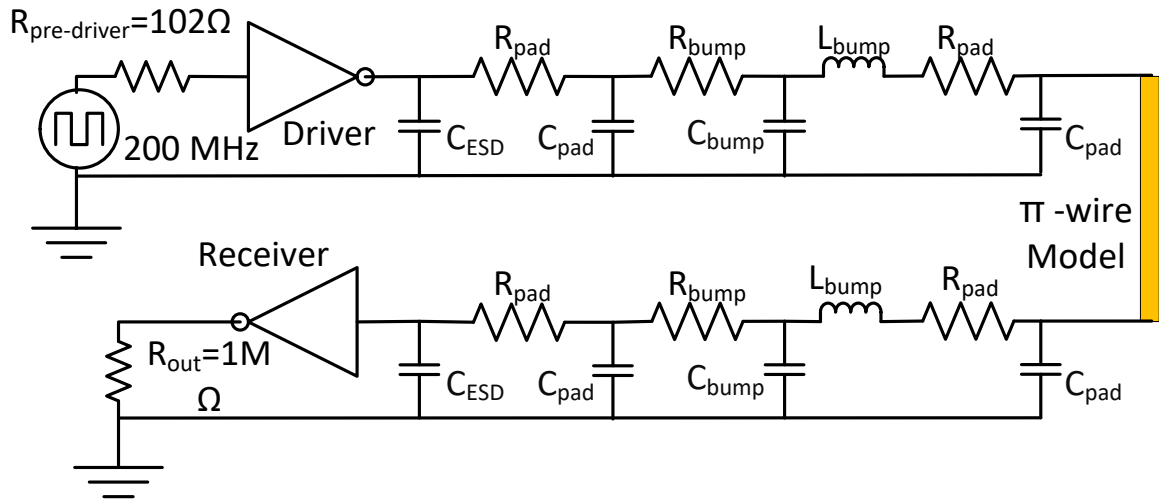
1 $M\Omega$ is added as termination impedance of the receiver output [85].

Table 16: Physical dimensions of each parameter of signaling models

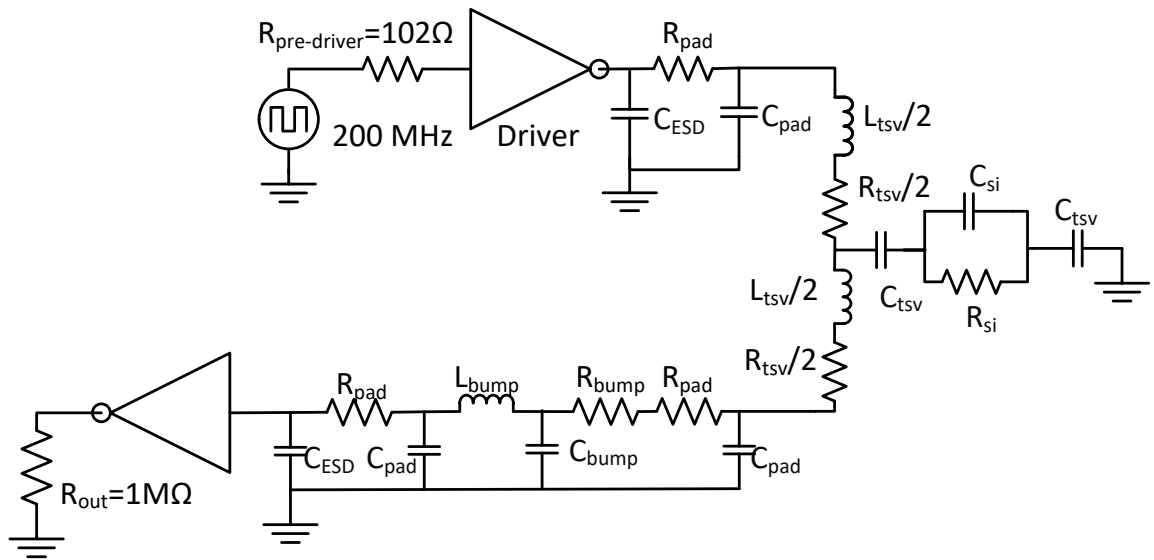
Parameter	value
Link wire length (mm)	0.1 ~ 5
Link wire pitch/thickness/width (μm)	1.6/2/0.8
Pad width (μm)	0.1 ~ 50
Microbump diameter & height	$0.8 \times$ Pad width
TSV diameter (μm)	0.1 ~ 10
TSV height	$15 \times$ TSV diameter
TSV dioxide thickness	$0.05 \times$ TSV diameter
Microbump/TSV pitch	$2 \times$ Microbump diameter

The physical dimensions of the I/Os and wires are summarized in Table 16. The wire specifications are based on the dimensions of the top global wires from NCSU FreePDK 45 nm [92]. The wire routing configuration is assumed to use a fan-in approach as demonstrated in [9, 93], therefore the wire pitch could be smaller than the microbump pitch.

The models and equations to estimate the parasitics of each parameter are summarized



(a)



(b)

Figure 80: Illustration of digital signal channels (a) 2.5-D integration (b) 3-D integration

Table 17: Equation for parasitics estimation

Component	Equation
Wire to substrate capacitance: C_{wire} [85]	$C_{wire} = \epsilon_0 \cdot \epsilon_{ox} \cdot \frac{W \times L}{t_{ox}}$ <p>W and L are wire width and length.</p>
Wire to wire coupling capacitance: C_{wire} [85]	$C_{wire} = \epsilon_0 \cdot \epsilon_{ox} \cdot \frac{t_{ox} \times L}{P_{wire}}$ <p>P_{wire} is the wire pitch</p>
Pad capacitance: C_{pad} [90]	$C_{pad} = \epsilon_0 \cdot \epsilon_{ox} \cdot \frac{W_p^2}{t_{ox}}$ <p>W_p is the pad width and t_{ox} is the ILD thickness</p>
Microbump to ground capacitance: C_{bump} [90, 95]	$C_{bump} = \epsilon_0 \cdot \epsilon_{underfill} \cdot \frac{2 \cdot \pi \cdot H_{bump}}{\text{arcosh}(\frac{P_{bump}}{D_{bump}})}$ <p>P_{bump} is the microbump pitch</p>
Microbump inductance: L_{bump} [90, 87]	$L_{bump} = \frac{\mu_0 \cdot \mu_{bump}}{2 \cdot \pi} \times H_{bump} \times \ln(\frac{2 \cdot P_{bump}}{D_{bump}})$
TSV dioxide capacitance: C_{ox} [90]	$C_{ox} = \epsilon_0 \cdot \epsilon_{ox} \cdot \frac{2 \cdot \pi \cdot H_{tsv}}{\ln(\frac{D_{tsv}}{D_{tsv} - 2 \cdot t_{ox}})}$ <p>t_{ox} is the dioxide thickness, H_{tsv} is the TSV height</p>
TSV depletion capacitance: C_{dep} [94]	$C_{ox} = \epsilon_0 \cdot \epsilon_{Si} \cdot \frac{2 \cdot \pi \cdot H_{tsv}}{\ln(\frac{D_{tsv} + W_{dep}}{D_{tsv}})}$ <p>W_{dep} is the depletion width</p>
TSV total capacitance: C_{tsv}	$C_{tsv} = \frac{C_{ox} \cdot C_{dep}}{C_{ox} + C_{dep}}$
TSV inductance: L_{tsv} [90]	$L_{tsv} = \frac{\mu_0 \cdot \mu_{tsv}}{2 \cdot \pi} \times H_{tsv} \times \ln(\frac{2 \cdot P_{tsv}}{D_{tsv}})$
TSV to substrate capacitance: C_{Si} [94]	$C_{ox} = \epsilon_0 \cdot \epsilon_{Si} \cdot \frac{2 \cdot \pi \cdot H_{tsv}}{\text{arcosh}(\frac{P_{tsv}}{D_{tsv}})}$ <p>t_{ox} is the dioxide thickness, H_{tsv} is the TSV height</p>
TSV to substrate resistance: R_{Si} [94]	$R_{Si} = \frac{\epsilon_0 \cdot \epsilon_{Si}}{C_{Si} \cdot \sigma_{Si}}$ <p>σ_{Si} is the silicon conductance</p>

and microbumps are assumed to be cylinder for simplicity.

For the driver and receiver designs, multiple driver stages with a constant fan-out of 4 (FO4) between stages were selected [87]. The minimum sized inverter driving four identical inverters is tuned to achieve equal rise and fall times. Energy-delay-product (EDP) was used to optimize the number of driver stages [96], which ranges from 1 to 5 stages in all the simulations. We anticipate the signal channels are used in applications similar to Wide I/O spec [62] and therefore, a low-frequency digital signal input of 200 MHz are used [85].

6.2 Comparison of integration platform latency, energy efficiency and bandwidth density

In this section, we develop the circuit models in *HSPICE* netlists, and simulate the 50%-to-50% propagation delay and total energy of the signal channels for all the 2.5-D and 3-D integration scenarios: bridge-chip-, interposer- and HIST- based 2.5-D and TSV- and monolithic-based 3-D. The device models are based on *ASU PTM* 45 nm HP library [84]. The version of *HSPICE* is 2017.03sp1, and the BSIM model for 45 nm library is level 54 version 4.5.

The I/O and wire dimensions used in the evaluation along with the latency and energy simulation results are summarized in Table 18. The interconnects of HIST consist of small microbumps between dice and the bridge-chip, therefore the dimensions are much smaller compared to other 2.5-D scenarios. For 2.5-D integration, we assume the chip-to-chip wire is 1 mm long [9] and for 3-D cases, the interconnect is a vertical via, which is 75 μm and 800 nm long for TSV- and monolithic-based 3-D IC cases, respectively.

From Table 18, we find the HIST shows better electrical performance than the bridge-chip and interposer because the pads and microbumps of HIST are much smaller than those of the bridge-chip and interposer, which results in a smaller capacitance. The total capacitance of a microbump and a pair of the pads for HIST are approximately 18X and 4X smaller compared to the bridge-chip and interposer cases, respectively. As a result, HIST

Table 18: Comparison of different integration platforms

	Bridge-chip [9, 93]	Interposer [6]	HIST [10]	TSV-3D [97, 98, 87]	Monolithic-3D [89, 99]
Microbump (pitch/diameter/height)	50/25/50 μm^\dagger	24/12/12 μm^*	8/4/4 μm	40/20/20 μm^*	0.4/0.2/0.2 μm
Pad size	37.5 μm	15 μm	5 μm	30 μm	0.3 μm
Link wire length	1 mm	1 mm	1 mm	75 μm	800 nm
Microbump capacitance	14.75 fF	3.54 fF	1.18 fF	5.90 fF	0.06 fF
Pad capacitance	9.72 fF	1.56 fF	0.17 fF	6.20 fF	~ 0 fF
Link wire capacitance	118.9 fF	118.9 fF	118.9 fF	31.5 fF	0.1 fF
Link latency with ESD	125.1 ps	118.6 ps	117.3 ps	98.9 ps	94.5 ps
Link energy with ESD	306.3 fJ/bit	267.4 fJ/bit	259.9 fJ/bit	176.2 fJ/bit	135.1 fJ/bit
Link latency without ESD	107.9 ps	101.0 ps	99.7 ps	75.9 ps	33.0 ps
Link energy without ESD	206.0 fJ/bit	167.0 fJ/bit	159.5 fJ/bit	76.2 fJ/bit	3.7 fJ/bit

[†] For bridge-chip case, the microbump is extended into the package substrate and has a higher value than its diameter.

