

UC San Diego

UC San Diego Electronic Theses and Dissertations

Title

System level design of power distribution network for mobile computing platforms

Permalink

<https://escholarship.org/uc/item/28q3304f>

Author

Shayan Arani, Amirali

Publication Date

2011

Peer reviewed|Thesis/dissertation

UNIVERSITY OF CALIFORNIA, SAN DIEGO

**System Level Design of Power Distribution Network
for Mobile Computing Platforms**

A dissertation submitted in partial satisfaction of the
requirements for the degree
Doctor of Philosophy

in

Computer Science (Computer Engineering)

by

Amirali Shayan Arani

Committee in charge:

Professor Chung-Kuan Cheng, Chair
Professor Arif Ege Engin
Professor Tara Javidi
Professor Andrew B. Kahng
Professor Ryan Kastner
Professor Dean Tullsen

2011

Copyright
Amirali Shayan Arani, 2011
All rights reserved.

The dissertation of Amirali Shayan Arani is approved,
and it is acceptable in quality and form for publication
on microfilm and electronically:

Chair

University of California, San Diego

2011

DEDICATION

To the memory of Baba Shayan, my inspiration for whole life.

To the best parents in the world, Farideh and Shahrokh.

To Yassi, my little sister for always caring about me.

To Behrokh for her unlimited love, companionship and kindness.

EPIGRAPH

*Like God, if this world I could control
Making new world would be my role
I would create the world anew, whole
Such that the free soul would attain desired goal.*

—Omar Khayyam, Persian Mathematician and Poet, 11th century.

TABLE OF CONTENTS

Signature Page	iii
Dedication	iv
Epigraph	v
Table of Contents	vi
List of Figures	x
List of Tables	xiii
Acknowledgments	xiv
Vita	xvii
Abstract of the Dissertation	xx
Chapter 1 Introduction	1
1.1 Motivation and Challenges	3
1.1.1 Limited Area and Form Factor	3
1.1.2 System Level Co-design Requirements	4
1.1.3 Cost	5
1.1.4 On-chip Variation	5
1.1.5 Heterogeneous System	6
1.2 Dissertation Contributions	7
Chapter 2 System Level Design and Analysis of Mobile Platforms PDN	11
2.1 Introduction	11
2.2 System PDN Design	15
2.2.1 PCB PDN Design	15
2.2.2 Package for PDN Design	17
2.2.3 Silicon PDN Design	17
2.3 PDN Analysis Methodologies	18
2.3.1 Frequency Domain Analysis	19
2.3.2 Frequency Domain Automation	21
2.3.3 Time Domain Analysis	22
2.3.4 PDN Optimization Using Time Domain Analysis	23
2.4 Impact of Logic on Performance	26
2.5 Summary	30

Chapter 3	Design and Reliability Aspects of 3D Integration PDN	31
	3.1 Introduction	31
	3.2 Overview of Reliability Challenges	33
	3.3 Time and Frequency Domain Analysis of the Three Di- dimensional Networks	34
	3.3.1 Parallel 3D PDN Analysis Flow	34
	3.3.2 Design Planning for the 3D Power Grid	36
	3.3.3 Frequency Domain Analysis of 3D PDN	38
	3.3.4 On-chip Inductance Scaling in 3D PDN	39
	3.3.5 On-chip Resistance Scaling in 3D PDN	39
	3.3.6 Sensitivity of Impedance to Stacked Layers	41
	3.4 Reliability-aware 3D PDN Model Considering Substrate Coupling	44
	3.4.1 Substrate Coupling Model in 3D PDN	44
	3.4.2 Frequency Dependent Through Silicon Via Model	44
	3.4.3 Frequency Dependent Substrate Coupling Model .	46
	3.4.4 Frequency Dependent Power Grid Model in Each Tier	48
	3.5 Reliability of Stacked Power Grid	49
	3.5.1 Electromigration Constraint	49
	3.5.2 Thermo-mechanical Reliability of 3D Stacking . .	51
	3.5.3 Reliability-aware Experimental Results and Anal- ysis	51
	3.6 Summary	55
	3.7 Acknowledgments	56
Chapter 4	What Would be a Worst Current Scenario?	57
	4.1 Resonance-aware Methodology for System Level Power Distribution Network Co-design	57
	4.1.1 Introduction	57
	4.1.2 Theoretical Background: Resonance-aware Mod- ulation	59
	4.1.3 Broadband Frequency Domain Co-design	62
	4.1.4 Resonance-aware Time Domain Methodology . . .	63
	4.1.5 Experimental Results	66
	4.2 Vector-based Rogue Wave Synthesis Algorithm Using Re- alistic Current Activity	70
	4.2.1 Complexity of Rogue Wave Synthesis Algorithm .	70
	4.2.2 Sensitivity of Voltage Noise to Different Q -factors	72
	4.3 Summary	77
	4.4 Acknowledgments	77

Chapter 5	Low Power Distribution Network Regulations	79
5.1	On-die LDO-based Design and Optimization of PDN under Worst Loading	80
5.1.1	Introduction	80
5.1.2	Problem Formulation	81
5.1.3	LDO Frequency Domain Characteristics	83
5.1.4	Analytical Worst-case Voltage Drop	85
5.1.5	Exact analytical value of LDO peak voltage drop	88
5.1.6	LDO-based Experimental Results and Tradeoff Analysis	88
5.2	Parallel Flow to Analyze the Impact of the Off-chip Voltage Regulators	91
5.2.1	Overview of Voltage Regulator Impact	91
5.2.2	Complete Distributed PDN Model	92
5.2.3	MNA-based Early PDN Analysis Flow	95
5.2.4	Modified Nodal Analysis-based Frequency Domain Analysis	96
5.2.5	Time Domain Signal Recovering Using Vector Fitting	98
5.2.6	VRM-based PDN Results and Analysis	99
5.2.7	Voltage Regulator Module Impact on Noise	102
5.3	Summary	105
5.4	Acknowledgments	105
Chapter 6	Power Distribution Impact on Performance	107
6.1	Estimation of Power Integrity Impact to Low Power Processor Performance through Pre-Silicon Simulation and Post-Silicon Measurements	108
6.1.1	Power Delivery Network Pre-Silicon Analysis	108
6.1.2	Silicon Measurement for Impact of Power Integrity on Processor Performance	111
6.1.3	Mid-High Frequency Impedance Resonance Impact on Performance	112
6.1.4	Low Frequency Impedance Impact on Performance	112
6.2	Worst-case Performance Prediction under Supply Voltage and Temperature Variation	116
6.2.1	Overview of Dynamic Supply Noise	116
6.2.2	Related Work	118
6.2.3	Implementation Flow	120
6.2.4	Modeling Methodology	121
6.2.5	Multivariate Adaptive Regression Splines	123
6.2.6	Accurate Cell Delay Modeling	124
6.2.7	Worst-case Performance Model	127

	6.2.8	Experimental Results and Validation	128
	6.3	Summary	130
	6.4	Acknowledgments	131
Chapter 7		Conclusions	133
	7.1	Thesis Summary	134
	7.2	Future Research Directions	135
	7.2.1	Architecture/Software/PDN Co-design (PDN-friendly Architectures)	135
	7.2.2	Variation-aware Timing-Voltage Analysis Integra- tion	136
	7.2.3	Distributed On-die Regulations and Power Man- agement	136
		Bibliography	137

LIST OF FIGURES

Figure 2.1:	Typical PDN system consists of voltage regulator, bulk caps, top side caps, back side caps mounted on the PCB, and die side caps mounted on the package.	12
Figure 2.2:	Right plot shows typical response of a PDN system to a sudden increase in current demand by the die. Left plot is an illustration of PDN impedance profile.	13
Figure 2.3:	Normalized sensitivity of NOR2 ring oscillator frequency to voltage supply variation.	14
Figure 2.4:	Multiple voltage domains voltage hotspots.	18
Figure 2.5:	Early PDN analysis flow in frequency domain.	20
Figure 2.6:	An example of floorplan spreadsheet for frequency domain early-stage analysis.	22
Figure 2.7:	Voltage noise time domain analysis integration in the physical design flow.	24
Figure 2.8:	Voltage noise time domain analysis for the complete PDN using lumped PCB two ports S parameter model.	25
Figure 2.9:	Typical functional block load current and voltage drop.	25
Figure 2.10:	Line of CMOS inverters powered by a noisy supply source with frequency F_N	27
Figure 2.11:	Averaging of effective V_{dd} over a long line of inverters.	28
Figure 3.1:	3D PDN stacked core and memory model.	32
Figure 3.2:	Tier to Tier (T2T) power connection model of 3D stacked die.	35
Figure 3.3:	Stacked PDN reliability-aware and noise co-analysis flow.	36
Figure 3.4:	3D PDN parasitic $RLCK$ model with package and VRM.	37
Figure 3.5:	Impedance spectrum of 3D stacked PDN versus via distribution density.	39
Figure 3.6:	Impedance magnitude of 3D PDN with and without on-chip inductance.	40
Figure 3.7:	On-chip inductance scaling impact on impedance magnitude in 3D stacking.	40
Figure 3.8:	Resistance scaling impact on impedance spectrum in 3D stacking.	41
Figure 3.9:	Impedance spectrum scaling with different stacked layers.	42
Figure 3.10:	Simulation CPU time as a function of number of processors for 3D stacking.	42
Figure 3.11:	Time domain noise magnitude of 3D stacking versus via density.	43
Figure 3.12:	Through silicon via modeling in HFSS [2].	45
Figure 3.13:	TSV $RLGC$ equivalent $RLCG(f)$ model with substrate coupling.	46
Figure 3.14:	TSV equivalent resistance R (D=diameter, P=pitch).	47
Figure 3.15:	TSV equivalent inductance L (D=diameter, P=pitch).	47
Figure 3.16:	Substrate equivalent shunt resistance (D=diameter, P=pitch).	48

Figure 3.17: Substrate equivalent shunt coupling (D=diameter, P=pitch).	49
Figure 3.18: Power grid model in each tier.	50
Figure 3.19: Normalized TSV thermo-mechanical failure rate(%).	52
Figure 3.20: Decoupling capacitor allocation tradeoff versus TSV.	53
Figure 3.21: Power noise of 3D stacking considering substrate coupling model.	54
Figure 3.22: Optimization cost function($Failure\ rate_{TSV} \times max\ \Delta V$).	55
Figure 4.1: Resonance-aware load current modulation impact on voltage variation (Top: modulated current, Middle: package current, Bottom: on-chip voltage variation).	60
Figure 4.2: RLC resonance with $R=40m$, $L=2nH$, and $C=10nF$	61
Figure 4.3: RLC resonance after the $2nH$ inductance is segmented with shunt capacitances.	63
Figure 4.4: Current waveform generated from individual wavelets.	65
Figure 4.5: Flow chart of the time domain worst case Z-aware PDN analysis.	65
Figure 4.6: Synthetic resonance-aware current generation.	67
Figure 4.7: Frequency energy content of the modulated vector current.	68
Figure 4.8: Multiple modes resonance-aware modulation.	69
Figure 4.9: Lumped model of multiple stage power distribution.	72
Figure 4.10: Synthesized rogue wave stimuli for multiple Q -factors set 1.	73
Figure 4.11: Synthesized rogue wave stimuli for multiple Q -factors set 2.	74
Figure 4.12: PDN Impedance profile change for different Q -factors.	75
Figure 4.13: Synthesize of rogue wave stimuli for multiple Q -factors set 3.	76
Figure 4.14: Impedance profile in different Q -factor PDN system.	78
Figure 5.1: Model of LDO-based power distribution.	82
Figure 5.2: LDO-PDN system worst-case noise optimization flow.	82
Figure 5.3: Lumped approximation of the LDO-PDN model.	84
Figure 5.4: LDO poles and zeros impact on the system loop gain.	84
Figure 5.5: Phase of the power distribution impedance $Z(f)$	85
Figure 5.6: Amplitude of the power distribution impedance $Z(f)$	86
Figure 5.7: Step response of the LDO based power distribution.	87
Figure 5.8: Worst-case voltage drop based on LDO step response	89
Figure 5.9: Peak voltage drop undershoot as a function of LDO power and on-chip decap.	90
Figure 5.10: Peak voltage drop overshoot as a function of LDO power and on-chip decap.	90
Figure 5.11: Early analysis flow for the complete PDN network.	92
Figure 5.12: Complete power distribution network model.	93
Figure 5.13: Voltage regulator model.	94
Figure 5.14: 2D partitioned model for package and on-chip grid.	95
Figure 5.15: Framework for calculating the voltage response of the P/G networks.	98