* Although smaller sized microbump is reported in [97], we use an average value for these two scenarios.

achieves a 6.2% and 1.1% reduction in latency and approximately 15.1% and 2.8% in link energy compared to bridge-chip and interposer cases, respectively.

The performance and energy efficiency trends can be better understood from Fig. 81, where we study the impact of pad size on latency and link energy. A key difference between conventional 2.5-D (bridge-chip and interposer case) and HIST is the pad size. Therefore, as we further scale the pad size to be comparable to on-chip dimensions, the performance of the interconnect will approach that of on-chip interconnects [86]. However, we also notice a diminishing return when the pad size is scaled below $5 \mu m$. When the pad size is reduced from $100 \mu m$ to $5 \mu m$, the link energy decreases by approximately 18.0% and the delay reduces by 7.81%; further decreasing the size of the pad will not bring significant benefits. To take full advantage of HIST and to prevent fabrication complexity, we need to identify the optimal pad size. For a 1 mm long wire, $5 \mu m$ is used.

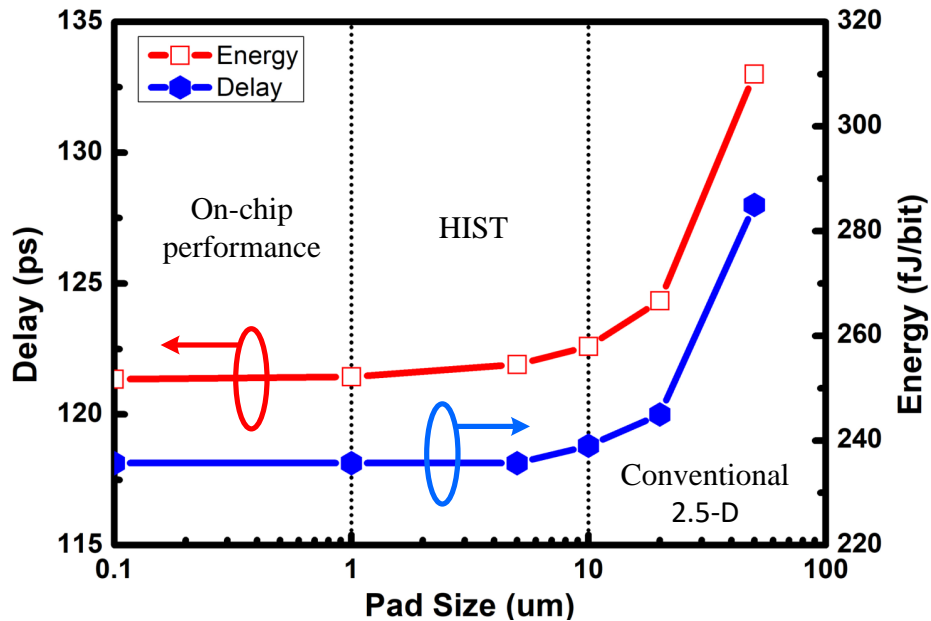


Figure 81: Impact of pad scaling on electrical performance of signal channel.

Due to the significantly shorter die-to-die wires compared to 2.5-D integration, 3-D integration exhibits smaller link latency and energy than 2.5-D designs. For the TSV 3-D case, even with ESD capacitors, the TSV-based 3-D design achieves approximately 15.7% smaller latency and 32.2% smaller energy than those of HIST. For the monolithic 3-D case,

due to the utilization of nanoscale vertical vias, the link latency and energy are approximate 19.4% and 48.0% smaller compared to HIST, respectively. Because of the on-chip like performance of monolithic 3-D [100, 101], ESD may not be necessary, and in which case, the latency and energy will be 71.8% and 98.6% smaller than HIST, respectively. However, 3-D designs generally require vertical vias and die-to-die bonding which poses additional complexity and challenges in fabrication and manufacturing.

Next, based on the above parameter specifications, we simulate the maximum data rate per link (F_{max}) using $6 \cdot \tau$ (τ is latency and $F_{max} > \frac{1}{6 \cdot \tau}$) settlement [86]. The bandwidth per millimeter (bandwidth density, BWD) is then calculated using the following equation:

$$BWD = \frac{1}{P_{bump}} \times F_{max} \text{ (Gbps/mm)} \quad (6.1)$$

where P_{bump} is the bump pitch with units of millimeter. We assume two rows of staggered bumps, of which half are for signals and the rest are for ground [86, 97, 93].

The simulated BWD of all signal channels are summarized in Table 19. Likewise, HIST achieves the largest BWD among all the 2.5-D solutions. Due to the extremely high I/O density along with low latency, the BWD of monolithic 3-D integration reaches 12.6 Tbps/mm. The high BWD demonstrates the potential of using monolithic 3-D to solve the communication bottlenecks of integrating many chips in a single package [101].

Table 19: Bandwidth density of each integration platform

	Bridge-chip	Interposer	HIST	TSV-3D	Monolithic-3D (no ESD)
I/Os per mm	20	41	125	25	2500
Max data rate (GHz)	1.38	1.41	1.41	1.69	5.05
BWD (Gbps/mm)	27.6	56.4	176.25	42.25	12,625

6.3 Impact of technology parameter scaling

In this section, we explore the impact of two aspects of technology scaling on 2.5-D and 3-D integrated systems. First, the transistors are fabricated using advanced processes for

improved device performance and energy efficiency; second, the length of chip-to-chip wires is also reduced.

6.3.1 Technology process scaling

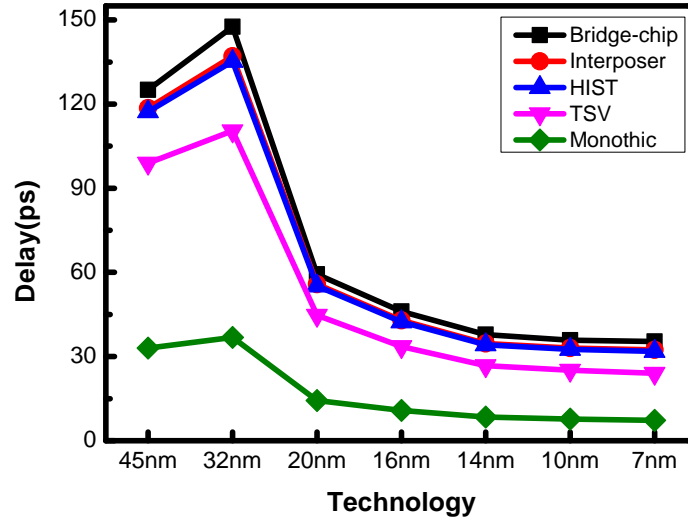
Device currents have been enhanced dramatically by FinFET technology beyond 32 nm. Moreover, energy efficiency is also improved through better channel engineering by using 3-D surround gates and supply voltage scaling [102, 103, 104, 105, 11]. Table 20 summarizes the device information of each process of the *PTM* libraries, of which the multi-gate (MG) devices use *HSPICE* model level 72 version 110. From 20 nm to 7 nm, the devices are built using FinFET technology and the transistor equivalent width is calculated as follows,

$$W = 2 \times Fin_{Height} + Fin_{Width} \quad (6.2)$$

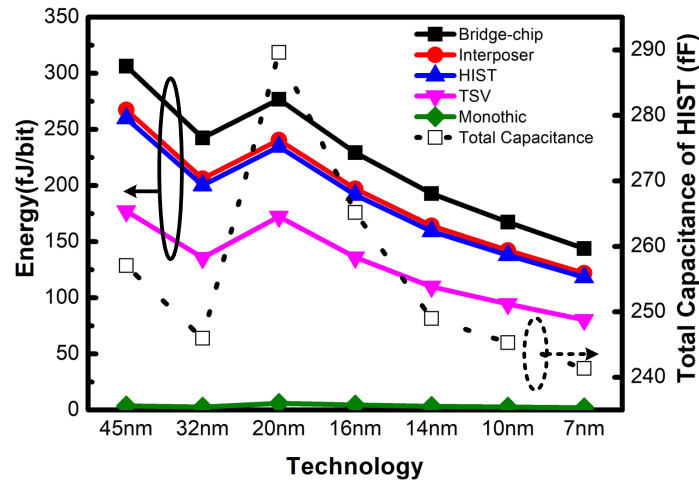
Table 20: Channel length, minimum inverter size, and supply voltage of each process technology using *PTM* device library

		45 nm	32 nm	20 nm	16 nm	14 nm	10 nm	7 nm
Device type		planar	planar	MG	MG	MG	MG	MG
Min channel length (nm)		50	36	24	20	18	14	11
Min inverter size (nm)	NMOS	50	36	71	64	56	51	43
	PMOS	97	66	79	68	57.5	51	43.5
VDD (V)		1.0	0.9	0.9	0.85	0.8	0.75	0.7

The delay and energy of signal channels as a function of process technology are plotted in Fig. 82. Both metrics decrease as process technology advances except at 32 nm, where the impact of VDD scaling (0.9 V at 32 nm against 1.0 V at 45 nm) dominates and results in an increased delay. Nevertheless, even if the VDD is scaled down with technology, we still observe a reduction in delay between the process technology due to device performance enhancement. In addition, for energy analysis, we plot the total capacitance in Fig. 82(b) (dotted line plotted to the right Y-axis) which shows a similar decreasing trend as energy, thus such energy reduction is not only due to VDD scaling.



(a)



(b)

Figure 82: Delay and energy of the signal channels implemented by different technology nodes.

Another important observation is that the relative difference in the electrical metrics of different 2.5-D integration approaches become larger, as shown in Fig. 83. This is because the total capacitance reduces with technology scaling (as shown in the dotted line of Fig. 82(b)) while the capacitance of the I/Os remain relatively constant [11]. Therefore, this shows the significance of scaling I/O dimensions for those multi-die systems using advanced technology process such as *Intel Stratix 10* (14 nm) [65]. Otherwise, if the I/O dimensions are not scaled properly with device technology, there may be minimal benefits

using advanced device technologies. For example, the bridge-chip-based 2.5-D platform at 16 nm is comparable in power efficiency (229.4 fJ/bit) to the HIST-based platform at 20 nm (234.6 fJ/bit).

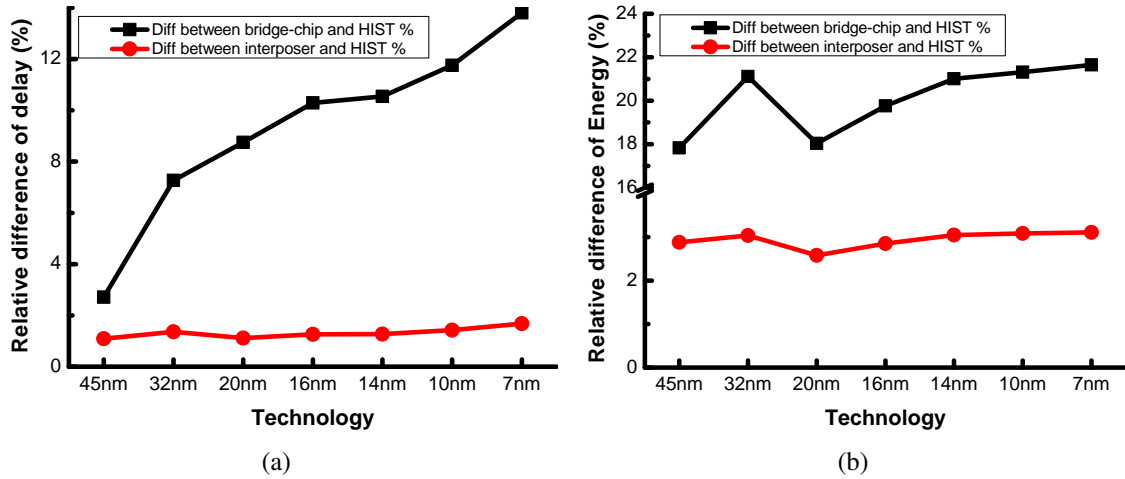
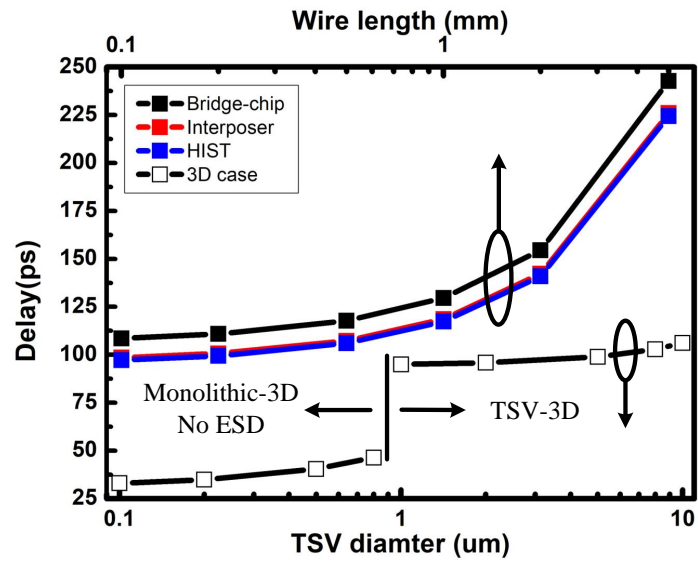


Figure 83: Relative difference between bridge-chip/interposer and HIST vs. device process technology (a) delay (b) energy.

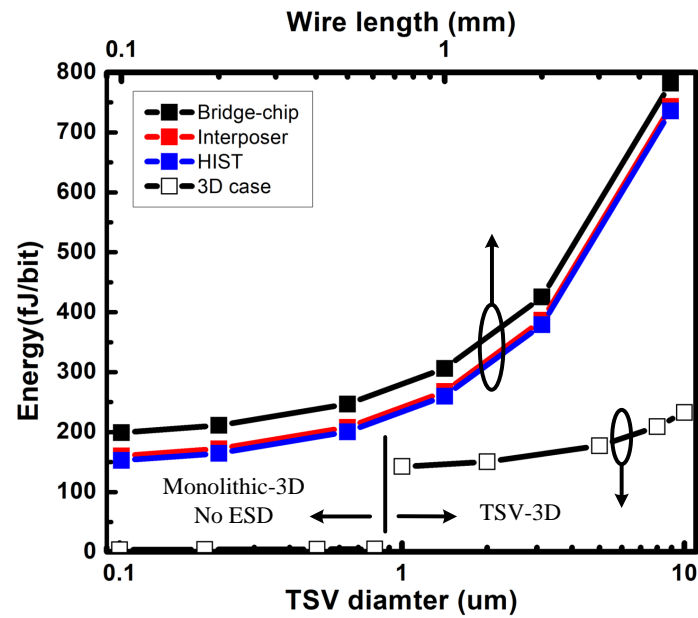
6.3.2 Impact of interconnect wire length in signal channels

The delay and energy of signal channels with varying wire lengths and via diameters are shown in Fig. 84. For 2.5-D integration platforms, the wire length varies from $100 \mu m$ to $5 mm$ and for 3-D cases, the via diameter varies from $100 nm$ to $10 \mu m$ with a fixed aspect ratio of 15. Note that for monolithic 3-D case, we do not include the ESD capacitor due to its monolithic fabrication process, therefore it attains ultra-low delay and energy. Likewise, for both 2.5-D and 3-D cases, with a scaled interconnect wire length, the electrical performance is improved in both delay and energy. Moreover, when the 2.5-D wire length is close to $100 \mu m$, the delay and energy of HIST 2.5-D is comparable to the TSV-3D case even with a TSV diameter of $1 \mu m$ because the parasitics of the channel are close and are dominated by ESD capacitors.

Similar to the technology scaling, when the wire is short ($< 500 \mu m$), the relative



(a)



(b)

Figure 84: The impact of interconnect scaling on (a) delay (b) energy

difference between bridge-chip/interposer and HIST cases becomes larger, as shown in Fig. 85. This is because as the wire length is reduced, the parasitics of microbumps and pads become comparable to that of the wires. For example, the relative difference in energy between HIST and bridge-chip is as high as 30.5% when using a wire length of 0.1 mm

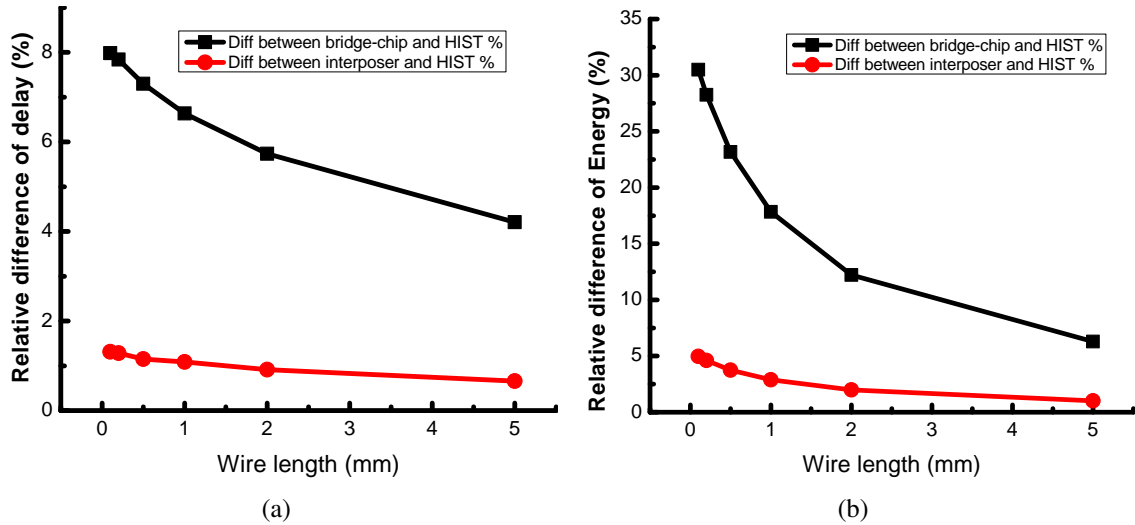


Figure 85: Relative difference between bridge-chip/interposer and HIST vs. wire lengths (a) delay (b) energy

but is only 6.3% with a wire length of 5 mm. Likewise, this shows the necessity of using smaller I/Os (HIST) especially when the interconnect wire is scaled down and the systems are built more tightly. On the other hand, as longer wires significantly lower the electrical performance, it is less appealing to design such channels beyond 2 mm as the delay and energy degrades very rapidly.

6.4 Impact of temperature on signaling in 2.5-D and 3-D integration

Higher temperature lowers the carrier mobility and threshold voltage (leads to higher leakage current), thus it may result in a larger delay and higher energy [106]. As previously shown in Chapter 3, there are thermal challenges for 2.5-D and 3-D integration and the integrated dice may experience higher temperatures than single-die packages. Therefore, without careful consideration of the thermal impact, the comparison between them may be incomplete. In this section, we focus on HIST-based 2.5-D and TSV-based 3-D and discuss the impact of temperature on signaling.

Fig. 86(a) and 86(b) show the thermal impact on delay and energy for HIST-based

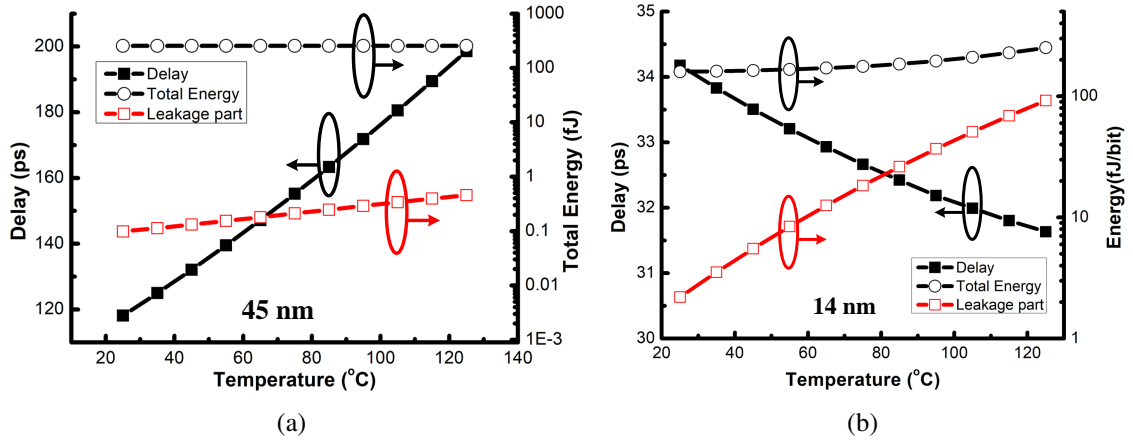


Figure 86: Impact of temperature on delay and energy of HIST-based 2.5-D integration (a) 45 nm (b) 14 nm

2.5-D with 45 nm and 14 nm technologies, respectively. Using 45 nm, as the temperature increases, the delay becomes larger due to device performance degradation. However, the total energy is not significantly impacted due to the leakage power being only a small portion of the total power. Using 14 nm technology, which is FinFET-based device, the trend is different due to the temperature inversion effect [104, 107, 108, 109, 110], in which the rate of mobility degradation is less than that of threshold voltage. Therefore the delay decreases by about 7.4% from 34.2 ps to 31.6 ps. Due to the extremely short channel at 14 nm, the leakage power is larger than that of the 45 nm devices and impacts the total energy dramatically. The total energy changes by approximately 57.9% from 25 °C to 125 °C.

6.4.1 2.5-D and 3-D signaling comparison revisit with the impact of temperature

Based on Chapter 3 Table 6, the temperature was set to 98 °C for the 2.5-D case and 125 °C for the 3-D case and their electrical performance was compared for 45 nm and 14 nm process technologies, respectively. The results using the 45 nm library are shown in Fig. 87(a) and Fig. 87(b) while the results using 14 nm library are shown in Fig. 87(c) and Fig. 87(d).