Figure 5.16: Magnitude of voltage response of PDN in frequency domain without VRM.	99
Figure 5.17: Recovered time domain voltage response comparison with HSPICE.	100
Figure 5.18: Frequency domain resonance peaks of the PDN with VRM. . .	102
Figure 5.19: Comparison of time domain voltage response with and without VRM.	102
Figure 5.20: Impedance profile for the 4×4 2D PDN partitions.	103
Figure 5.21: PDN voltage response over chip operation current profile. . . .	104
Figure 6.1: System impedance sensitivity due to package layer scaling and die side cap scaling ¹	109
Figure 6.2: Miscorrelation between commercially available equivalent circuit model and physical model we developed ²	110
Figure 6.3: Impedance measurement correlation with pre-silicon simulation using three main harmonics.	111
Figure 6.4: Impact of systematic increase of bulk cap on improving the low frequency impedance on F_{max} across different processes.	114
Figure 6.5: Measurement results for <i>Low</i> frequency voltage droop impact on F_{max} for multiple clock cycles and loops.	115
Figure 6.6: Accurate worst-case performance-driven power distribution network optimization flow.	118
Figure 6.7: Implementation flow.	120
Figure 6.8: Delay of an inverter cell versus noise slew, for different input slew values.	125
Figure 6.9: Impact of supply voltage noise offset on cell delay.	126
Figure 6.10: Sample inverter delay and output slew models in 65nm.	127

LIST OF TABLES

Table 3.1:	Simulation time of 3D PDN versus number of processors on clustered FWgrid multiple cores for proposed flow versus HSPICE .	41
Table 4.1:	Setup table for current modulation of the processor.	66
Table 4.2:	Sensitivity of noise to Q -factor setting of synthesized rogue wave.	72
Table 5.1:	Simulation time and speedup of proposed frequency flow versus HSPICE and Cadence Spectre for 300nsec using 1 and 4 cores. .	101
Table 5.2:	Simulation time and speedup of proposed frequency flow versus HSPICE and Cadence Spectre for 300nsec using 8 and 16 cores.	101
Table 6.1:	Measured F_{max} sensitivity to systematical removal of package capacitors.	113
Table 6.2:	Bulk decap impact on the performance.	115
Table 6.3:	List of parameters used for performance modeling.	121
Table 6.4:	Model stability versus random selection of the training set. . . .	128
Table 6.5:	Comparison of our proposed worst-case performance model and SPICE for an inverter chain. Rank values are out of 30720 configurations.	130
Table 6.6:	Comparison of proposed worst-case performance model and SPICE for a 2-input NAND chain. Rank values are out of 30720 configurations.	130
Table 6.7:	Comparison of proposed worst-case performance model and SPICE for a mixed inverter-NAND chain. Rank values are out of 30720 configurations.	130

ACKNOWLEDGMENTS

First and foremost, I would like to thank my advisor, Prof. Chung-Kuan Cheng, for his guidance and support through all aspects of my graduate research. His vision, breadth and depth of knowledge and intuition will always inspire me. I certainly feel privileged and grateful for having the opportunity to work directly with Prof. Cheng and learn lessons that will remain with me for the rest of my life. Prof. Cheng always encourages me to explore new and diverse topics in research and to see what we can improve. I feel blessed for such an opportunity.

I would like to thank all of my thesis committee members, in alphabetical order, Professor Arif Ege Engin, Professor Tara Javidi, Professor Andrew B. Kahng, Professor Ryan Kastner, and Professor Dean Tullsen, for providing a lot of insightful suggestions and constructive discussions regarding my research work and PhD dissertation.

Words can not describe my sincere thanks to my managers in Qualcomm, Dr. Christopher Pan, Lew Chua-Eoan and Matt Severson for providing me the industrial opportunity to participate in many diverse and interesting projects. Their confidence in me and their encouragement help me to grow and to shape the future of my research career.

Special thanks to Dr. Mikhail Popovich, Prof. Wenjian Yu, Yi Zhu, Wanning Zhang, Xiang Hu, Vincent Peng, Yulei Zhang, Renshen Wang, Shih-Hung Weng, Peng Du, Kambiz Samadi, Vasileios Kontronis, Shervin Sharifi, Houman Homayoun, Thomas Toms, John McDonalds, Ling Zhang, Haikun Zhu, Sorin Dobre, Kevin Bowles, Jason Xiaoming Chen, Xiaohua Kong, Qualcomm Qpower members, UCSD Computer Science and Engineering Department faculties, staffs and students and all the folks whom I had the opportunity to work with.

Finally, I owe many thanks to my family for all their support through the years. First and foremost, I am grateful to my parents, Farideh Ghezeli and Shahrokh Shayan for their unconditional love and sacrifices. For their encouragement and wisdom, without which I would have never made it this far. I would like to specially thank my little sister, Yassaman for being the best sister of the world and caring about me in all ups and downs. Above all, I would like to specially

thank my caring significant other Behrokh for her unlimited love, companionship and kindness. This dissertation is dedicated to them.

The material in this thesis is based on the following publications. Chapter 3 is based on the following publications:

- A. Shayan, X. Hu, H. Peng, W. Zhang, M. Popovich, L. Chua-Eoan, C.K. Cheng, “Power Distribution Co-design for Nanoscale Stacked Silicon ICs”, *IEEE Conference on Electrical Performance of Electronic Packaging (EPEP)*, 2008.
- A. Shayan, X. Hu, M. Popovich, A.E. Engin, C.K. Cheng, “3D Stacked Power Distribution Considering Substrate Coupling”, *International Conference on Computer Design (ICCD)*, 2009.
- A. Shayan, X. Hu, H. Peng, W. Yu, T. Toms, M. Popovich, X. Chen, C.K. Cheng, “Reliability-aware Through Silicon Via Planning for Nanoscale Stacked Silicon ICs”, *Design Automation and Test in Europe (DATE)*, 2009.

The dissertation author was the primary researcher and author, and the co-authors involved in the above publications directed, supervised, and assisted in the research which forms the basis for that material.

Chapter 4 is based on the following publication:

- A. Shayan, K. Bowles, S. Dobre, M. Popovich, X. Chen, C. Pan , “Resonance-aware Modulation Methodology for System Level Power Distribution Co-design”, *IEEE Conference on Electrical Performance of Electronic Packaging (EPEP)*, 2009.

The dissertation author was the primary researcher and author, and the co-authors involved in the above publications directed, supervised, and assisted in the research which forms the basis for that material.

Chapter 5 is based on the following publications:

- A. Shayan, X. Hu, H. Peng, W. Yu, W. Zhang, C.K. Cheng, M. Popovich, X. Chen, L. Chua-Eaon, Xiaohua Kong, “Parallel Flow to Analyze the Impact

of the Voltage Regulator Model in Nanoscale Power Distribution Network”, *The International Symposium on Quality Electronic Design (ISQED)*, 2009.

- A. Shayan, X. Hu, C.K. Cheng, W. Yu, C. Pan, “Linear Dropout Regulator based Power Distribution Design under Worst Loading”, *IEEE International Conference on ASIC*, 2011.

The dissertation author was the primary researcher and author, and the co-authors involved in the above publications directed, supervised, and assisted in the research which forms the basis for that material.

Chapter 6 is based on the following publications:

- A. Shayan, C. Pan, M. Popovich, K. Bowles, “Estimation of Power Integrity Impact to Low Power Processor Performance through Pre-Silicon Simulation and Post-Silicon Measurements”, *AMSE InterPack*, 2011.
- Chung-Kuan Cheng, Andrew B. Kahng, Kambiz Samadi and Amirali Shayan, “Worst-case Performance Prediction under Supply Voltage and Temperature Variation”, *ACM/IEEE System-Level Interconnect Prediction (SLIP)*, 2010, pp. 91–96.

The dissertation author was the researcher and co-author of both papers. My co-authors (Dr. Christopher Pan, Dr. Mikhail Popovich, Kevin Bowles, Prof. Chung-Kuan Cheng, Prof. Andrew B. Kahng, and Dr. Kambiz Samadi) have all kindly approved the inclusion of the aforementioned publications in my thesis.

For the rest of the publications, my co-authors (Prof. Wenjian Yu, Prof. Ege Engin, Lew Chua-Eoan, Christopher Pan, Mikhail Popovich, Xiaoming Chen, Xiaohua Kong, Wanping Zhang, Xiang Hu, He Peng, Sorin Dobre, Kevin Bowles, Thomas Toms, and Du Peng) have all kindly approved the inclusion of the aforementioned publications in my thesis.

VITA

- 2005 B. S. in Electrical Engineering, University of Tehran, Iran.
- 2008 M. S. in Computer Science, University of California, San Diego, La Jolla, CA, USA.
- 2011 Ph. D. in Computer Science (Computer Engineering), University of California, San Diego, La Jolla, CA, USA.

PUBLICATIONS

- H. Homayoun, V. Kontorinis, A. Shayan, T. Lin, D. Tullsen, “Dynamically Heterogeneous Cores Through 3D Resource Pooling”, *IEEE International Symposium on High Performance Computer Architecture (HPCA)*, 2012.
- A. Shayan, X. Hu, C.K. Cheng, W. Yu, C. Pan, “Linear Dropout Regulator based Power Distribution Design under Worst Loading”, *IEEE International Conference on ASIC*, 2011.
- C. Pan, A. Shayan, M. Popovich, K. Bowles, “Estimation of Power Integrity Impact to Low Power Processor Performance through Pre-Silicon Simulation and Post-Silicon Measurements”, *AMSE InterPack*, 2011.
- L. Chua-Eoan, B. Andreev, C. Pan, A. Shayan, X. Kong, M. Popovich, M. Calle, I.-J. Chang, “On-chip Sensor For Measuring Dynamic Power Supply Noise Of The Semiconductor Chip”, *US Patent*, US20110193589, Aug 11, 2011.
- C.-K. Cheng, A. B. Kahng, K. Samadi, A. Shayan, “Worst-case Performance Prediction Under Supply Voltage and Temperature Variation”, *ACM/IEEE System-Level Interconnect Prediction (SLIP)*, 2010.
- P. Du, X. Hu, A. Shayan, X. Chen, A. E. Engin, C.-K. Cheng, “Worst-case Noise Prediction With Non-zero Current Transition Times for Early Power Distribution System Verification”, *International Symposium on Quality Electronic Design (ISQED)*, 2010.
- S. Dobre, A. Shayan, M. Popovich, K. Bowles, X. Chen, C. Pan, “Package/PCB-aware On-die Power Grid Noise Analysis”, *IEEE Design Automation Conference, User Track*, 2010.
- S. Dobre, K. Sajid, A. Shayan, R. Jalilzainali, S. Dundigal, S. Evan, “ESD Verification and ESD-aware Design Optimization for Complex System-On-a-Chip Design”, *IEEE Design Automation Conference, User Track*, 2010.

- W. Zhang, L. Zhang, A. Shayan, W. Yu, X. Hu, Z. Zhu, E. Engin, and C.K. Cheng, “On-chip Power Network Optimization with Decoupling Capacitors and Controlled-ESRs”, *Asia and South Pacific Design Automation Conference (ASPDAC)*, 2010.
- X. Hu, W. Zhao, P. Du, A. Shayan, C.K. Cheng, “An Adaptive Parallel Flow for Power Distribution Network Simulation Using Discrete Fourier Transform”, *Asia and South Pacific Design Automation Conference (ASPDAC)*, 2010.
- A. Shayan, X. Hu, M. Popovich, A.E. Engin, C.K. Cheng, “3D Stacked Power Distribution Considering Substrate Coupling”, *International Conference on Computer Design (ICCD)*, 2009.
- A. Shayan, X. Hu, H. Peng, W. Yu, T. Toms, M. Popovich, X. Chen, C.K. Cheng, “Reliability-aware Through Silicon Via Planning for Nanoscale Stacked Silicon ICs”, *Design Automation and Test in Europe (DATE)*, 2009.
- A. Shayan, X. Hu, H. Peng, W. Yu, W. Zhang, C.K. Cheng, M. Popovich, X. Chen, L. Chua-Eaon, Xiaohua Kong, “Parallel Flow to Analyze the Impact of the Voltage Regulator Model in Nanoscale Power Distribution Network”, *The International Symposium on Quality Electronic Design (ISQED)*, 2009.
- A. Shayan, K. Bowles, S. Dobre, M. Popovich, X. Chen, C. Pan, “Resonance-aware Modulation Methodology for System Level Power Distribution Co-design”, *IEEE Conference on Electrical Performance of Electronic Packaging (EPEP)*, 2009.
- Vasileios Kontorinis, A. Shayan, R. Kumar, D. Tullsen, “Reducing Peak Power with a Table-Driven Adaptive Processor Core”, *IEEE/ACM International Symposium on Microarchitecture (MICRO)*, 2009.
- W. Zhang, Y. Zhu, W. Yu, A. Shayan, R. Wang; Z. Zhu; C.K. Cheng, “Noise Minimization During Power-up Stage for a Multi-domain Power Network”, *Asia and South Pacific Design Automation Conference (ASPDAC)*, 2009.
- X. Hu, W. Zhao, P. Du, Y. Zhang, A. Shayan, C. Pan, A.E. Engin, C.K. Cheng, “On the Bound of Time-domain Power Supply Noise based on Frequency-domain Target Impedance”, *System Level Interconnect Prediction (SLIP)*, 2009.
- W. Zhang, W. Yu, X. Hu, A. Shayan, A.E. Engin, C.K. Cheng, “Predicting the Worst-case Voltage Violation in a 3D Power Network”, *System Level Interconnect Prediction (SLIP)*, 2009.
- A. Shayan, X. Hu, H. Peng, W. Zhang, M. Popovich, L. Chua-Eaon, C.K. Cheng, “3D Power Distribution Co-design for Nanoscale Stacked Silicon ICs”, *IEEE Conference on Electrical Performance of Electronic Packaging (EPEP)*, 2008.