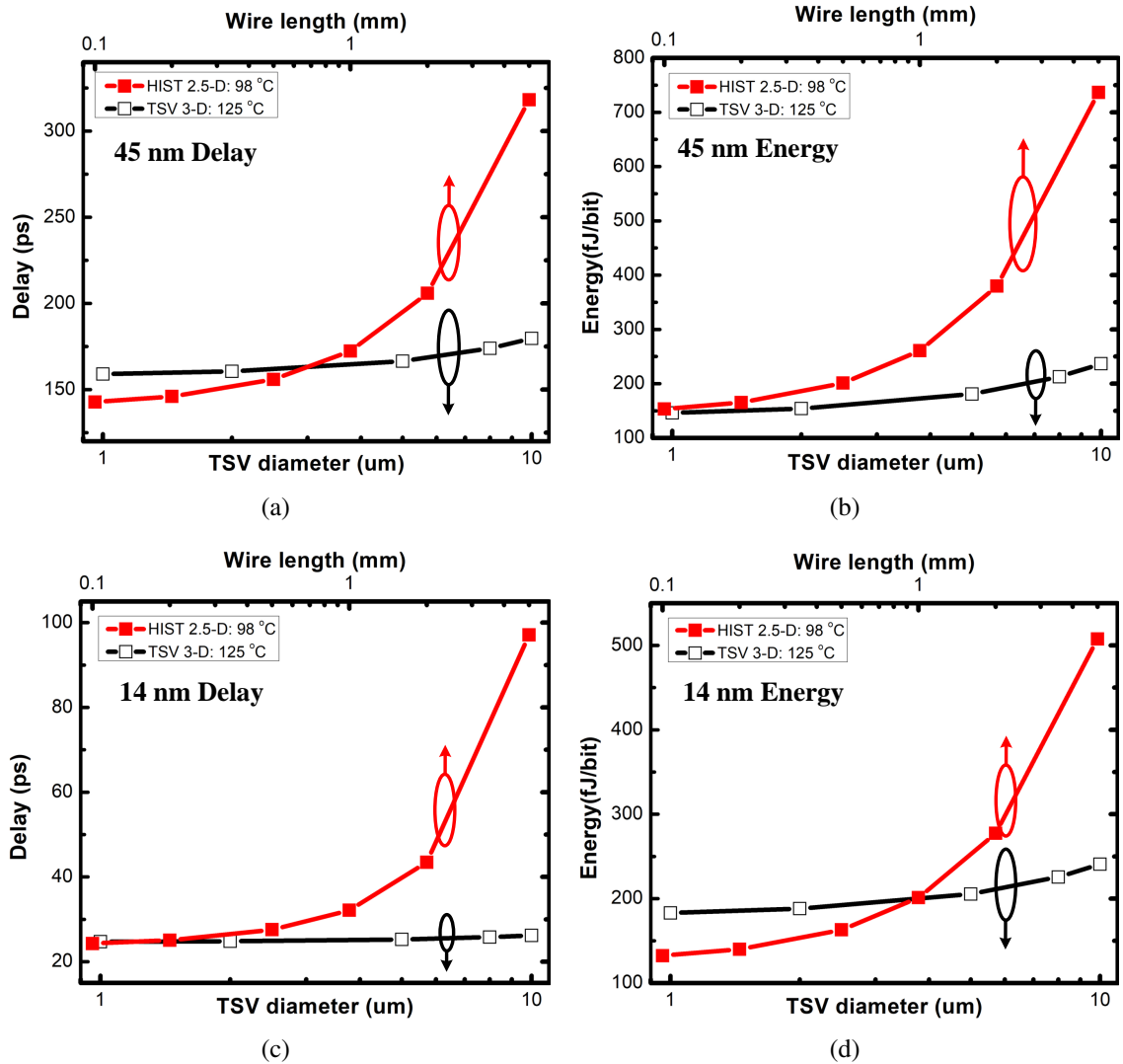


Figure 87: Comparison of HIST-based 2.5-D and TSV-based 3-D integration (a) delay using 45 nm library (b) energy using 45 nm library (c) delay using 14 nm library (d) energy using 14 nm library

For 45 nm, 2.5-D shows better delay than 3-D case when the wire length is below 0.7 mm. As a comparison to Fig. 84(a), we observe TSV-3D is better than HIST 2.5-D using any wire length when not considering thermal impact. Therefore, for TSV 3-D integration, to compensate for the negative impact of elevated temperature on delay, smaller and shorter TSVs must be used. From Fig. 87(b), TSV 3-D is better than HIST 2.5-D in energy because for 45 nm, leakage power is relatively a small portion which results in the total energy being almost independent of temperature.

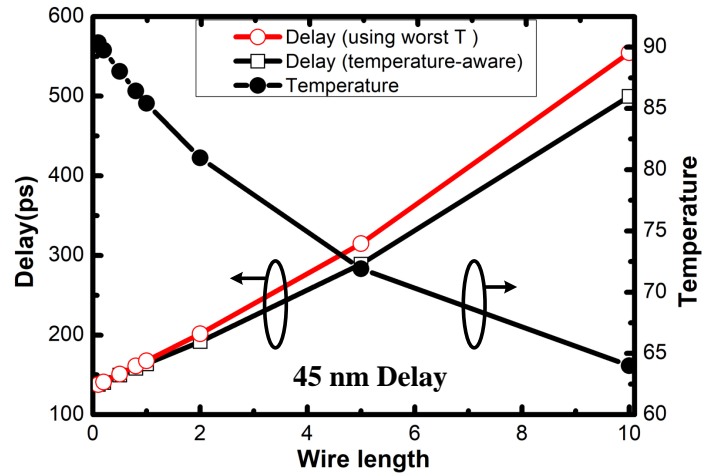
For 14 nm, due to the temperature inversion effect, the observation and conclusion are different. In terms of delay, 3-D shows superior performance compared to 2.5-D, which is similar to the trend shown in Fig. 84(a). However, HIST 2.5-D shows power benefits over TSV 3-D case when the wire length is below 1 mm. This result can be explained from Fig. 86(b), where leakage power percentage increases dramatically from 80°C to 125°C .

In summary, HIST 2.5-D is capable of exhibiting better speed and power efficiency for 45 nm and 14 nm (with a moderate wire length) than TSV 3-D case, respectively. Without carefully considering the thermal impact and the fact that 3-D integration leads to an elevated temperature, the evaluation of the two types of integration may not be accurate.

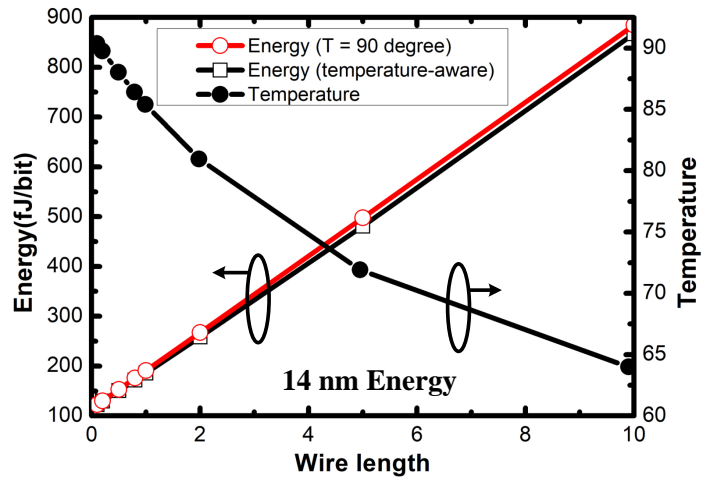
6.4.2 Thermal and electrical tradeoffs on die spacing in 2.5-D integration

Based on the discussions in Chapter 3 Section 3.6.4, there are thermal and electrical tradeoffs for die spacing. As the spacing between dice increases, the junction temperature of both dice decreases. Moreover, the rate of temperature reduction is significantly larger for the low-power die than that for the high power die. On the other hand, as the die spacing increases, the communication latency increases and power efficiency decreases because the wires have larger parasitics. By taking the temperature into consideration, we focus on the low power die (memory die) and investigate whether or not the temperature benefits obtained from larger spacing could compensate the electrical penalty in both 45 nm (focused on delay), and 14 nm (focused on energy) technologies.

The results are shown in Fig. 88. For both libraries, we consider a baseline case under a constant worst temperature of 90°C for the digital signal channels. Although we observe a trend of increased benefits from the lower temperature (by increasing the die spacing) for the delay at 45 nm and energy at 14 nm, respectively, the effect is small relative to the effect of longer wire length on delay and energy efficiency. Moreover, for the 14 nm library, the leakage power change is not large (less than 15%) when the temperature is reduced from 90°C to 64°C , and thus, larger spacing brings minimal electrical benefits. In summary,



(a)



(b)

Figure 88: Thermal and electrical tradeoffs for (a) delay using 45 nm library (b) energy using 14 nm library

even though larger spacing would lower the system temperature, especially for the lower power dice, the delay and energy penalty is still larger than the thermal benefits and it is not optimal to sacrifice the communication performance for temperature benefits.

6.5 Summary

In this Chapter, we develop circuit models to benchmark digital communication channels of 2.5-D and 3-D integration platforms. The delay, energy per bit and bandwidth density

of each integration platform are compared. HIST-based 2.5-D integration shows better performance compared to other 2.5-D approaches because of the ultra low parasitics of microbumps and pads. Monolithic-based 3-D ICs using nanoscale vias show superior performance compared to conventional TSV-based 3-D ICs due to the small dimensions of the vias.

Moreover, the thermal impact on 2.5-D and 3-D integration is investigated. Because the TSV 3-D case exhibits a higher temperature, HIST 2.5-D is capable of attaining better speed and power efficiency for 45 nm and 14 nm (with a moderate wire length) than TSV 3-D case, respectively. In addition, the thermal and electrical tradeoffs of die spacing in 2.5-D integration are investigated. Increasing die spacing lowers system temperature, but the thermal benefits may not compensate delay and energy penalties.

CHAPTER 7

CONCLUSION AND FUTURE DIRECTIONS

2.5-D and 3-D integration are emerging technologies with the potential to offer significant benefits in communication bandwidth, footprint reduction, power efficiency and heterogeneous functionalities to keep up with the rapidly evolving requirements of machine learning, cloud computing, IoT applications and Moore's Law. In this thesis, the thermal and power delivery challenges in 2.5-D and 3-D integration have been presented. Thermal and PDN modeling frameworks are developed to enable the design space exploration and optimize the temperature and power delivery of microelectronic systems. Research efforts have been made through the following projects:

1. Thermal isolation technologies are explored and developed to address the thermal coupling issues in heterogeneous 3-D ICs.
2. Thermal evaluation and benchmarking are performed focusing on bridge-chip 2.5-D integration.
3. PDN benchmarking is conducted to understand and address the challenges arising for bridge-chip based 2.5-D integration
4. A thermal-PDN co-simulation framework is built to accurately model temperature, supply voltage and power dissipation along with the interactions between them.
5. Digital signal channels of 2.5-D and 3-D integration are evaluated and benchmarked.

In this Chapter, we conclude our presented work and summarize the contribution of the above projects. Several future directions and potential works will be briefly discussed.

7.1 Summary of the presented work

This thesis has five major parts, of which the conclusions and contributions are summarized as follows:

First, a novel heterogeneous 3-D integration architectures with interposer embedded microfluidic cooling, air-gap isolation and extended heat spreader for the isolated dice is proposed as a potential technology for thermal decoupling issues in heterogeneous 3D ICs. In the evaluated memory-processor stack with air-gap isolation, the memory temperature is reduced by $32.79\text{ }^{\circ}\text{C}$ compared to conventional bonding with underfill. To maintain the thermal benefits of an air gap, the TSVs should be carefully designed by taking the thermal effects of their number, diameter and layout into account.

Second, a comprehensive thermal study for 2.5-D integration focusing on bridge-chip based technology is performed to identify the thermal limits and challenges in such integration approaches. A CPU-FPGA-DRAM assembly is used as an application example. Bridge-chip 2.5-D integration is compared to interposer and non-embedded bridge-chip 2.5-D integration. Compared to bridge-chip 2.5-D integration, interposer 2.5-D integration offers a modest improvements in terms of maximum die junction temperature due to better heat spreading in the interposer layer. Bridge-chip 2.5-D integration is also compared to TSV and monolithic 3-D integrations and shows improved thermal response due to smaller power density. Through parametric study of bridge-chip based 2.5-D, the impact of die thickness mismatch and die spacing are investigated in 2.5-D systems. We conclude that the die dissipating the largest power should be the thickest in a multi-die package. Larger lateral spacing between dice reduces the temperature at the unrealistic expense of communication power efficiency. Therefore, it is necessary to consider this tradeoff when selecting appropriate die spacing.

Third, a PDN modeling framework for emerging heterogeneous 2.5-D integration platforms is presented. Validation using IBM power grid benchmarks shows the IR-drop and

transient analysis have a maximum relative error of less than 7.29% and 0.67%, respectively. Next, the framework is used to evaluate interposer and bridge-chip based 2.5-D integration platforms. The simulation results show that interposer based 2.5-D integration may exhibit a worse power supply noise due to TSV parasitics. In bridge-chip based 2.5-D integration, under the assumption that the bridge chips underneath the active dice block access to package power/ground planes, some power delivery challenges are highlighted. Minimizing the overlap region between bridge chip and active dice, using multiple bridge chips instead of a single large bridge and inserting through bridge chip vias help to mitigate PSN.

Fourth, a thermal and power delivery network (PDN) co-simulation framework is presented for single-die and emerging multi-die configurations that incorporates the interactions between temperature, supply voltage, and power dissipation. The temperature dependencies of wire resistivity and leakage power are considered and the supply voltage dependencies of power dissipation are modeled. Starting with a reference power dissipation, the framework is capable of evaluating the temperature distribution and PDN noise simultaneously and eventually updating the power dissipation based on the thermal and supply voltage distributions. The simulation results of an example two-tier 3-D stack show that prior models considering only part of the interactions between temperature, supply voltage and power dissipation have a maximum error of 7.66%, 9.79%, 4.64 % in IR-drop, transient power supply noise, and temperature, respectively.

Fifth, we present signaling evaluation framework to benchmark the communication links of 2.5-D and 3-D integration platforms. The impact of technology scaling, pad size, and interconnect length are shown. HIST platforms show significant latency and power efficiency improvement compared to bridge-chip and interposer based platforms. Moreover, thermal impact on signaling is discussed for 2.5-D and 3-D integration.

7.2 Summary of the future directions

There are several opportunities and directions to extend and advance the work of this thesis.

First, one of the advantages of the thermal and PDN modeling framework is the flexibility to add more advanced algorithms [111]. The models developed for this thesis have not yet been optimized for simulation speed and memory efficiency. However, as more details and parameters are added into the models, the problem size will increase dramatically, and the advanced computational algorithms may become necessary to implement. Spectral graph sparsification-based PCG algorithm [112] and machine learning based acceleration techniques [113] have been preliminarily explored, which achieved over 10X speedups. We can further explore those algorithms and implement the thermal and PDN models for very large-scale circuits (>1 billion unknowns).

Second, based on the presented thermal and PDN co-simulation framework, we are able to perform thermal and PDN evaluation for more emerging technologies such as fan-out wafer-level packaging (FOWLP), as shown in Fig. 89. FOWLP routes the power/ground I/Os for the top die in the periphery therefore, the power delivery path for the top die is longer and contains larger parasitics. On the other hand, since FOWLP exploits mold material, which is a poor thermal conductor, for redistribution layer formation, there will be thermal challenges for such 3-D stack.

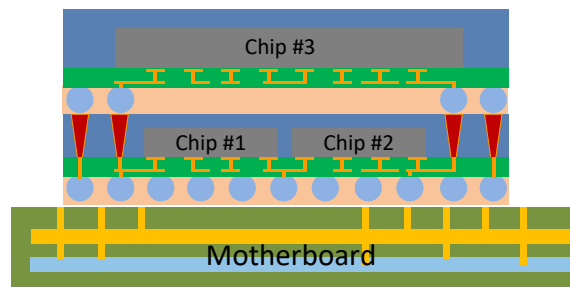


Figure 89: A 3-D chip stack using FOWLP technology.

Third, the PDN modeling framework could be extended with a more accurate package-level model considering the packaging layout and metal layer geometry. As the package

keeps shrinking and becomes comparable to the die in FOWLP technology, the coupling between package and on-die wires becomes non-negligible and it is necessary to accurately extract package parasitics considering the packaging routing, decap placement, vias and etc. To accomplish this work, a package-chip co-design framework has to be developed and package parasitics extraction flow needs to be established.

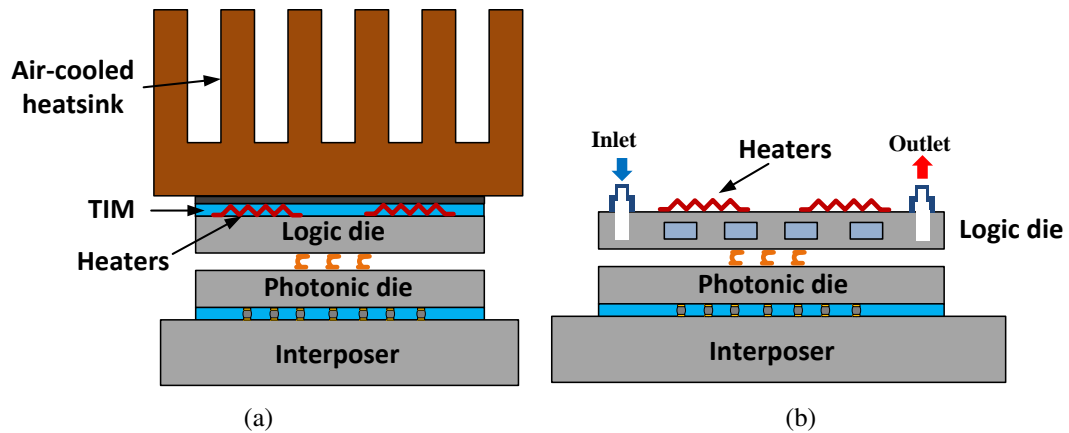
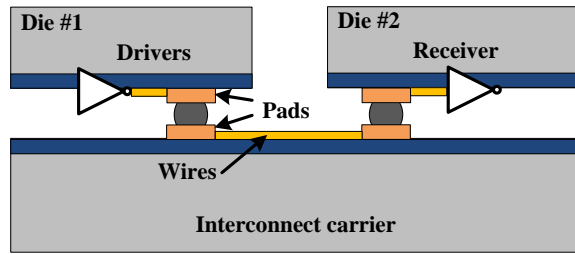


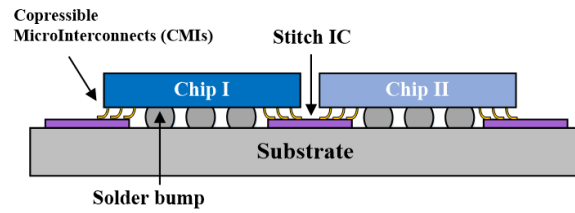
Figure 90: Two architectures for nanophotonics-based systems with thermal isolation technologies (a) using air-cooled heat sink (b) using microfluidic-cooled heat sink.

Fourth, we can explore the thermal impact on heterogeneous integration systems of CMOS chip with nanophotonics. Microring resonators are temperature sensitive and complex thermal compensation circuits are usually required for reliable optical signal propagation [114]. Next, we will investigate the thermal benefits of using our proposed thermal isolation technologies, as shown in Fig. 90. Due to the high thermal resistance of air isolation layer between processor and nanophotonics, the nanophotonics die will experience a much lower temperature variation even when the processor die experiences a high-duty switching activities. It is important to understand how this system takes advantage of thermal isolation technologies and quantitatively analyze the electrical benefits.

Last, based on the signaling evaluation framework, we could design test circuitries to demonstrate the electrical benefits of HIST system and compare with conventional 2.5-D integrated systems. The I/O parts of each integrated die will be emulated by CMOS chips



(a)



(b)

Figure 91: Demonstrated HIST platforms with active chips (a) two active chips emulating driver and receiver circuitries (b) emulated HIST system

and these chips will be assembled in a test carrier substrate, as shown in Fig. 91.

REFERENCES

- [1] Altera, *Enabling Next-Generation Platforms Using Altera's 3D System-in-Package Technology*.
- [2] J. Jeddelloh and B. Keeth, "Hybrid memory cube new dram architecture increases density and performance," in *Proc. Symposium on VLSI Technology*, 2012, pp. 87–88.
- [3] S. Keckler, W. Dally, B. Khailany, M. Garland, and D. Glasco, "Gpus and the future of parallel computing," *Proc. Annual Int. Symp. Microarchitecture*, vol. 31, no. 5, pp. 7–17, 2011.
- [4] A. Putnam, A. Caulfield, E. Chung, D. Chiou, K. Constantinides, J. Demme, H. Esmailzadeh, J. Fowers, G. Gopal, J. Gray, M. Haselman, S. Hauck, S. Heil, A. Hormati, J.-Y. Kim, S. Lanka, J. Larus, E. Peterson, S. Pope, A. Smith, J. Thong, P. Xiao, and D. Burger, "A reconfigurable fabric for accelerating large-scale data-center services," in *Proc. IEEE Int. Symp. on Computer Architecture*, 2014, pp. 13–24.
- [5] J. Cong, M. A. Ghodrati, M. Gill, B. Grigorian, K. Gururaj, and G. Reinman, "Accelerator-rich architectures: Opportunities and progresses," in *Proc. ACM Design Automation Conf.*, San Francisco, CA, USA: ACM, 2014, 180:1–180:6, ISBN: 978-1-4503-2730-5.
- [6] C. Erdmann, D. Lowney, A. Lynam, A. Keady, J. McGrath, E. Cullen, D. Breathnach, D. Keane, P. Lynch, M. D. L. Torre, R. D. L. Torre, P. Lim, A. Collins, B. Farley, and L. Madden, "A heterogeneous 3d-ic consisting of two 28 nm fpga die and 32 reconfigurable high-performance data converters," *Solid-State Circuits, IEEE Journal of*, vol. 50, no. 1, pp. 258–269, 2015.
- [7] D. H. Kim, K. Athikulwongse, M. B. Healy, M. M. Hossain, M. Jung, I. Khorosh, G. Kumar, Y. J. Lee, D. L. Lewis, T. W. Lin, C. Liu, S. Panth, M. Pathak, M. Ren, G. Shen, T. Song, D. H. Woo, X. Zhao, J. Kim, H. Choi, G. H. Loh, H. H. S. Lee, and S. K. Lim, "Design and analysis of 3d-maps (3d massively parallel processor with stacked memory)," *IEEE Trans. Computers*, vol. 64, no. 1, pp. 112–125, 2015.
- [8] S. Panth, K. Samadi, Y. Du, and S. K. Lim, "High-density integration of functional modules using monolithic 3d-ic technology," in *Proc. IEEE Asia and South Pacific Design Automation Conf.*, 2013, pp. 681–686.