A. Shayan, “Online Thermal-aware Scheduling for Multiple Clock Domain CMPs”, *IEEE International SOC Conference (SOCC)*, 2007.

A. Shayan, Y. Zhu, Y.N. Cheng, C.K. Cheng, S.F. Lin, P.S. Chen, “Exploring Cardioneural Signals from Noninvasive ECG Measurement”, *IEEE International Conference on Bioinformatics & Bioengineering (BIBE)*, 2007.

A. Shayan, Y. Zhu, W. Zhang, T.P. Jung, J.R. Duann, C.K. Cheng, “Spatial Density Reduction in the Study of the ECG Signal using Independent Component Analysis”, *International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC)*, 2007.

ABSTRACT OF THE DISSERTATION

**System Level Design of Power Distribution Network
for Mobile Computing Platforms**

by

Amirali Shayan Arani

Doctor of Philosophy in Computer Science (Computer Engineering)

University of California, San Diego, 2011

Professor Chung-Kuan Cheng, Chair

Providing a reliable *power distribution network* (PDN) is a critical design challenge for mobile system on chip platforms. A well-designed power distribution network should be robust enough to support chipset performance while avoids eroding product profit margins through excessive design guardbanding. The solution space between these two requirements is small for PDN designs. On one hand, an inadequate PDN design can lead to test failures, missed performance targets, and intermittent functional problems in the field. On the other hand, some of the more direct PDN improvements such as adding on-die regulators, on-package discrete decoupling capacitors, and package layers increase die and package size, and could cost tens to hundreds of millions of dollars per product line. Mobile platform PDN

design is challenging due to limited form factor, heterogeneous congested blocks with different design specifications, and adoption of multiple low-power techniques and modes. Therefore, it is important to develop a set of PDN design methodologies and analysis tools that can guide the product development from product inception through test and debug.

This dissertation focuses on different aspects of reliable power distribution network design for mobile computing platforms. First, we propose an early-stage power distribution modeling framework to analyze the power distribution during design cycle. We consider the complete closed-loop system from voltage regulation module, printed circuit board, package and silicon die for the co-simulation. Subsequently, an enhanced time domain and frequency domain analysis flow is proposed.

For assessing the performance of the power distribution in the presence of multiple functional power modes, we introduce a worst-case current loading generation. The current generation algorithm synthesizes the functional vector load based on anti-resonance (i.e., resonance-aware) and rogue wave to gain more realistic worst-case voltage variation.

We investigate the impact of power distribution variation on the performance of mobile processors. We estimate the impact of power integrity considerations on low-power processor performance through pre-silicon simulation and post-silicon measurements. We present a predictive performance model under voltage and temperature variations to guide the designers in the early stages of design.

As part of this effort, new emerging technologies for design of power distribution are investigated. A reliability-aware model for 3D stacked chips is developed. The model considers the complete system including *Through Silicon Via* (TSV), substrate noise and stacked dies from both time- and frequency-domain perspectives. Finally, we discuss the recent efficient direction towards on-die regulations for design and optimization of the Linear Dropout based PDN under worst performance. In summary, the complete framework of this thesis aims to provide the means for designing robust power distribution for current evolving mobile computing platforms.

Chapter 1

Introduction

Designing robust power distribution networks has become one of the foremost challenges for performance and reliability of mobile computing platforms. Recent advances in wireless communication bandwidth and development of new mobile applications push for integration of multiple power hungry computing units on mobile *system-on-a-chip* (SOC) solutions such as smart phones. On one hand, the portability requirement of mobile platforms increases the popularity of small form factor and small-scale batteries. On the other hand, advancement in increasingly complex applications such as multimedia and web, which exploit the new higher bandwidth modems (i.e., LTE), requires dealing with power hungry applications with diverse specifications. Limited lifetime of portable batteries and on-chip variations require the adoption of multiple low power modes and techniques to offer longer battery life without sacrificing performance. Thus, a proper strategy for designing power distribution in portable devices is critical in delivering reliable power for a heterogeneous system on a chip.

The power distribution system in portable platforms is comprised of close feedback loop of the battery, *voltage regulation module* (VRM), *printed circuit board* (PCB), package, and on-die grid. Power distribution design is intended to guarantee the reliable performance of the functional units and to eliminate soft and hard failures by providing adequate power. The ultimate goal is to meet the target performance and also eliminate the over-designing cost of the PDN, saving millions of dollars per product line. This thesis tries to tackle co-design of power

distribution units and focuses on different aspects of it that need attention.

Development of efficient analysis tools and methodologies is necessary for both the early stage design of the PDN and also the final signoff on the complete system. PDN methodology and tools evolve gradually from more efficient tools in the early stages to more accurate ones in the final signoff phase. One of the main emphases of this thesis is that these design tools need to assess the performance of the power delivery both in time and in frequency domains. Co-simulation of the time and frequency domains is a key requirement for successful design. From the frequency domain, we can identify the weakness of the PDN loop and identify needed resource location. From the time domain analysis, we evaluate the voltage drop and gain a better understanding of the performance impact of the power delivery design. Co-simulation of both domains is one of the main directions governing this dissertation.

Further, co-design of power distribution units and architectural blocks is required for building a robust bridge between two worlds as performance and technology scales in nanometer eras. For one, considering the mutual interaction of architecture cores and power distribution is important since architecture blocks draw random current from network and as a result we see voltage variation across the chip. The voltage variation across the power distribution will in turns affect the maximum achievable performance of the core. One of the main missions of this thesis is to stress the need for PDN-architecture co-design considerations moving forward. As of today, most of the co-design only involve hardware, architecture and power distribution physical design. Moving forward, it is highly anticipated that we see an inevitable demand for involving the software community for the co-design process.

The rest of this chapter is organized as followed: In Section 1.1, we will highlight the motivation of this thesis. Main challenges facing the designers of mobile system on chip PDN are discussed. Section 1.2 provides an overview of the contributions of this thesis. Problems tackled in this this thesis are briefly outlined in this chapter.

1.1 Motivation and Challenges

Maintaining predictable performance with controlling voltage variation during *functional* and *test* modes is important for all mobile system-on-chip designs. Satisfying this goal requires addressing the following challenges, some of which are exclusively seen by mobile platforms and some are in common with high performance cores.

1.1.1 Limited Area and Form Factor

Portability and form factor impose restrictive area limitations for the design of power distribution in mobile systems. The area restriction stems from different factors:

- i. Final application and use model, i.e., smart phones are trending towards smaller form factors.
- ii. Technology scaling, i.e., metal and power *redistribution layer*'s (RDL) pitch and width are getting smaller in below *22nm* technologies.
- iii. Cost saving results in further area reduction.

Each section of the total PDN feedback loop will be penalized because of the area restrictions in the following forms:

1. Voltage regulation: Off-die regulations such as *switch mode power supply* (SMPS) require large inductance and bulk capacitors. On the other hand, on-die *linear dropout regulators* (LDO) need fair enough real estate for decoupling capacitors allocation.
2. Board routing: Board routing is congested and decap placement is mostly restricted to *top side caps* (TSC) due to the design geometry and form factor limitations.
3. Package: Additional package power routing layer and on-package capacitors are one of the most expensive solutions.

4. Silicon: On-die real estate for decap allocation and metal layers for power routing is very limited.

Overall PDN designers for mobile platforms are dealing with very limited system area and resources and need to make the best utilization for designing a robust power distribution.

1.1.2 System Level Co-design Requirements

One of the foremost challenges of a mobile computing platform PDN is the interaction of different elements of the closed loop power delivery during design cycles. The power distribution network suffers from imbalanced impedance profile with anti-resonance peaks, which are observed due to the interface and boundaries of the PDN structures such as interface of package and board. The overlap of anti-resonance with loading demand of the functional block will lead to failure if the energy spectrum of the current load collides with the anti-resonance frequency. As a result, a large voltage variation is seen in the form of an under-shoot or an over-shoot on the supply or as a ground bounce. Eliminating the interface anti-resonance is not easy since each element of the PDN is ready in a different phase of design.

The physical design of the silicon die is usually ready earlier than the package and board during the design cycle. Thus, the interaction of the PDN loop interface is not possible or very accurate early on. Therefore, PDN system designers should estimate the performance of the system based on the limited history of previous product designs and come up with an estimation methodology to derive the co-design.

In addition to physical design, architecture plays a dominant role in PDN design. Today, concurrent understanding of the PDN characteristics and architecture is mandatory, while conventional PDN designs were fairly independent of the architectural aspects. The co-design is not limited to full loop physical design. Instead, there is a great performance and cost benefit if architecture and software designers assist in PDN design by controlling the peak current value, average current demand, and frequency content. Software and hardware architects should

come up with efficient strategies to ensure certain ON-OFF power mode sequences are limited during functional and test modes. Complexity of power distribution design in current state of the art technology requires a joint co-design effort by the board, packaging, and chip community as well as the architecture and software community to prevent failures during the mutual interactions of PDN and system on a chip.

1.1.3 Cost

There is a fine line between having a cost effective power distribution and failure of the processor due to inadequate power delivery. The main contributors to the power distribution cost are: (i) Silicon area for routing and allocating decap, (ii) Package layers and on-package caps, (iii) Board layers and via formation, and (iv) Regulator design. We can look at the cost of power distribution in mobile platforms from different angles. Possible cost options for designing the power delivery of the mobile system on a chip are:

- I. Max Performance:** Designing a high performance core with maximum achievable performance (F_{max}). This way PDN cost will be part of the performance cost, meaning that higher performance and F_{max} is achievable with additional resource cost of the PDN.
- II. Cost Effective:** The goal of cost effective PDN design is to merely satisfy the minimum performance requirements with a cost effective PDN resource planning.

1.1.4 On-chip Variation

On-chip variation (OCV) is another main challenge ahead of designing mobile computing platforms. As low power processor cores reach GHz range similar to high performance cores, the on-chip variation becomes more pronounced. First, the functional block applications such as multimedia and web are power hungry and draw a substantial current from the PDN network. Further, portable low power processors in GHz are facing the same thermal variation challenges that

high performance cores faced conventionally. The lack of a cooling system such as a fan will boost the spatial and temporal variation of temperature across mobile phones. Also, reduction of the supply voltage further increases the on chip *Process/Temperature/Voltage* (PVT) variations. Process variation is more observable in low voltages. Thus, controlling the temporal and spatial variation of voltage and temperature is one of the main missions of mobile power delivery. The OCV will affect the design margins and performance of the system. To reduce timing failures, designers need to add sufficient timing margins i.e., hold and setup margins and derating factors which introduce additional area and performance penalty. Timing closure of the mobile SOC will be a challenging task. Therefore a unified framework is required to address timing and power integrity simultaneously. The reduction of the supply voltage in the mobile computing system on chip reduces performance and yield significantly. Based on the *ITRS* prediction, integrated circuits will reach Metal 1 (*M1*) half pitch of below $22nm$ on chip, core voltage of below $0.7V$, and frequency of over $10GHz$ in the near future. The circuit delay and performance is more prone to error and thus results in a tighter noise margin. Tight noise margins require special attention with regards to single digit mV voltage drop and thus require complete understanding of the full PDN design.

1.1.5 Heterogeneous System

Power distribution in portable platforms is part of a heterogeneous system on chip. Different functional blocks are integrated together, combining low power cores, high performance processors, digital and analog units on a chip. The memory and core also require different voltage level specs due to $V_{cc_{min}}$ requirement of the memory. The current direction is to adopt split power rails for memory and core. Multiple voltage domains are necessary for satisfying each block's different power requirements. Different low power modes and techniques will increase the non-uniformity in the structure of the PDN. For example, to eliminate the substantial leakage power in below $22nm$, multiple threshold devices such as foot-switches and head-switches are required to be placed across the power grid. Combination of all the above factors will increase the non-uniformity in the power delivery scheme

and make it a more challenging task for system level designers.

The final goal of the system level PDN designer is to consider all the aforementioned trade off factors and design a reliable and portable SOC power delivery network. In the next section, we will introduce an overview of this dissertation's contributions that aim to tackle system level design challenges of the mobile SOC power delivery.

1.2 Dissertation Contributions

In this dissertation, we tackle different aspects of designing reliable power distribution networks for mobile system on chips.

Early Stage PDN System Analysis Flow

In Chapter 2, we introduce early-stage modeling and analysis flow, required to derive the power distribution network system along the design cycle. A complete list of system level issues for designing a robust power distribution is discussed.

- **Early-Stage Modeling:** First, we present initial modeling of the PDN for performing different *what-if* early-stage scenarios.
- **Time-Frequency Co-simulation:** Next, we introduce an efficient power distribution co-simulation methodology to perform time domain simulation and frequency domain analysis based on the package/board resonance behavior.

3D Stacking Power Distribution Analysis

In Chapter 3, we investigate *3D stacked chip* as a potential emerging technology ahead of scaling. Moving towards 3D stacking is one of the directions in mobile system on chip design because of the form factor and technology scaling limitations. The additional dimension is beneficial for integration of the heterogeneous mobile system on chip.

- **TSV and Substrate Coupling Model:** We propose a frequency-based modeling to study the impact of substrate coupling in 3D stacking. The modeling

takes into account multiple tiers and is tailored for co-simulation of time and frequency domains.

- **Reliability-aware Design:** The failure mechanism of TSV and reliability aspects of the 3D stacking structure is elaborated. We propose a reliability-aware optimization flow that combines the joint effect of thermo-mechanical reliability and voltage noise in 3D stacking. The ultimate goal of proposed optimization framework is to minimize the voltage variation across the die and maximize the reliability. The proper resource allocation in the stacked structure is analyzed.

Worst-case Current Prediction

A realistic worst-case current loading is needed to assess the performance of the power distribution network along the design phase. In Chapter 4, we propose a vector-based worst-case current generation methodology. The intention is to estimate the voltage variation under resonance and rogue wave phenomena.

- **Resonance-aware:** A resonance-aware current stimulus synthesis is described. Chances are high that a functional core has a current frequency spectrum aligned with the PDN resonance. Here, we propose a package/board resonance aware current generation flow, based on synthesis of the block vectors in different power modes and frequencies.
- **Vector-based Synthesis of Rogue Wave:** We propose an algorithm to construct the rogue wave envelope based on the functional vectors. The rogue wave current is based on the step response of the system and is intended to highlight different regions of resonance. The algorithm is flexible such that it can preserve different temporal interval windows of the architecture block sequence. The proposed vector-based rogue wave current solution will enable the PDN designers to have a fair assessment of the PDN performance under realistic worst-case loading scenarios.

On-die Voltage Regulation

In Chapter 5, we propose an optimization flow for linear-dropout-regulator-based power distribution design and optimization. The impact of the off-chip voltage regulation is evaluated through an efficient parallel frequency-domain-based flow.

- **LDO based On-die Regulation:** Mobile computing system on chips are no longer efficient when using a conventional passive off-chip regulation method. Multiple power specifications, adoption of power management techniques such as *dynamic voltage and frequency scaling (DVFS)*, *Dynamic Power Gating (DPM)*, and area/performance requirements necessitate a customized on-die regulation. Faster response time and independence from an off-chip passive network make on-die regulation an appealing solution. In this chapter, we propose a linear-dropout-regulator-based (LDO-based) design of the power distribution under worst loading. The proposed methodology is based on identifying the dominant poles and zeros of the system. We calculate the exact analytical step response. The step response and rogue wave current are used to calculate the worst voltage drop. An area and power optimization flow is proposed to minimize the voltage variation across the die.
- **Parallel Flow for VRM Impact:** We propose an efficient parallel flow for assessing the impact of the voltage regulator in the system design. We highlight the significance of inclusion of the voltage regulator in the design analysis for the first time. Failure to do so will lead to optimistic results and final processor failure which is the topic of Chapter 6.

Power Integrity Impact on Performance

Conventional power distribution design process was independent of the performance of the cores. In Chapter 6, we investigate the impact of power distribution on performance. The correlation of analysis with the silicon measurement is discussed in detail. A model for prediction of performance under worst voltage and temperature variation is presented.

- Performance Silicon Correlation: All the analysis and design flows should be correlated with the silicon measurement to be accurate. In this chapter, we estimate the performance of the low power processor through a pre-silicon analysis and a post-silicon measurement. We show that the power distribution can impact the performance up to 15%.
- Predictive Performance Model: We introduce a predictive performance model under worst voltage and temperature variation.

The remainder of this thesis is organized as followed: Chapter 2 provides details of our early-stage PDN modeling and analysis flow. Chapter 3 is the study on 3D integration from reliability, noise and substrate coupling perspectives. In Chapter 4, we propose a vector-based worst-case current generation algorithm built upon resonance-aware and rogue wave phenomena. Chapter 5 discusses our proposed LDO-based design and optimization of the PDN under worst loading. An efficient flow for highlighting the impact of the regulators is detailed. Chapter 6 covers our correlation study on power integrity impact on performance of low power processors through pre-silicon analysis and post-silicon measurement. A predictive performance model is proposed. Finally, Chapter 7 summarizes the main results of this thesis and provides hints for potential future research directions.

Chapter 2

System Level Design and Analysis of Mobile Platforms PDN

Power distribution design is a critical task that impacts chip functionality and cost significantly. System level characteristics of mobile PDN and its design physical challenges are topic of Chapter 2. We describe our proposed analysis flow for PDN co-simulation in both time domain and frequency domain. Strategies for design of a typical mobile ‘power delivery comprised of die, package and board is discussed.

2.1 Introduction

Providing adequate power delivery is a critical challenge for mobile chipsets. A well designed power delivery network is robust enough to support chipset performance, while at the same time not overdesigned to erode product profit margins. The solution space between these two requirements is small for PDN designs. On one side, an inadequate PDN design can lead to test failures, missed performance targets, and intermittent functional problems in the field. On the other hand, some of the more direct PDN improvements, such as adding on-die decoupling capacitors, on-package discrete decoupling capacitors, and more package layers, increase die and package size, and can cost tens to hundreds of millions of dollars per product line. It is therefore important to develop a set of PDN analysis tools

that can guide the product development from product inception through test and debug.

A typical power delivery system is shown in Figure 2.1. Regulated voltage comes out of the *voltage regulator* (VR) on the printed circuit board and is filtered first by the bulk caps and then by *top side caps* (TSCs) and *back side caps* (BSCs). After the supply moves from to PCB to package, it is further decoupled with *die side caps* (DSCs) and on-die capacitance before delivered to transistors to power logic operations. The capacitors are used as local charge reservoirs that are placed incrementally closer to the silicon logic gates. The capacitors are needed because traces and vias on the PCB and package become inductive with frequency and prevent quick changes in supply current. For example, inductive path from back side caps is shortened to PCB vias and package vias and inductive path from die side caps is shortened to package vias and routing. So typically die side caps are more effective than back side caps in a well designed package.

While it is desirable to put as much capacitance as possible close to silicon, spatial constraints are tighter closer to the die and available capacitance decreases due to limits in X , Y , and Z dimensions.

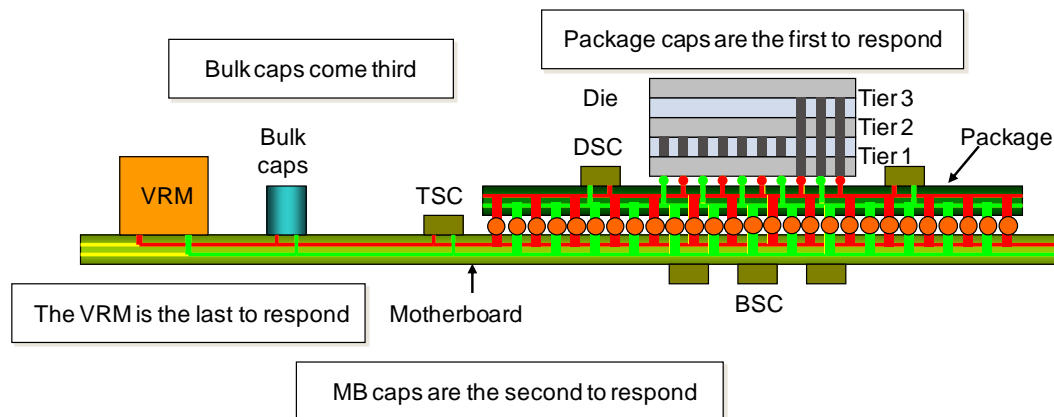


Figure 2.1: Typical PDN system consists of voltage regulator, bulk caps, top side caps, back side caps mounted on the PCB, and die side caps mounted on the package.

Die side caps, for example, are very often limited to small 0201 devices with

capacitance of about hundreds of nF . In some *package on package* (PoP) configurations, Z -height constraints can prevent any die side caps decoupling solution. On-die intentional decoupling capacitors are typically smaller in capacitance and are already taking a considerable fraction of the die. In a typical , for example, decoupling capacitors, fillers, and spacers take close to 75% of standard cell area and only 25% are for logic. Typical total intentional and intrinsic decoupling capacitance available for the core supply on a $45nm$ is around hundred of nF . Figure 2.2

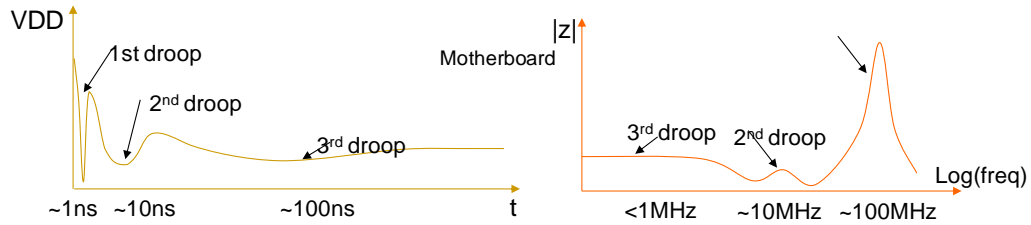


Figure 2.2: Right plot shows typical response of a PDN system to a sudden increase in current demand by the die. Left plot is an illustration of PDN impedance profile.

illustrates the typical response of a PDN system to a sudden increase in current demand from silicon. Due to inductive current bottlenecks in the package routing, charge stored in on-die decaps is depleted first and leads to a sharp drop in supply voltage. This is often referred to as the first droop. The supply voltage recovers when charge in die side caps, back side caps, and top side caps reaches silicon, but drops again when charge in these caps is depleted. This is often referred to as the second droop.

When charge from board level bulk caps arrive, the supply recovers again, but is followed by the third droop before the voltage regulator can respond. Once voltage regulator responds to the current demand increase, the supply reaches a steady state at a lower voltage level with a DC drop determined by the finite VR output resistance and board/package/die resistances. The voltage droops can be seen in a frequency domain impedance plot shown in the plot on the right hand side of Figure 2.2.

The voltage droop events highlighted in Figure 2.2 have direct impact to

silicon performance and power. Logic gate propagation delay is a function of voltage seen at the gate. Therefore, higher voltage leads to faster logic gates and higher maximum operating frequency or F_{max} . From Figure 2.3 we observe, if silicon must function during large current transient events, then the lowest voltage reached must be the voltage used to characterize and sign off designs. During steady state operation, logic gates see higher voltage.

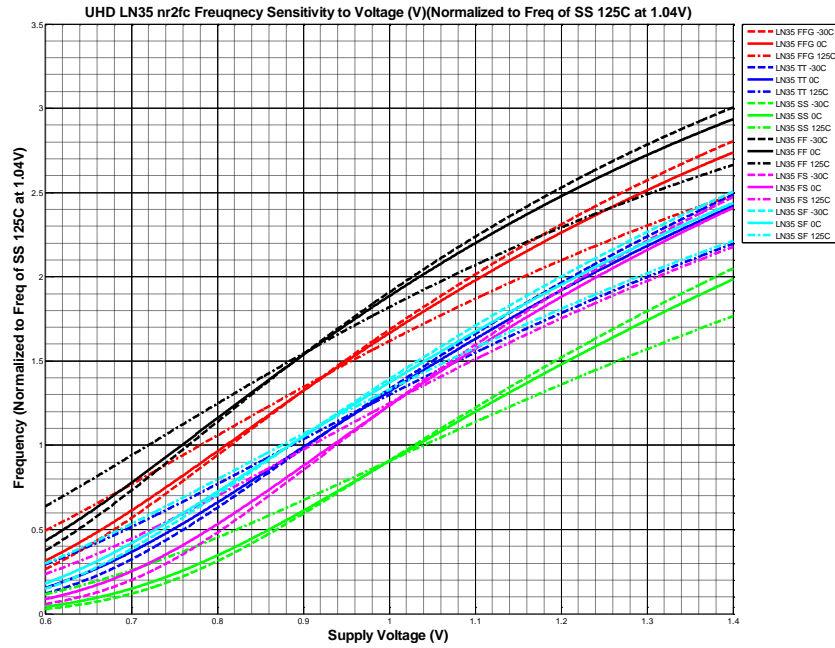


Figure 2.3: Normalized sensitivity of NOR2 ring oscillator frequency to voltage supply variation.