- [9] R. Mahajan, R. Sankman, N. Patel, D. W. Kim, K. Aygun, Z. Qian, Y. Mekonnen, I. Salama, S. Sharan, D. Iyengar, and D. Mallik, "Embedded multi-die interconnect bridge (emib) – a high density, high bandwidth packaging interconnect," in *2016 IEEE 66th Electronic Components and Technology Conference (ECTC)*, 2016, pp. 557–565.
- [10] X. Zhang, P. K. Jo, M. Zia, G. S. May, and M. S. Bakir, "Heterogeneous interconnect stitching technology with compressible microinterconnects for dense multi-die integration," *IEEE Electron Device Letters*, vol. 38, no. 2, pp. 255–257, 2017.
- [11] ITRS, *International Technology Roadmap for Semiconductors*, 2013.
- [12] Y. Zhang, A. Dembla, and M. S. Bakir, "Silicon micropin-fin heat sink with integrated tsvs for 3-d ics: Tradeoff analysis and experimental testing," *IEEE Transactions on Components, Packaging and Manufacturing Technology*, vol. 3, no. 11, pp. 1842–1850, 2013.
- [13] T. E. Sarvey, Y. Zhang, Y. Zhang, H. Oh, and M. S. Bakir, "Thermal and electrical effects of staggered micropin-fin dimensions for cooling of 3d microsystems," in *Fourteenth Intersociety Conference on Thermal and Thermomechanical Phenomena in Electronic Systems (ITherm)*, 2014, pp. 205–212.
- [14] T. E. Sarvey, Y. Zhang, L. Zheng, P. Thadesar, R. Gutala, C. Cheung, A. Rahman, and M. S. Bakir, "Embedded cooling technologies for densely integrated electronic systems," in *Custom Integrated Circuits Conference (CICC), 2015 IEEE*, 2015, pp. 1–8.
- [15] L. Zheng, Y. Zhang, X. Zhang, and M. S. Bakir, "Silicon interposer with embedded microfluidic cooling for high-performance computing systems," in *2015 IEEE 65th Electronic Components and Technology Conference (ECTC)*, 2015, pp. 828–832.
- [16] H. Oh, Y. Zhang, T. E. Sarvey, G. S. May, and M. S. Bakir, "Tsvs embedded in a microfluidic heat sink: High-frequency characterization and thermal modeling," in *2016 IEEE 20th Workshop on Signal and Power Integrity (SPI)*, 2016, pp. 1–4.
- [17] B. Black, M. Annavaram, N. Brekelbaum, J. DeVale, L. Jiang, G. H. Loh, D. McCaule, P. Morrow, D. W. Nelson, D. Pantuso, P. Reed, J. Rupley, S. Shankar, J. Shen, and C. Webb, "Die stacking (3d) microarchitecture," in *Proc. Annual Int. Symp. Microarchitecture*, 2006, pp. 469–479.
- [18] S. K. Samal, S. Panth, K. Samadi, M. Saedi, Y. Du, and S. K. Lim, "Fast and accurate thermal modeling and optimization for monolithic 3d ics," in *Proc. ACM Design Automation Conf.*, San Francisco, CA, USA: ACM, 2014, 206:1–206:6.

- [19] H. Wei, T. F. Wu, D. Sekar, B. Cronquist, R. F. Pease, and S. Mitra, "Cooling three-dimensional integrated circuits using power delivery networks," in *Proc. IEEE Int. Electron Devices Meeting*, 2012, pp. 14.2.1–14.2.4.
- [20] H. Oprins and E. Beyne, "Generic thermal modeling study of the impact of 3d-interposer material and thickness options on the thermal performance and die-to-die thermal coupling," in *Proc. IEEE Intersociety Conf on Thermal and Thermo-mechanical Phenomena in Electronic Systems.*, 2014, pp. 72–78.
- [21] X. Zhang, J. K. Lin, S. Wickramanayaka, S. Zhang, R. Weerasekera, R. Dutta, K. F. Chang, K.-J. Chui, H. Y. Li, D. S. Wee Ho, L. Ding, G. Katti, S. Bhattacharya, and D.-L. Kwong, "Heterogeneous 2.5d integration on through silicon interposer," *Applied Physics Reviews*, vol. 2, no. 2, 021308, 2015.
- [22] M. S. Gupta, J. L. Oatley, R. Joseph, G.-Y. Wei, and D. M. Brooks, "Understanding voltage variations in chip multiprocessors using a distributed power-delivery network," in *Proc. Design, Automation and Test in Europe*, Nice, France, 2007, pp. 624–629.
- [23] X. Zhang, T. Tong, S. Kanev, S. K. Lee, G.-Y. Wei, and D. Brooks, "Characterizing and evaluating voltage noise in multi-core near-threshold processors," in *Proc. Int. Symp. on Low Power Electronics and Design*, 2013, pp. 82–87.
- [24] R. Zhang, K. Wang, B. Meyer, M. Stan, and K. Skadron, "Architecture implications of pads as a scarce resource," in *Proc. IEEE Int. Symp. on Computer Architecture*, 2014, pp. 373–384.
- [25] H. Zhuang, S.-H. Weng, J.-H. Lin, and C.-K. Cheng, "Matex: A distributed framework for transient simulation of power distribution networks," in *Proc. ACM Design Automation Conf.*, 2014, pp. 1–6.
- [26] L. Zheng, Y. Zhang, and M. S. Bakir, "Full-chip power supply noise time-domain numerical modeling and analysis for single and stacked ics," *IEEE Transactions on Electron Devices*, vol. 63, no. 3, pp. 1225–1231, 2016.
- [27] H. He and J.-Q. Lu, "Modeling and analysis of pdn impedance and switching noise in tsv-based 3-d integration," *Electron Devices, IEEE Transactions on*, 2015.
- [28] J. Xie and M. Swaminathan, "Electrical and thermal cosimulation with nonconformal domain decomposition method for multiscale 3-d integrated systems," *Components, Packaging and Manufacturing Technology, IEEE Transactions on*, vol. 4, no. 4, pp. 588–601, 2014.
- [29] Y. Zhang and M. S. Bakir, "Integrated thermal and power delivery network co-simulation framework for single-die and multi-die assemblies," *IEEE Transactions*

on Components, Packaging and Manufacturing Technology, vol. 7, no. 3, pp. 434–443, 2017.

- [30] H. Su, F. Liu, A. Devgan, E. Acar, and S. Nassif, “Full chip leakage-estimation considering power supply and temperature variations,” in *Proc. Int. Symp. on Low Power Electronics and Design*, 2003, pp. 78–83.
- [31] A. Sridhar, A. Vincenzi, D. Atienza, and T. Brunschwiler, “3d-ice: A compact thermal model for early-stage design of liquid-cooled ics,” *IEEE Transactions on Computers*, vol. 63, no. 10, pp. 2576–2589, 2014.
- [32] W. Huang, M. R. Stan, K. Skadron, K. Sankaranarayanan, S. Ghosh, and S. Velusam, “Compact thermal modeling for temperature-aware design,” in *Proceedings of the 41st Annual Design Automation Conference*, San Diego, CA, USA: ACM, 2004, pp. 878–883, ISBN: 1-58113-828-8.
- [33] Y. Shao, Z. Peng, and J.-F. Lee, “Thermal-aware dc ir-drop co-analysis using non-conformal domain decomposition methods,” *Proceedings of the Royal Society of London A: Mathematical, Physical and Engineering Sciences*, vol. 468, no. 2142, pp. 1652–1675, 2012.
- [34] Y. Liu, R. Dick, L. Shang, and H. Yang, “Accurate temperature-dependent integrated circuit leakage power estimation is easy,” in *Proc. Design, Automation and Test in Europe*, 2007, pp. 1–6.
- [35] G. M. Link and N. Vijaykrishnan, “Thermal trends in emerging technologies,” in *Proceedings of the 7th International Symposium on Quality Electronic Design*, ser. ISQED '06, Washington, DC, USA: IEEE Computer Society, 2006, pp. 625–632, ISBN: 0-7695-2523-7.
- [36] Y. Zhan and S. S. Sapatnekar, “A high efficiency full-chip thermal simulation algorithm,” in *Proceedings of the 2005 IEEE/ACM International Conference on Computer-aided Design*, ser. ICCAD '05, San Jose, CA: IEEE Computer Society, 2005, pp. 635–638, ISBN: 0-7803-9254-X.
- [37] J. Xie and M. Swaminathan, “Electrical-thermal co-simulation of 3d integrated systems with micro-fluidic cooling and joule heating effects,” *IEEE Trans. Compon. Packag. Manuf. Technol.*, vol. 1, no. 2, pp. 234–246, 2011.
- [38] H. Qian, H. Liang, C.-H. Chang, W. Zhang, and H. Yu, “Thermal simulator of 3d-ic with modeling of anisotropic tsv conductance and microchannel entrance effects,” in *Proc. IEEE Asia and South Pacific Design Automation Conf.*, 2013, pp. 485–490.