Figure 2.3 illustrates the F_{max} sensitivity to supply voltage for different process and temperature corners, channel lengths and threshold voltages (V_T). Because of $P = C \cdot V^2 \cdot F$, transient voltage drop forces higher supply voltage setting for a given F_{max} . It is therefore critical to design good PDN to minimize voltage drops within the entire bandwidth where current can have significant spectral content.

The rest of this chapter is organized as follows. Section 2.2 contains an

overview of system PDN design. In Section 2.3, time and frequency domain PDN analysis flows are explained in detail. In Section 2.4, PDN noise impact to logic performance is discussed. We conclude the summary in Section 2.5.

2.2 System PDN Design

We discuss system design considerations for PDN in this section. Power delivery spans the entire system, starting from the output of the voltage regulator, through PCB and package power shapes and decoupling solutions, to on-chip power distribution network, and finally reaching the silicon logic gates. PDN modeling difficulties lie in:

1. The size of the entire interconnect system that can reach many centimeters.
2. The varying interconnect features that measure in centimeters on PCB and micrometers on silicon.
3. The distributive and non-uniform nature of the system.

In this section, we will give an overview of design issues at PCB, package, and silicon levels.

2.2.1 PCB PDN Design

In a mobile phone or computer system, the PCB houses both the voltage regulator and the chipset. Significant amount of board routing resources are used for delivering power from the voltage regulator to the chipset. We discuss different components placed on board and their impact on noise here.

Voltage Regulator

A typical voltage regulator needs a compliment of bulk and *multiple layer ceramic chip* (MLCC) capacitors at its output to hold voltage level steady. It also needs a finite amount of time to respond to a load current transient. A load transient that is faster than a few hundred kilohertz is normally outside of the

ability for the voltage regulator to respond. It is therefore also necessary to place faster MLCC capacitors close to chipset power pins to supply charge during fast transient events.

Board Routing Impact on Power Noise

Distance from the output of the VR and the chipset package power pins is an important parameter for PCB PDN design. Long distance between them can lead to large board resistance and DC voltage drop. Due to various board placement constraints, the regulator is usually placed away from the chipset package power supply pins and a power shape is used to delivery power from the VR to the package power pins. One thing to note is that as power shape gets close to package pins, its continuity is greatly compromised by the anti-pads from IO vias. This often leads to high inductance and resistance for the power shape, and care should be taken when developing package pin map.

Decoupling Capacitors

MLCC capacitors are needed to be placed close to the chipset package power pins to supply charge during fast transient events. Effectiveness of the capacitors is determined by the equivalent series inductances from the capacitors to the package pins and through package substrate routing. Logical locations for these decoupling capacitors are on the top side of the PCB immediately adjacent to the chipset package and on the backside of the PCB immediately under the package pins. Top side capacitor placement adds at least a few millimeters lateral distance from capacitor terminals to package power pins, so they must be placed as close as possible to the chipset. For capacitors placed on the backside of the PCB, stored charge must go through PCB vias, which can be between few *mm* in height. Because ground vias do not provide as ideal of a current return path, PCB vias can be a significant source of parasitic inductance. This cuts down the effectiveness of BSCs even though they may be directly benefit the package and the die.

2.2.2 Package for PDN Design

Package is a critical component of a PDN network. Before the onset of a steep rise in the need for computing power, package design had emphasized on smallest form factor and high level of integration. Recently, however, high performance cores were designed in response to the new level of computing demand from smart phones and *mobile internet devices* (MID). Average current and instantaneous current requirements from these cores are significantly higher than other subsystems. To avoid becoming a performance bottleneck, a new PDN aware design approach is being developed for the new generation of high performance mobile platforms.

A PDN-aware design flow should be adapted and attention is paid up front to ensure continuous power shapes are use to deliver power to core. New generation of the high performance mobile cores have very demanding transient current profiles. Thus, care has to be taken so I/O routing will not cut into these power shapes.

2.2.3 Silicon PDN Design

PDN characteristics are affected significantly with silicon floorplan and constraints from physical design. Thus, system design needs to be evolved constantly with timely feedback from PDN during design cycle. Early PDN feedback guarantees that system performance would be predictable with minimum design change and cost. On-die PDN design is challenging because of the following issues:

- Fragmented *Power/Ground* (PG) mesh and cut off islands from analog or digital domains will generate voltage hotspots region that starve from lack of homogenous power delivery (e.g., Figure 2.4). This issue is pronounced more for multiple power domains with split power rails.
- Hard macros under the shadow of broken or segmented RDL, low bump densities and too sparse distance will face increase in the local die voltage hotspots. The adoption of comb structure RDL, double stripe RDL and both direction RDL reduce these local hotspots.

- Unbalanced bump distribution, sparse bump assignment and I/O and power bump allocation competitions will increase die IR drop drastically. We can reduce the noise with PDN-aware floorplanning strategies along the design process.
- Bump impedance map could potentially guide optimum placement of the macros and memories.

All mentioned challenges enforce the need for PDN analysis feedback during evolution of floorplanning and along the design.

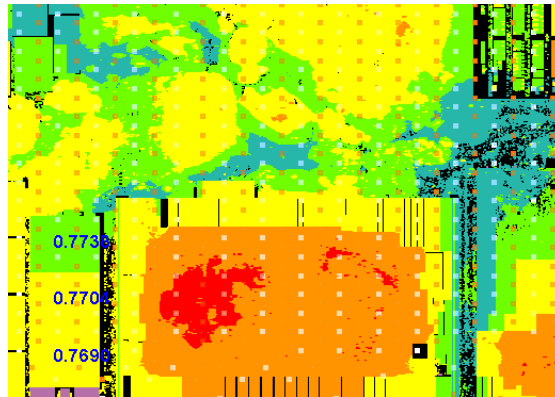


Figure 2.4: Multiple voltage domains voltage hotspots.

2.3 PDN Analysis Methodologies

We explain the application of time domain and frequency domain analysis for mobile PDN design in this section. PDN modeling for mobile chipsets spans from board to die and is complicated because of the following unique properties:

1. Large problem size: Board PDN network is electrically large with power shapes ranging in 10's of millimeters.
2. Wide range of feature sizes: Feature sizes found in a PDN network varies from tens of millimeter on board, to tens of microns on package, to nanoscale on silicon.

3. Fully coupled system: Power delivery shapes and lines are fully coupled in a single system. Accurate modeling of PDN at board/package and package/die resonances is critical to ensure sound design.
4. Heterogenous model: Highly non-uniform die-level capacitive and current loading.

Accurate modeling of a complete PDN system requires tools that are capable of handling highly distributed board, package, and silicon level interconnect parasitics. PDN is also a *linear time invariant* (LTI) interconnecting system with current flows that have tightly synchronized components in time and frequency domains.

2.3.1 Frequency Domain Analysis

Frequency domain analysis for a typical PDN system is discussed here. For a linear time invariant system, it is prudent to analyze it in the frequency domain to find system resonances. Once resonant frequencies are determined, the impedance specifications can be developed (DC impedance, low/mid frequency impedance, and the maximum allowed magnitude of the resonant peak) to meet specified voltage margins. The PDN analysis in frequency domain can also be used to guide time domain analysis, reconstructing worst-case current stimuli to hit the system resonances (resulting in the worst-case voltage drop).

Traditionally, PDN is analyzed when 90% of the system design is done (very late in the design cycle). Unfortunately, at this stage major PDN modifications and fixes cannot be done. A novel design methodology is therefore required to analyze PDN early in the design cycle, permitting several iterations until design is finalized. An early PDN analysis flow in frequency domain has been developed as a part of this thesis, providing vital design feedback before the design is finalized, i.e., pinmap is frozen, RDL is routed, floorplan is placed, and package and board are designed.

The flow is based on electromagnetic field solver, providing impedance profile as a function of frequency. Package is connected to board with several options

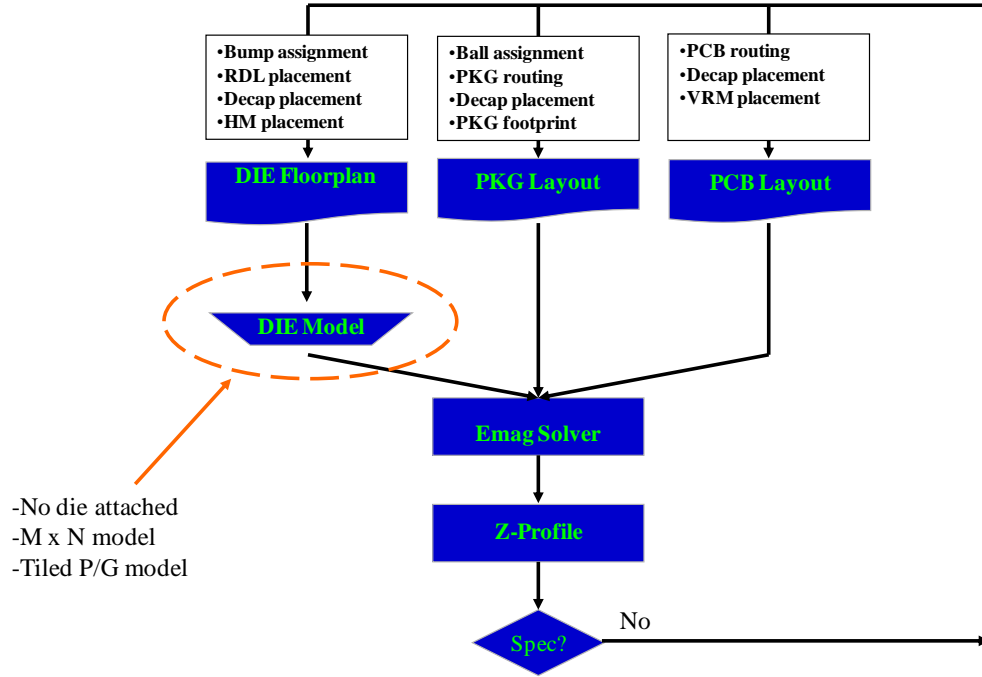


Figure 2.5: Early PDN analysis flow in frequency domain.

available to model die (no die attached, $M \times N$ partitioned model, or tile-based die model). The analysis is iterated until impedance specifications are met. Short turn around time permits several *what-if* iterations per day.

The developed early analysis PDN methodology is summarized in Figure 2.5. The flow is built around electromagnetic field solver. Package is attached to board and impedance as a function of frequency is plotted. It is important to put all major components of a PDN together (VRM, board, package, and die) to accurately capture system resonances.

Despite being easy to use, detailed model becomes available late in the design cycle, when 90% of the design is done. On contrary, tile model is based on *RLCK* parasitic of an on-chip power distribution grid extracted from power grid physical structure. Tile-based die model can be built early in the design cycle when only plan of record for power grid is available.

A tile size of $100\mu m \times 100\mu m$ has been chosen as a tradeoff between accuracy

and computational complexity. The early PDN analysis methodology has been developed to handle power delivery system with two power supplies, i.e., split rails such as core and memory. The developed methodology can be applied to PDN with multiple power supply voltages. A physical structure of a core tile is shown in Figure 3.18. Note that only upper metal layers from M3 are analyzed. No regular on-chip power distribution grid exists on lower metal layers (only local power delivery in standard cells).

RLCK parasitic parameters of core and memory tiles in three directions (left – right, back – front, and top – bottom) are extracted using *Q3D* [2] electromagnetic extractor. Maximum frequency of extraction is set to $1GHz$.

2.3.2 Frequency Domain Automation

We introduce our frequency domain automation flow applied during system design. Scripts have been developed to stitch tiles together based on floorplan, RDL connectivity and package bumps. Three Excel spreadsheets with bumps coordinates, floorplan and RDL stripes are populated and fed to the script. An example of floorplan spreadsheet is illustrated in Figure 2.6. Early stage analysis of the PDN is an iterative process which evolves during design cycles. To synchronize with the latest physical design database and populate the latest PDN models, a frequency domain automation flow is developed. Floorplan information, bump assignments and RDL routings are read from physical design database and updated PDN spreadsheet is generated. The PDN spreadsheet contains user friendly template sheets with the latest updates which facilitates the practice of various *what-if* scenarios.

On-die intrinsic and intentional decoupling capacitors play an important role on defining the die/package anti-resonance frequency. We extract the decoupling capacitors for each block from the hard macros spice models and divide them among the tiles.

The tile-based die model is generated as follows. First, the tile-based model of on-chip power distribution grid is created. RDL with extracted parasitics is then placed on top. Note that either routed RDL database or dummy RDL placement

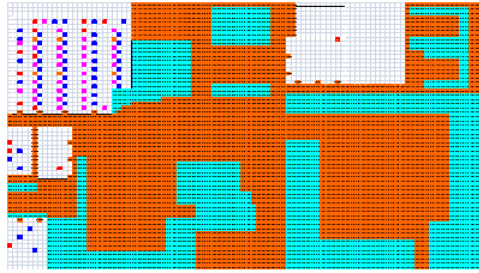


Figure 2.6: An example of floorplan spreadsheet for frequency domain early-stage analysis.

can be used (permitting early *what-if* iterations). Finally, tile-based die model is attached to the package bumps using script based on package pin mapping that contains names and co-ordinates of each bump). On-chip decoupling capacitance (both intrinsic and intentional) with effective *series resistance* (ESR) is also included in the tile-based die model. The total on-chip decoupling capacitance available for each voltage domain is equally distributed among tiles.

Developed methodology has been successfully tested on the industry based mobile testcases. Unacceptably high impedance in critical locations has been determined. The PDN early analysis flow in frequency domain has been correlated with time domain flow based on spice, as discussed in next section. Based on provided feedback, dramatic reduction in power supply noise (voltage drop) has been observed.

2.3.3 Time Domain Analysis

Time domain analysis is one the most computational challenging tasks for the PDN analysis. Analyzing a complete digital device with hundreds of millions of devices (millions of instances) for multiple functional and test modes in time domain is a very complex computational problem. PDN time domain analysis of the chipset is centered on the voltage droops for all the instances and the power mesh on the die contrary to the frequency domain analysis of the PDN which is centered on the package and board components of the PDN. Dynamic voltage

noise analysis of the PDN uses a high capacity transient spice engine as simulation engine and provides full visibility of the voltage noise for all the instances on the die to the designer. Current waveforms at the package bumps and the total current flowing through all the power and ground nets are generated during this transient simulation. Time domain dynamic analysis validates voltage supply assumptions at the device level made during design timing closure, the robustness of power delivery network and the results of the PDN frequency domain analysis. Because we are analyzing voltage and current waveforms, the results of the analysis are easier to correlate and validate with silicon measurements. We have developed a comprehensive time domain analysis flow which is used for PDN optimization and voltage noise signoff. This flow has been developed incrementally starting with the analysis of the die stand alone, followed by the analysis of the die with the package and recently the analysis of the voltage noise on the die with the package and the board integrated in the simulation test bench. The time domain analysis flow has been tested for main mobile platforms.

2.3.4 PDN Optimization Using Time Domain Analysis

Time domain analysis flow requires as inputs real design data (DEF, LEF, .lib, and current signatures) and is suitable for bottom up design flow implementation. Our time domain analysis follows the design implementation stages from floorplan stage to final timing sign off as summarized in Figure 2.7.

Time domain analysis verification flow follows PDN development stages which usually include: die power mesh implementation and optimization, bump assignment and die package analysis and co-design together with package board PCB co-design. We have developed a complete analysis flow which has been tested on industry test cases and includes time domain analysis with *RLC* package model and lumped two ports *S* parameter model as described in Figure 2.8.

The *S* parameter PCB model is generated using frequency domain analysis for the layout PCB. An *S2P* touchstone file is generated and connected to the *RLC* package wrapper. An example of the impedance profile for an un-optimized PCB is presented in Figure 2.2.

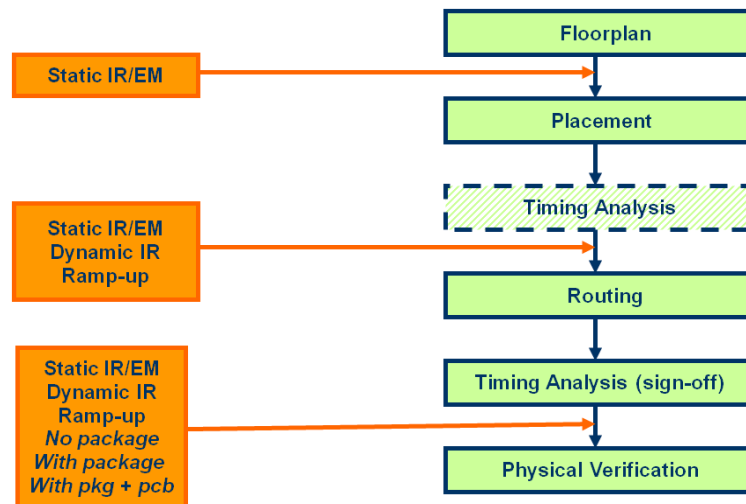


Figure 2.7: Voltage noise time domain analysis integration in the physical design flow.

For the time domain analysis a stimulus is required to drive the simulation. Stimulus can be a VCD vector (FSDB) which is generated by the frontend designers or can be generated using a vectorless analysis driven by the average power of the design.

VCD file can reflect multiple functional modes or test (ATPG) modes of the design. Because the VCD vectors are not easy to generate for the full chip design, vectorless vectors are used most of the time to stimulate the PDN for functional modes analysis.

The results of the dynamic analysis are: voltage droop maps as presented in Figure 2.9.

Based on the static and dynamic voltage droop results the power delivery network is optimized during the design implementation stages from floorplan to the final signoff. This includes PG mesh optimization, bump placement optimization, on die decoupling capacitors optimization and validation and package design. For low power designs rush current analysis is performed for the power and ground gates at the block level and top level.

Due to the fact that simulation is driven by design stimulus, chipset voltage

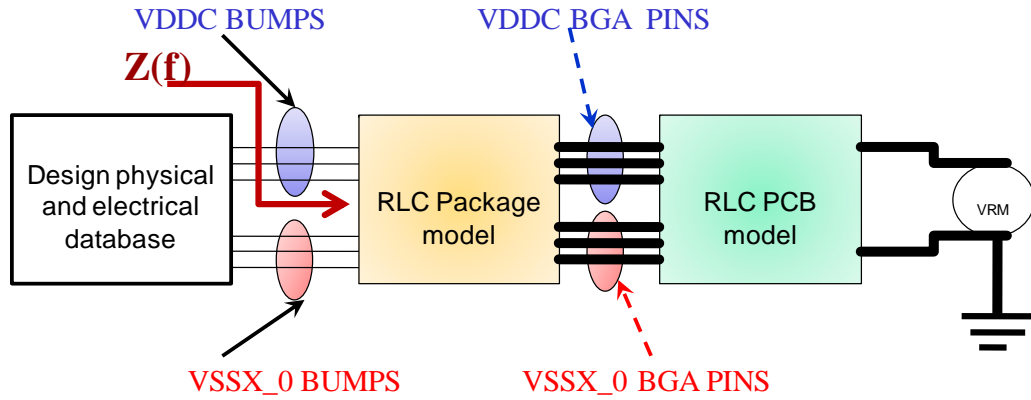


Figure 2.8: Voltage noise time domain analysis for the complete PDN using lumped PCB two ports S parameter model.

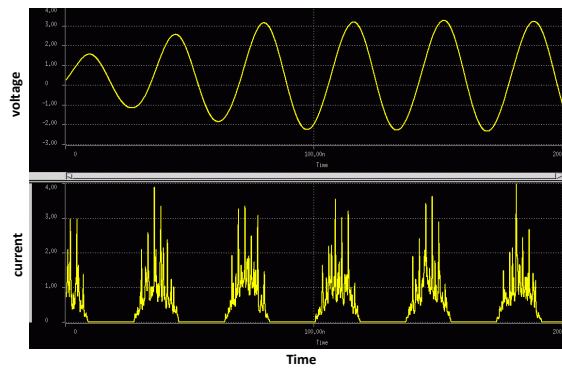


Figure 2.9: Typical functional block load current and voltage drop.

noise failures can be validated using our current analysis flow. Time domain voltage noise analysis in correlation with power analysis drives fundamental design architecture modifications for low power and high performances SOC's designs. Different types of architectures for a SOC design can be optimized to take into consideration PDN requirements, to minimize the cost of PDN and to optimize the system for EMI between the digital and RF subsystems. The main cost factors for PDN are: BEOL metal stack on-die, on-die decoupling capacitors, package substrate, on package decoupling capacitors, number of layers for the board and on board decoupling capacitors. In additions, an accurate time domain voltage

noise analysis can help DVFS SOC implementation and optimization.

2.4 Impact of Logic on Performance

Power supply noise is tightly coupled to the CMOS circuits implementing the core functionality of the chip. We discuss noise impact on performance in this section. Switching CMOS devices produce dynamic noise on the PDN, which in turn impacts the performance of all circuits on the die, including those producing the voltage noise. The performance impact can be described in terms of the primary parameters of the voltage noise: amplitude V_N and frequency F_N .

A low frequency voltage source ($F_N \ll F_{\text{circuit}}$) with amplitude ΔV_N affects the current of MOS devices based on the following expression:

$$I_D \sim k(W/L) \cdot (V_{DD} + \Delta V_N - V_{TH})^\alpha \quad (2.1)$$

where k is a technology-specific trans-conductance parameter, W and L are the channel width and length of the MOS transistor, V_{TH} is the effective threshold voltage, and α is an empirical parameter between Equations (2.1) and (2.2), which describes the relative impact of the overdrive factor ($V_N - V_{TH}$). CMOS circuit delay represents the time to charge and discharge logic and interconnects capacitances, and is correspondingly inversely proportional to the device current expressed in Equation (2.1). The non-linear relationship between circuit delay and power supply voltage can be expressed by Equation (2.2) which is illustrated on Figure 2.3.

$$T_D \sim C_{load}/(V_{DD} + \Delta V_N - V_{TH})^\alpha \quad (2.2)$$

Figure 2.3 shows, reduced power supply voltage increases the sensitivity to voltage noise, which is explained by the squeezed overdrive factor ($V_N - V_{TH}$). This is particularly pronounced on CMOS devices with high-threshold voltage V_{TH} , which are typical for the low power (low leakage) process technologies. For devices with lower threshold voltage V_{TH} , the overdrive factor ($V_N - V_{TH}$) is increased and the impact of ΔV_N voltage noise is reduced.

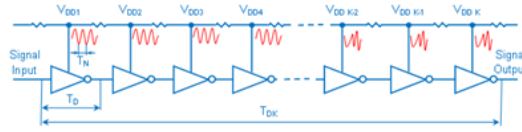


Figure 2.10: Line of CMOS inverters powered by a noisy supply source with frequency F_N .

With the increasing importance of low-voltage low-power operation modes, the significance of voltage noise delay impact is more pronounced and a number of variation-aware timing signoff approaches are being adopted such as *advance on chip variation* (AOCV) to address these challenges.

The impact of voltage noise frequency F_N on circuit delay is less documented and correspondingly more interesting to explore. For illustrative purposes, the impact of a sinusoidal power supply source on a delay line of K CMOS inverters is described as shown on Figure 2.10.

The phase of the voltage noise is considered random with respect to the arrival time of the input signal into the delay line. A number of parameters and their relations characterize the impact of the power supply noise on the delay of the inverter chain. The relationships between the insertion delay of single delay element T_D and the overall line delay of K inverters T_{DK} to the voltage noise period T_N ($T_N = 1/F_N$) define a number of cases to be considered:

- 1. $T_N \ll T_D$:** The voltage noise has no significant impact on the inverter unit delay as the circuit switches over a period of time incorporating many noise periods, effectively canceling the noise impact.
- 2. $T_N \gg T_D$:** The voltage noise is purely random and constant for the duration of the CMOS gate switching. This case was described at the beginning of this chapter and the dependency is illustrated in Figure 2.3. Supply voltage impacts circuit performance and voltage noise sensitivity.
- 3. $T_N \approx T_D$:** Voltage noise period is on the order of the circuit delay, in which case the impact on a single inverter can be considered random, but the correlation to the delay of the following circuit is very strong. This scenario is illustrated

on Figure 2.11. Each individual inverter (unit delay circuit) experiences a unique voltage noise impact, which may be described by an effective power supply voltage V_{eff} , defined as the steady supply voltage at which the same inverter would have the same delay as with the noisy power supply. Approximate values for the effective voltages of the stage delays are illustrated in Figure 2.11. During the switching period of an individual inverter, the effective voltage is constant and the delay impact is as shown in Figure 2.3. Supply voltage impact on circuit performance and voltage noise sensitivity. An alternative approach to consider the effective stage voltages is as weight factors on the stage inverter delays. It is important to note, however, that along a delay line of multiple similar circuits, the average V_{eff} tends to cancel out the individual stage variations and the mean delay $E\{T_{DLK}\}$ approaches the ideal (clean power) delay of the line $K \times T_D$. Therefore, a long line of similar circuits becomes power supply noise insensitive when $T_N \approx T_D$ and K is large, or equivalently $V_{DLK} \gg T_N$. Unfortunately, this desirable effect is not valid when the delay line is composed of various circuits with wide variation in their circuit characteristics and insertion delays. In that case, the relative phase of voltage noise waveform with respect to the input signal is important since it defines the individual effective stage voltages, correspondingly defining the weight factors of the different delays along the chain. Since the phase of the noise waveform is generally considered random, the impact of that noise on a line of diverse circuits is also random and correspondingly requires additional timing margin.