- [39] Z. Wan, H. Xiao, Y. Joshi, and S. Yalamanchili, “Co-design of multicore architectures and microfluidic cooling for 3d stacked ics,” *Microelectron. J.*, vol. 45, no. 12, pp. 1814–1821, Dec. 2014.
- [40] A. Ziabari, E. K. Ardestani, J. Renau, and A. Shakouri, “Fast thermal simulators for architecture level integrated circuit design,” in *2011 27th Annual IEEE Semiconductor Thermal Measurement and Management Symposium*, 2011, pp. 70–75.
- [41] K. Athikulwongse, M. Pathak, and S. K. Lim, “Exploiting die-to-die thermal coupling in 3d ic placement,” in *Proceedings of the 49th Annual Design Automation Conference*, ser. DAC ’12, San Francisco, California: ACM, 2012, pp. 741–746, ISBN: 978-1-4503-1199-1.
- [42] Z. Wan, Y. J. Kim, and Y. Joshi, “Compact modeling of 3d stacked die inter-tier microfluidic cooling under non-uniform heat flux,” in *ASME 2012 International Mechanical Engineering Congress and Exposition*, American Society of Mechanical Engineers, 2012, pp. 911–917.
- [43] H. Oprins, V. O. Cherman, B. Vandeveld, G. V. der Plas, P. Marchal, and E. Beyne, “Numerical and experimental characterization of the thermal behavior of a packaged dram-on-logic stack,” in *Electronic Components and Technology Conference (ECTC), 2012 IEEE 62nd*, 2012, pp. 1081–1088.
- [44] J. Liu, B. Jaiyen, Y. Kim, C. Wilkerson, and O. Mutlu, “An experimental study of data retention behavior in modern dram devices: Implications for retention time profiling mechanisms,” in *Proceedings of the 40th Annual International Symposium on Computer Architecture*, ser. ISCA ’13, Tel-Aviv, Israel, 2013, pp. 60–71.
- [45] J. Liu, B. Jaiyen, R. Veras, and O. Mutlu, “Raidr: Retention-aware intelligent dram refresh,” in *Proceedings of the 39th Annual International Symposium on Computer Architecture*, ser. ISCA ’12, Portland, Oregon, 2012, pp. 1–12.
- [46] Z. Li, M. Mohamed, X. Chen, E. Dudley, K. Meng, L. Shang, A. R. Mickelson, R. Joseph, M. Vachharajani, B. Schwartz, and Y. Sun, “Reliability modeling and management of nanophotonic on-chip networks,” *IEEE Transactions on Very Large Scale Integration (VLSI) Systems*, vol. 20, no. 1, pp. 98–111, 2012.
- [47] C. Sun, M. T. Wade, Y. Lee, J. S. Orcutt, L. Alloatti, M. S. Georgas, A. S. Waterman, J. M. Shainline, R. R. Avizienis, S. Lin, *et al.*, “Single-chip microprocessor that communicates directly using light,” *Nature*, vol. 528, no. 7583, pp. 534–538, 2015.
- [48] M. Fish, P. McCluskey, and A. Bar-Cohen, “Thermal isolation within high-power 2.5d heterogenously integrated electronic packages,” in *2016 IEEE 66th Electronic Components and Technology Conference (ECTC)*, 2016, pp. 1847–1855.

- [49] L. Zheng, Y. Zhang, and M. S. Bakir, "A silicon interposer platform utilizing microfluidic cooling for high-performance computing systems," *IEEE Transactions on Components, Packaging and Manufacturing Technology*, vol. 5, no. 10, pp. 1379–1386, 2015.
- [50] Y. Demir and N. Hardavellas, "Parka: Thermally insulated nanophotonic interconnects," in *Proceedings of the 9th International Symposium on Networks-on-Chip*, ser. NOCS '15, Vancouver, BC, Canada: ACM, 2015, 1:1–1:8, ISBN: 978-1-4503-3396-2.
- [51] Y. Zhang, Y. Zhang, and M. S. Bakir, "Thermal design and constraints for heterogeneous integrated chip stacks and isolation technology using air gap and thermal bridge," *IEEE Transactions on Components, Packaging and Manufacturing Technology*, vol. 4, no. 12, pp. 1914–1924, 2014.
- [52] C. Zhang, H. S. Yang, and M. S. Bakir, "Highly elastic gold passivated mechanically flexible interconnects," *IEEE Transactions on Components, Packaging and Manufacturing Technology*, vol. 3, no. 10, pp. 1632–1639, 2013.
- [53] D. Kearney, "A numerical model of an inter-strata liquid cooling solution for a 3d ic architecture," in *Thermal Investigations of ICs and Systems (THERMINIC), 2010 16th International Workshop on*, 2010, pp. 1–6.
- [54] Y. Zhang, T. E. Sarvey, and M. S. Bakir, "Thermal challenges for heterogeneous 3d ics and opportunities for air gap thermal isolation," in *3D Systems Integration Conference (3DIC), 2014 International*, 2014, pp. 1–5.
- [55] U. Kang, H.-J. Chung, S. Heo, S.-H. Ahn, H. Lee, S.-H. Cha, J. Ahn, D. Kwon, J. H. Kim, J.-W. Lee, H.-S. Joo, W.-S. Kim, H.-K. Kim, E.-M. Lee, S.-R. Kim, K.-H. Ma, D.-H. Jang, N.-S. Kim, M.-S. Choi, S.-J. Oh, J.-B. Lee, T.-K. Jung, J.-H. Yoo, and C. Kim, "8gb 3d ddr3 dram using through-silicon-via technology," in *Solid-State Circuits Conference - Digest of Technical Papers, 2009. ISSCC 2009. IEEE International*, 2009, 130–131,131a.
- [56] Intel, *Intel Core™ i7 Processor Families for the LGA2011-0 Socket, Thermal Mechanical Specification and Design Guide*.
- [57] S. Li, J. H. Ahn, R. Strong, J. Brockman, D. Tullsen, and N. Jouppi, "Mcpat: An integrated power, area, and timing modeling framework for multicore and manycore architectures," in *Proc. Annual Int. Symp. Microarchitecture*, 2009, pp. 469–480.
- [58] Y. Zhang, "Hybrid microfluidic cooling and thermal isolation technologies for 3d ics," PhD thesis, Georgia Institute of Technology, May 2015.

- [59] Y. Zhang, Y. Zhang, T. Sarvey, C. Zhang, M. Zia, and M. Bakir, "Thermal isolation using air gap and mechanically flexible interconnects for heterogeneous 3-d ics," *IEEE Transactions on Components, Packaging and Manufacturing Technology*, vol. 6, no. 1, pp. 31–39, 2016.
- [60] Y. Zhang, T. E. Sarvey, Y. Zhang, M. Zia, and M. S. Bakir, "Numerical and experimental exploration of thermal isolation in 3d systems using air gap and mechanically flexible interconnects," in *2016 IEEE International Interconnect Technology Conference / Advanced Metallization Conference (IITC/AMC)*, 2016, pp. 83–85.
- [61] C. Zhang, H. S. Yang, H. D. Thacker, I. Shubin, J. E. Cunningham, and M. S. Bakir, "Mechanically flexible interconnects with contact tip for rematable heterogeneous system integration," *IEEE Transactions on Components, Packaging and Manufacturing Technology*, vol. 6, no. 11, pp. 1587–1594, 2016.
- [62] JEDEC, *Jedec standard for wide i/o single data rate*.
- [63] D. H. Kim, K. Athikulwongse, M. B. Healy, M. M. Hossain, M. Jung, I. Khorosh, G. Kumar, Y. J. Lee, D. L. Lewis, T. W. Lin, C. Liu, S. Panth, M. Pathak, M. Ren, G. Shen, T. Song, D. H. Woo, X. Zhao, J. Kim, H. Choi, G. H. Loh, H. H. S. Lee, and S. K. Lim, "Design and analysis of 3d-maps (3d massively parallel processor with stacked memory)," *IEEE Trans. Computers*, vol. 64, no. 1, pp. 112–125, 2015.
- [64] J. Cong, G. Luo, and Y. Shi, "Thermal-aware cell and through-silicon-via co-placement for 3d ics," in *Proc. ACM Design Automation Conf.*, ser. DAC '11, San Diego, California: ACM, 2011, pp. 670–675, ISBN: 978-1-4503-0636-2.
- [65] Altera, *Leveraging HyperFlex Architecture in Stratix 10 Devices to Achieve Maximum Power Reduction*.
- [66] Altera, *PowerPlay Early Power Estimators (EPE) and Power Analyzer (Stratix IV and Stratix V)*.
- [67] K. Sikka, J. Wakil, H. Toy, and H. Liu, "An efficient lid design for cooling stacked flip-chip 3d packages," in *Proc. IEEE Intersociety Conf on Thermal and Thermo-mechanical Phenomena in Electronic Systems.*, 2012, pp. 606–611.
- [68] L. Zheng, Y. Zhang, G. Huang, and M. S. Bakir, "Novel electrical and fluidic microbumps for silicon interposer and 3-d ics," *IEEE Transactions on Components, Packaging and Manufacturing Technology*, vol. 4, no. 5, pp. 777–785, 2014.
- [69] N. H.K.S. M. Alam and S. Hassoun, "System-level comparison of power delivery design for 2d and 3d ics," in *2009 IEEE International Conference on 3D System Integration*, 2009, pp. 1–7.

- [70] S. J. Park, N. Natu, and M. Swaminathan, "Analysis, design, and prototyping of temperature resilient clock distribution networks for 3-d ics," *IEEE Transactions on Components, Packaging and Manufacturing Technology*, vol. 5, no. 11, pp. 1669–1678, 2015.
- [71] Z. Zeng, X. Ye, Z. Feng, and P. Li, "Tradeoff analysis and optimization of power delivery networks with on-chip voltage regulation," in *Proceedings of the 47th Design Automation Conference*, ser. DAC '10, Anaheim, California: ACM, 2010, pp. 831–836, ISBN: 978-1-4503-0002-5.
- [72] C. Pan, S. Mukhopadhyay, and A. Naeemi, "System-level chip/package co-design for multi-core processors implemented with power-gating technique," in *2014 IEEE 23rd Conference on Electrical Performance of Electronic Packaging and Systems*, 2014, pp. 11–14.
- [73] S. R. Nassif, "Power grid analysis benchmarks," in *Proc. IEEE Asia and South Pacific Design Automation Conf.*, Seoul, Korea, 2008, pp. 376–381.
- [74] J. D. Meindl, "Interconnect opportunities for gigascale integration," *IEEE Micro*, vol. 23, no. 3, pp. 28–35, 2003.
- [75] Y. Kim, J. Cho, K. Kim, H. Kim, J. Kim, S. Sitaraman, V. Sundaram, and R. Tummala, "Analysis and optimization of a power distribution network in 2.5d ic with glass interposer," in *2014 International 3D Systems Integration Conference (3DIC)*, 2014, pp. 1–4.
- [76] Y. Kim, J. Cho, K. Kim, V. Sundaram, R. Tummala, and J. Kim, "Signal and power integrity analysis in 2.5d integrated circuits (ics) with glass, silicon and organic interposer," in *IEEE Electronic Components and Technology Conf.*, 2015, pp. 738–743.
- [77] A. Fourmigue, G. Beltrame, and G. Nicolescu, "Efficient transient thermal simulation of 3d ics with liquid-cooling and through silicon vias," in *Proc. Design, Automation and Test in Europe*, 2014, pp. 1–6.
- [78] H. H. Chen and D. D. Ling, "Power supply noise analysis methodology for deep-submicron vlsi chip design," in *Proceedings of the 34th Design Automation Conference*, 1997, pp. 638–643.
- [79] M. Eireiner, D. Schmitt-Landsiedel, P. Wallner, A. Schne, S. Henzler, and U. Fiedler, "Adaptive circuit block model for power supply noise analysis of low power system-on-chip," in *2009 International Symposium on System-on-Chip*, 2009, pp. 013–018.