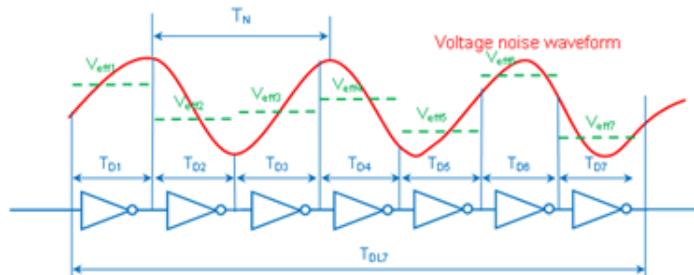


Figure 2.11: Averaging of effective V_{dd} over a long line of inverters.

Above discussion about voltage noise frequency impact was limited within the assumption of constant noise amplitude. Typical on-chip power supply noise exhibits significant variations in terms of both amplitude as well as spectral content as shown in Figure 2.9. The peaks of the noise waveforms typically represent the moments where large number of storage elements switch synchronously based on a clock edge. The pseudo-random nature of these noise variations produces corresponding random circuit delays and may lead to potential silicon timing failures. As it is difficult to comprehensively model all possible voltage noise scenarios, a conventional design practice is to apply a timing uncertainty margin to provide signoff robustness. Alternative approaches are currently being explored to reduce the timing margins by more intelligent utilization of the power grid information to narrow down the timing signoff assumptions on voltage noise amplitude and spectral content.

To provide some perspective, a typical inverter delay with ideal nominal power supply in $45nm$ CMOS technology is around $\sim 20ps$. Based on the discussion above, high-frequency noise above $200GHz$ ($T_N = 5ps$) would not impact the circuit delay significantly as several noise periods would fit within the switching time of the inverter and effectively cancel out. Voltage noise with frequency $F_N < 10GHz$ ($T_N > 100ps$) produces random delay variation based on the random phase alignment of the noise waveform with respect to the input signal. As illustrated in the previous section, however, today's integrated circuits have power grid noise in the range $0 - 2GHz$ with significant spectral energy contained within $100MHz$. For such low-frequency voltage noise, typical data and clock paths on the order of $1-10ns$ exhibit pseudo-random delay variations, which are hard to model deterministically. The impact of power supply noise on such paths of multiple logic gates may be modeled as correlated random variables based on the PDN extracted information. Temporal and spectral analysis should be combined to describe the probability density functions of these random variables and provide a comprehensive perspective on the realistic voltage noise impact on timing delays.

2.5 Summary

In this chapter, we provided an overview of the PDN system level design and related issues. We then outlined a frequency domain flow and a time domain flow we developed to accurately analyze PDN characteristics. Both of the frequency domain and time domain flows were evaluated on the industry testcases.

Chapter 3

Design and Reliability Aspects of 3D Integration PDN

Power distribution network design for stacked dies faces more reliability challenges compared with conventional SOC. Substrate coupling noise among TSV and PDN grid, thermo-mechanical stress and electromigration are pronounced more as a result of stacking. We detail a comprehensive modeling of TSV and stacked power grid with frequency dependent parasitic in Chapter 3. The analytical model considers the impact of the substrate coupling between the TSVs and in each die. A frequency domain analysis flow is introduced that incorporate frequency dependent parasitics. Design of a reliable power distribution network is formulated as an optimization problem to minimize power noise under reliability and electromigration constraints. Experimental results demonstrate the efficacy of the problem formulation and methodology.

3.1 Introduction

Early power delivery planning is a crucial 3D chip design step as noise margin becomes tighter in below $22nm$ technology with high clock frequency. It is important to perform this task early on along with TSV allocation in order to prevent from logic and chip failures [7].

About 30-40% of the available die to die vias in 3D chip are allocated

for power delivery and the rest are assigned for signals [60]. The supply current flows through the inductive solder bumps and narrows through silicon vias with considerable inductance parasitic, and results in significant *simultaneous switching noise* (SSN). Stacked chip requires less than the original 2D design footprint and as a result current density per pin is larger. Thus, power distribution network in 3D systems needs to be accurately modeled and designed for the higher current density. Huang *et al.* in [22] proposed an analytical PDN model for 3D stacking considering via inductance and stacked decoupling capacitors. However Huang’s model does not include the effect of on-chip inductance. As the clock frequency reaches GHz range, di/dt event and SSN from the switching dice will be as significant as the IR drop. Therefore, in the GHz range clock frequency on-chip inductance in the PDN model need to be considered for a realistic voltage drop assessment.

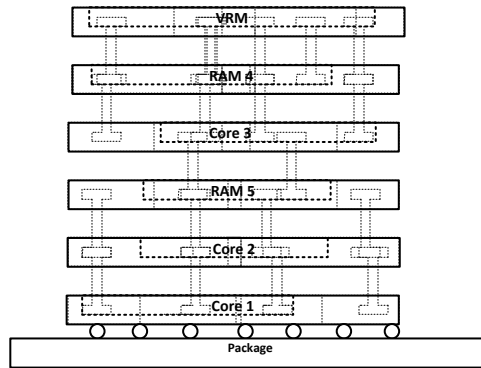


Figure 3.1: 3D PDN stacked core and memory model.

A 3D PDN model is presented considering the effect of on-chip inductance in chapter 3. On-chip dc-dc converter is adopted as a potential alternative to off-chip VRM regulation. Higher circuit density in 3D PDN results in even larger circuit netlist. Larger netlist makes the simulation of huge 3D PDN a challenging task. Each time domain sample point of the transient simulation depends on the results from previous time points. Thus, it will take a long time to simulate whole large networks sequentially. A parallel flow that adopts *Message Passing Interface* (MPI) on clustered Linux machines is developed which significantly reduces simulation time.

3.2 Overview of Reliability Challenges

Three dimensional integrated circuits introduce a technology potential to enhance the performance, functionality and device packaging density. Stacking enables integration of heterogenous functional blocks such as digital and DSP cores, memory, RF and analog modules [22, 27]. Having multiple active layers stacked on top of each other provide more flexibility for design to reduce cost and power. The increase in current demand and faster switching frequency of the stacking, introduce more severe power integrity and reliability challenges compared with conventional system on-chip. Substrate coupling modeling is a key issue in 3D due to its significant impact on the performance of the analog circuits in the stacked chip. This thesis present a comprehensive frequency dependent model for through silicon via and stacked power layers with substrate coupling in between.

We propose a reliability-aware TSV planning considering thermo-mechanical stress. Our model considers the impact of frequency dependent TSV and substrate parasitics on the voltage variation. We extend the reliability model to include *electromigration* (EM) constrains in the analysis. We trade block-out space in favor of decoupling capacitor area in this problem formulation. We propose in this chapter a unified methodology to address the problem of designing optimal 3D stacking under reliability constraints. One typically with the following properties and constraints:

1. Minimum IR drop across the power grid in layers and including substrate coupling among the TSVs and tiers.
2. Satisfy electromigration maximum current density constraints.
3. Maximizing lifetime and thermo-mechanical reliability of the chip and finally.
4. Maximizing area for routability and decoupling capacitors allocation.

Electromigration [25] and thermo-mechanical [45] stress are the root cause of major long term failure problems in 3D stacked ICs and is the focus of this chapter. For the second part of this chapter, we first, present our TSV and stacked grid electrical modeling scheme that considers substrate coupling. Few benefits of

frequency domain are the the fact that we can incorporate frequency dependent parasitic and that we can highlight resonant peaks. Our extracted model demonstrates variation of peak in frequency spectrum which is not observed in fixed parasitic model. Time domain transient response is then recovered with vector fitting [47]. We analyze power noise and reliability of stacked grid here. Design of stacked power delivery is defined as an optimization formulation to obtain minimum noise under reliability constraints.

Rest of this chapter is organized as follows: Parallel 3D PDN analysis flow is introduced in Section 3.3.1. Modeling for 3D power delivery is discussed in Section 3.3.2. In Section 3.3.3, we will discuss our frequency domain analysis result for the 3D PDN based on the grid design parameters. For the second part of this chapter, we discuss reliability-aware analysis of the 3D die stacking: Details of 3D power distribution model are presented in Section 3.4.1. Substrate coupling model among TSV and each tier is extracted as a frequency dependent electrical model. We discuss our frequency domain based flow which is applied for power integrity analysis of 3D PDN model in Section 3.3.1. In Section 3.5 electromigration and thermo-mechanical reliability aspects of stacking is elaborated. Section 3.5.3 presents our formulation for design of a reliable stacked IC power distribution along with experimental results; and finally Section 3.6 concludes summary of this chapter.

3.3 Time and Frequency Domain Analysis of the Three Dimensional Networks

3.3.1 Parallel 3D PDN Analysis Flow

Design of 3D power delivery requires a custom analysis flow with proper model of TSV and grid as an input that guides designers to make right design choices. We discuss our co-simulation flow that can highlight resonance peaks and maximum voltage drop of stacked die. Most of the conventional simulators fail to simulate entire system in a reasonable time because of the large scale of the

3D PDN and special inputs. Figure 3.3 depicts our efficient parallel processing package for analysis of the entire 3D power distribution network in both frequency and time domain.

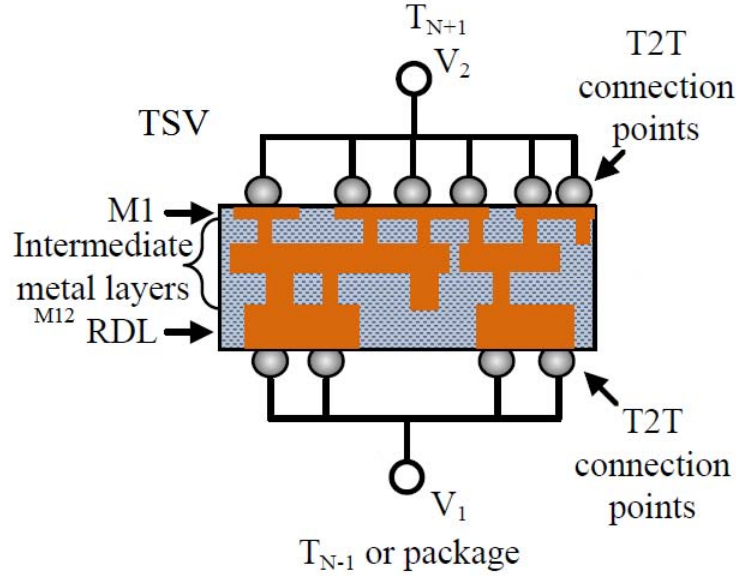


Figure 3.2: Tier to Tier (T2T) power connection model of 3D stacked die.

We apply *vector fitting* (VF) technique to convert back frequency domain results into time domain. We are in an era that clock frequency of the SOC is reaching GHz range and as a result SSN from di/dt cannot be ignored. To solve the noise in the stacked die, the frequency dependent $RLGC(f)$ of the entire system is inserted into iterative linear solver and is converted back using vector fitting. Conventional vector fitting technique [60] has unacceptable error margin and the recovered time domain result has a large error (ΔVF) compared with HSPICE:

$$\Delta V_{VF} = V(t) - V_{VF}(t) \quad (3.1)$$

where ΔVF is the deviation of fitted time domain approximation from the original time domain signal. We enhanced vector fitting algorithm to reduce the error. We apply the remainder of the vector fitting ΔVF iteratively into VF process and perform the vector fitting process until ΔVF reaches acceptable error margin rate of 10^{-16} .

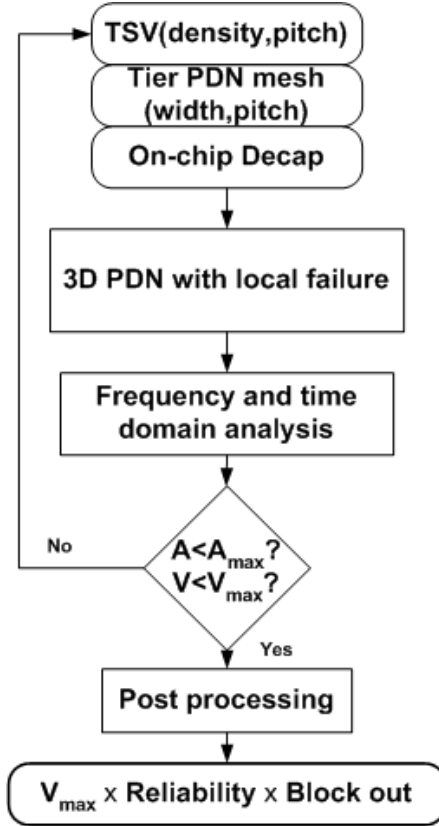


Figure 3.3: Stacked PDN reliability-aware and noise co-analysis flow.

We implement our parallel flow using MPI and run the flow on FWgrid infrastructure (<http://fwgrid.ucsd.edu>) with multiprocessors to speedup simulation time. We run the parallel flow on clustered Linux machine and our results demonstrate speedup of up to $22\times$ times with single processor and more than $430\times$ times by using up to 200 processors over HSPICE [20] transient simulation. Parallel processing reduces PDN simulation time significantly from hours to less than hundreds of seconds.