- [80] J. Gu, J. Keane, and C. H. Kim, "Modeling, analysis, and application of leakage induced damping effect for power supply integrity," *IEEE Transactions on Very Large Scale Integration (VLSI) Systems*, vol. 17, no. 1, pp. 128–136, 2009.
- [81] Y. Ma, K. Wang, S. Dong, Y. Wang, and X. Hong, "Thermal effects of leakage power in 3d ics," in *The 2010 International Conference on Green Circuits and Systems*, 2010, pp. 578–583.
- [82] S. Sarangi, G. Ananthanarayanan, and M. Balakrishnan, "Lightsim: A leakage aware ultrafast temperature simulator," in *Proc. IEEE Asia and South Pacific Design Automation Conf.*, 2014, pp. 855–860.
- [83] T. Liu, C. C. Chen, and L. Milor, "Accurate standard cell characterization and statistical timing analysis using multivariate adaptive regression splines," in *Sixteenth International Symposium on Quality Electronic Design*, 2015, pp. 272–279.
- [84] ASU, *Predictive Technology Model*.
- [85] M. A. Karim, P. D. Franzon, and A. Kumar, "Power comparison of 2d, 3d and 2.5d interconnect solutions and power optimization of interposer interconnects," in *2013 IEEE 63rd Electronic Components and Technology Conference*, 2013, pp. 860–866.
- [86] S. Jangam, S. Pal, A. Bajwa, S. Pamarti, P. Gupta, and S. S. Iyer, "Latency, bandwidth and power benefits of the superchips integration scheme," in *2017 IEEE 67th Electronic Components and Technology Conference (ECTC)*, 2017, pp. 86–94.
- [87] X. Wu, W. Zhao, M. Nakamoto, C. Nimmagadda, D. Lisk, S. Gu, R. Radojcic, M. Nowak, and Y. Xie, "Electrical characterization for intertier connections and timing analysis for 3-d ics," *IEEE Transactions on Very Large Scale Integration (VLSI) Systems*, vol. 20, no. 1, pp. 186–191, 2012.
- [88] M. M. Shulaker, T. F. Wu, A. Pal, L. Zhao, Y. Nishi, K. Saraswat, H. S. P. Wong, and S. Mitra, "Monolithic 3d integration of logic and memory: Carbon nanotube fets, resistive ram, and silicon fets," in *Proc. IEEE Int. Electron Devices Meeting*, 2014, pp. 27.4.1–27.4.4.
- [89] R. Abbaspour, D. K. Brown, and M. S. Bakir, "Fabrication and electrical characterization of sub-micron diameter through-silicon via for heterogeneous three-dimensional integrated circuits," *Journal of Micromechanics and Microengineering*, vol. 27, no. 2, p. 025 011, 2017.
- [90] J. Kim, J. S. Pak, J. Cho, E. Song, J. Cho, H. Kim, T. Song, J. Lee, H. Lee, K. Park, S. Yang, M. S. Suh, K. Y. Byun, and J. Kim, "High-frequency scalable electrical model and analysis of a through silicon via (tsv)," *IEEE Transactions on Components, Packaging and Manufacturing Technology*, vol. 1, no. 2, pp. 181–195, 2011.

- [91] X. Zhang, V. Kumar, H. Oh, L. Zheng, G. S. May, A. Naeemi, and M. S. Bakir, "Impact of on-chip interconnect on the performance of 3-d integrated circuits with through-silicon vias: Part ii," *IEEE Transactions on Electron Devices*, vol. 63, no. 6, pp. 2510–2516, 2016.
- [92] Nangate, *Nangate 45nm Open Cell Library*.
- [93] H. Braunisch, A. Aleksov, S. Lotz, and J. Swan, "High-speed performance of silicon bridge die-to-die interconnects," in *Proc. IEEE Electrical Performance of Electronic Packaging and Systems*, 2011, pp. 95–98.
- [94] Y. Peng, T. Song, D. Petranovic, and S. K. Lim, "Silicon effect-aware full-chip extraction and mitigation of tsv-to-tsv coupling," *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, vol. 33, no. 12, pp. 1900–1913, 2014.
- [95] I. Ndip, B. Curran, K. Lobbecke, S. Guttowski, H. Reichl, K. D. Lang, and H. Henke, "High-frequency modeling of tsvs for 3-d chip integration and silicon interposers considering skin-effect, dielectric quasi-tem and slow-wave modes," *IEEE Transactions on Components, Packaging and Manufacturing Technology*, vol. 1, no. 10, pp. 1627–1641, 2011.
- [96] S. Abbaspour, M. Pedram, and P. Heydari, "Optimizing the energy-delay-ringing product in on-chip cmos line drivers," in *Fourth International Symposium on Quality Electronic Design, 2003. Proceedings.*, 2003, pp. 261–266.
- [97] T. Kondo, N. Takazawa, Y. Takemoto, M. Tsukimura, H. Saito, H. Kato, J. Aoki, K. Kobayashi, S. Suzuki, Y. Gomi, S. Matsuda, and Y. Tadaki, "3-d-stacked 16-mpixel global shutter cmos image sensor using reliable in-pixel four million microbump interconnections with 7.6-um pitch," *IEEE Transactions on Electron Devices*, vol. 63, no. 1, pp. 128–137, 2016.
- [98] S. V. Huylenbroeck, M. Stucchi, Y. Li, J. Slabbekoorn, N. Tutunjan, S. Sardo, N. Jourdan, L. Bogaerts, F. Beirnaert, G. Beyer, and E. Beyne, "Small pitch, high aspect ratio via-last tsv module," in *2016 IEEE 66th Electronic Components and Technology Conference (ECTC)*, 2016, pp. 43–49.
- [99] P. Batude, M. Vinet, A. Pouydebasque, C. L. Royer, B. Previtali, C. Tabone, J. M. Hartmann, L. Sanchez, L. Baud, V. Carron, A. Toffoli, F. Allain, V. Mazzocchi, D. Lafond, O. Thomas, O. Cueto, N. Bouzaida, D. Fleury, A. Amara, S. Deleonibus, and O. Faynot, "Advances in 3d cmos sequential integration," in *2009 IEEE International Electron Devices Meeting (IEDM)*, 2009, pp. 1–4.
- [100] K. Acharya, K. Chang, B. W. Ku, S. Panth, S. Sinha, B. Cline, G. Yeric, and S. K. Lim, "Monolithic 3d ic design: Power, performance, and area impact at 7nm," in

2016 17th International Symposium on Quality Electronic Design (ISQED), 2016, pp. 41–48.

- [101] M. M. Shulaker, G. Hills, R. S. Park, R. T. Howe, K. Saraswat, H. S. P. Wong, and S. Mitra, “Three-dimensional integration of nanotechnologies for computing and data storage on a single chip,” *Nature*, vol. 547, no. 7661, pp. 74–78, 2017.
- [102] M. Bohr, “The new era of scaling in an soc world,” in *2009 IEEE International Solid-State Circuits Conference - Digest of Technical Papers*, 2009, pp. 23–28.
- [103] K. J. Kuhn, “Cmos scaling for the 22nm node and beyond: Device physics and technology,” in *Proceedings of 2011 International Symposium on VLSI Technology, Systems and Applications*, 2011, pp. 1–2.
- [104] X. Huang, W.-C. Lee, C. Kuo, D. Hisamoto, L. Chang, J. Kedzierski, E. Anderson, H. Takeuchi, Y.-K. Choi, K. Asano, V. Subramanian, T.-J. King, J. Bokor, and C. Hu, “Sub-50 nm p-channel finfet,” *IEEE Transactions on Electron Devices*, vol. 48, no. 5, pp. 880–886, 2001.
- [105] T.-J. King, “Finfets for nanoscale cmos digital integrated circuits,” in *Proceedings of the 2005 IEEE/ACM International Conference on Computer-aided Design*, ser. ICCAD ’05, San Jose, CA: IEEE Computer Society, 2005, pp. 207–210, ISBN: 0-7803-9254-X.
- [106] K. Roy, S. Mukhopadhyay, and H. Mahmoodi-Meimand, “Leakage current mechanisms and leakage reduction techniques in deep-submicrometer cmos circuits,” *Proceedings of the IEEE*, vol. 91, no. 2, pp. 305–327, 2003.
- [107] E. Cai and D. Marculescu, “Tei-turbo: Temperature effect inversion-aware turbo boost for finfet-based multi-core systems,” in *Proceedings of the IEEE/ACM International Conference on Computer-Aided Design*, ser. ICCAD ’15, Austin, TX, USA: IEEE Press, 2015, pp. 500–507, ISBN: 978-1-4673-8389-9.
- [108] W. Lee, K. Han, Y. Wang, T. Cui, S. Nazarian, and M. Pedram, “Tei-power: Temperature effect inversion-aware dynamic thermal management,” *ACM Trans. Des. Autom. Electron. Syst.*, vol. 22, no. 3, 51:1–51:25, Apr. 2017.
- [109] S. Y. Kim, Y. M. Kim, K. H. Baek, B. K. Choi, K. R. Han, K. H. Park, and J. H. Lee, “Temperature dependence of substrate and drain currents in bulk finfets,” *IEEE Transactions on Electron Devices*, vol. 54, no. 5, pp. 1259–1264, 2007.
- [110] Y. Zu, W. Huang, I. Paul, and V. J. Reddi, “Ti-states: Processor power management in the temperature inversion region,” in *2016 49th Annual IEEE/ACM International Symposium on Microarchitecture (MICRO)*, 2016, pp. 1–13.

- [111] W. Wahby, L. Zheng, Y. Zhang, and M. S. Bakir, “A simulation tool for rapid investigation of trends in 3-dic performance and power consumption,” *IEEE Transactions on Components, Packaging and Manufacturing Technology*, vol. 6, no. 2, pp. 192–199, 2016.
- [112] X. Zhao, L. Han, and Z. Feng, “A performance-guided graph sparsification approach to scalable and robust spice-accurate integrated circuit simulations,” *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, vol. 34, no. 10, pp. 1639–1651, 2015.
- [113] S. J. Park, B. Bae, J. Kim, and M. Swaminathan, “Application of machine learning for optimization of 3-d integrated circuits and systems,” *IEEE Transactions on Very Large Scale Integration (VLSI) Systems*, vol. 25, no. 6, pp. 1856–1865, 2017.
- [114] K. Padmaraju, D. F. Logan, T. Shiraishi, J. J. Ackert, A. P. Knights, and K. Bergman, “Wavelength locking and thermally stabilizing microring resonators using dithering signals,” *Journal of Lightwave Technology*, vol. 32, no. 3, pp. 505–512, 2014.

VITA

Yang Zhang was born in Xi'an, Shannxi province, China, in October 1990. He received his B.S. degrees in Microelectronics and Math (secondary major) from Peking University, China, in 2012. He also received the M.S. degree in Electrical and Computer Engineering from Georgia Institute of Technology, Atlanta, USA in 2015, where he is currently a PhD candidate.

From 2013 to present, he has been a graduate research assistant in the Georgia Tech Integrated 3-D System laboratory supervised by Dr. Muhannad S. Bakir. His primary research is in the area of 2.5-D and 3-D IC design, modeling and optimization with a concentration on thermal and power delivery network analysis. His other research interests include physical design flow development and signal integrity for emerging system integration technologies.