3.3.2 Design Planning for the 3D Power Grid

We explore efficient via placement to reduce the power noise in 3D structure. Power/ground wires in SOC are routed by using orthogonal interconnection. Peak and average current increases with increase in high performance processing demand. Current flows through the inductive solder bumps and through silicon

vias and results in severe SSN and power integrity noise. Die to die through silicon vias serves both for signal routing and power delivery. Dummy vias are adopted as heat conductance between the stacked die [42] to balance temperature hotspots.

We study on-die regulation via stacked VRM solution in 3D model. We analyze the noise impact of placement of the on-chip regulator in different tiers of the stacked die from bottom to center.

3D Power Distribution Model

We model power delivery with a stacked 3D RLC grid, as shown in Figure 3.4. In each tier of stacking a $10 \times 10 \text{mm}^2$ active die is modeled with mesh RLC network and the current sink stimuli in the center, representing the active switching transistors. Through silicon vias with the length of $200 \mu\text{m}$ connect each layer of the stacked dies and is modeled as a RL in series. RL values are based on the via process technology. TSV inductance will cause SSN and need to be accurately modeled. We detail our on-chip VRM allocation in next section.

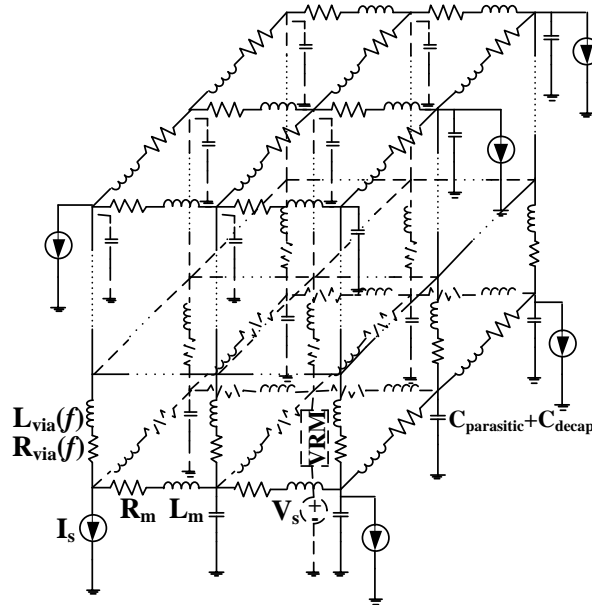


Figure 3.4: 3D PDN parasitic $RLCK$ model with package and VRM.

On-chip Voltage Regulator

To reduce di/dt event that causes supply voltage variation, an integrated on-chip voltage regulation in 3D stacked IC is adopted. In conventional power delivery systems, mounted VRM on the board has critical challenges: long interconnect path, parasitic inductance of the package-board which generate SSN, large decap size and large number of the power/ground pins required by chip. Mentioned requirements take expensive area of system and make packaging design more complex. Stacked dc-dc converter in 3D structure benefits from minimum interconnect parasitics, wider bandwidth and easy distribution to multiple domains [51].

It is practical to integrate the regulation circuitry and use small, discrete external inductors mounted close to the die as switching mode power supply. A large number of these external inductors would be required, at least one per PDN node. Inductance component can be embedded in the package as well. In addition, these physically small components have limited inductance [7]. We use the close loop impedance of the VRM model as input in to the analytical flow and simulate full system.

3.3.3 Frequency Domain Analysis of 3D PDN

We discuss our experimental results from the analysis using our proposed 3D PDN model. The sensitivity of the through silicon vias and PDN grid to different parameters are explored in this section.

Through Silicon Via Distribution Density

Each tier PDN mesh is divided into equal quadrant based on the power consumption spatial distribution and pitch. We allocate two via to the corner of each quadrant in the PDN simulation. The via distribution is uniform in each tier. We vary via distribution density from 10% to 80% of the total allocable locations in grid. In Figure 3.5, we observe that for low frequency higher via density results in less output impedance, but in high frequency the PDN with higher via density has larger impedance peaks. This is because reduced resistance increases the Q

factor, which leads to higher anti-resonance peak.

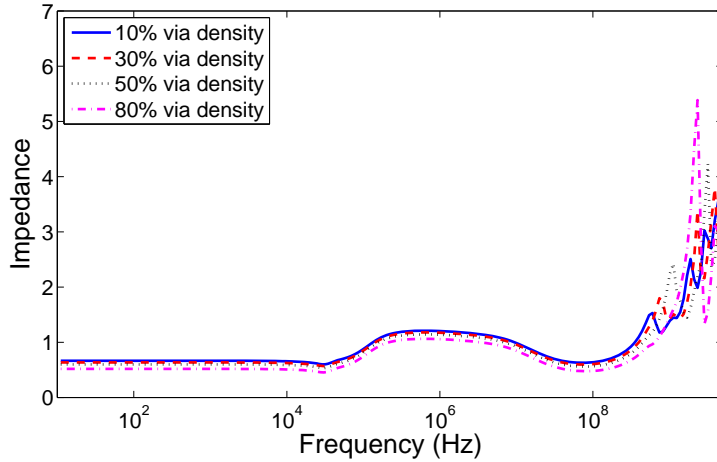


Figure 3.5: Impedance spectrum of 3D stacked PDN versus via distribution density.

3.3.4 On-chip Inductance Scaling in 3D PDN

Figure 3.6 illustrates our 3D PDN impedance spectrum with and without on-chip inductance. For high frequencies in the range of GHz the impedance increases up to 8 times with on-chip inductance compared to RC on-chip power grid. This tells us that on-chip inductance should no longer be ignored when operation frequency increases to GHz .

We scale 3 orders of magnitude on-chip inductance in the 3D P/G mesh. Low frequency impedance magnitude of all three scales are same but as the frequency increases the impedance magnitude is increased up to $10\times$ as depicted in Figure 3.7.

3.3.5 On-chip Resistance Scaling in 3D PDN

Figure 3.8 shows that as the resistance of the grid increases in low frequencies, 3D PDN impedance magnitude is increased towards DC. This can be explained by the fact that as the frequency increases the damping effect of the R becomes dominant and compensates with increasing low frequency impedance.

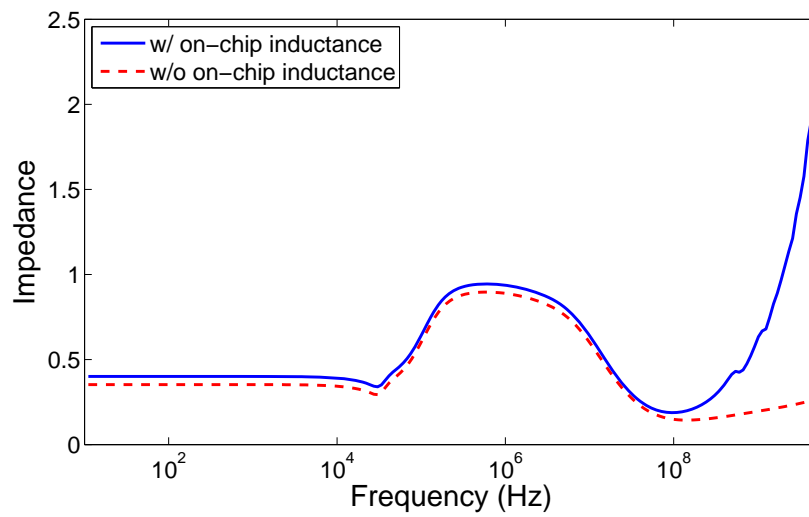


Figure 3.6: Impedance magnitude of 3D PDN with and without on-chip inductance.

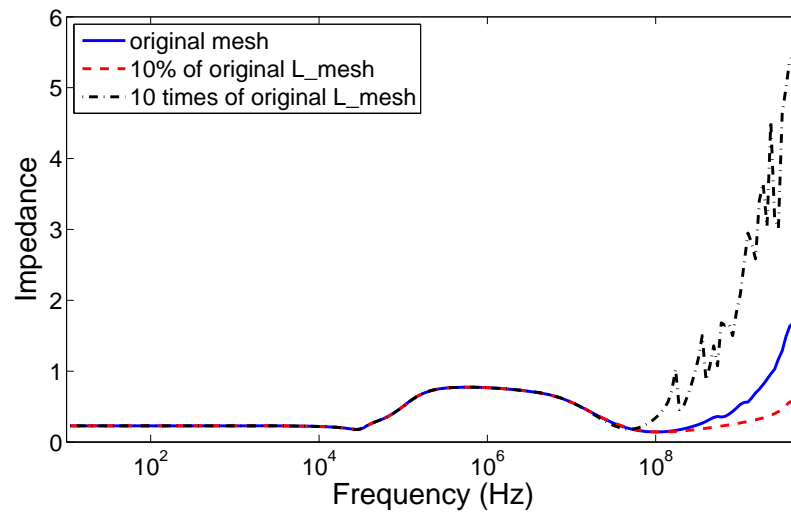


Figure 3.7: On-chip inductance scaling impact on impedance magnitude in 3D stacking.

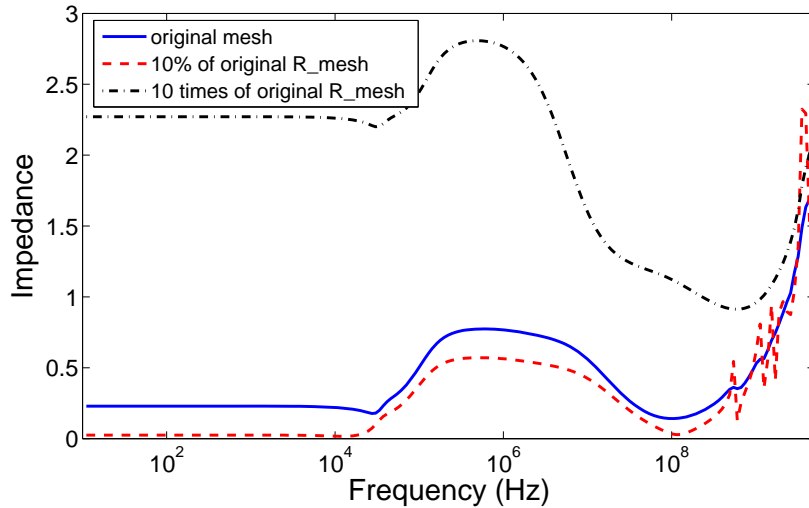


Figure 3.8: Resistance scaling impact on impedance spectrum in 3D stacking.

3.3.6 Sensitivity of Impedance to Stacked Layers

We show how increasing the number of stacked layers from 5 to 10 affects impedance magnitude and voltage noise. We run the flow on 3D power testcase with x and y dimension of 32×119 nodes in the grid and 5 to 15 stacked layers. We observe in Figure 3.9 that the impedance magnitude peak is reduced with more stacked layers. Figure 3.10 shows simulation time versus number of processors for both cases on FWgrid machine with up to 200 processors where simulation time is reduced significantly to less than a minute. Same simulations of case 1 and 2 in HSPICE took $9911sec$ and $110479sec$ respectively for $300nsec$ duration (Table 3.1).

Table 3.1: Simulation time of 3D PDN versus number of processors on clustered FWgrid multiple cores for proposed flow versus HSPICE

sec	HSPICE	No. of processors						
		1 p.	4 p.	16 p.	32 p.	64 p.	128 p.	200 p.
case1	9911.32	1531.80	342.20	136.70	68.70	38.57	21.35	23.01
case2	110479.36	2635.06	808.90	214.40	126.50	77.00	52.30	28.75

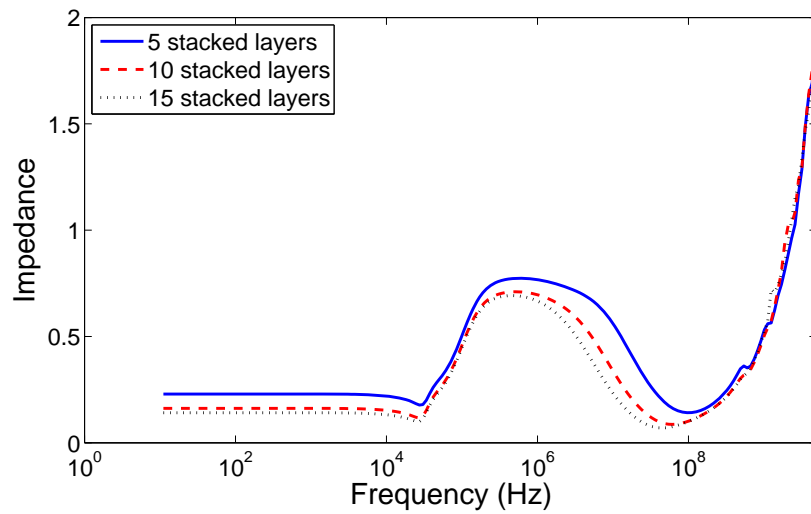


Figure 3.9: Impedance spectrum scaling with different stacked layers.

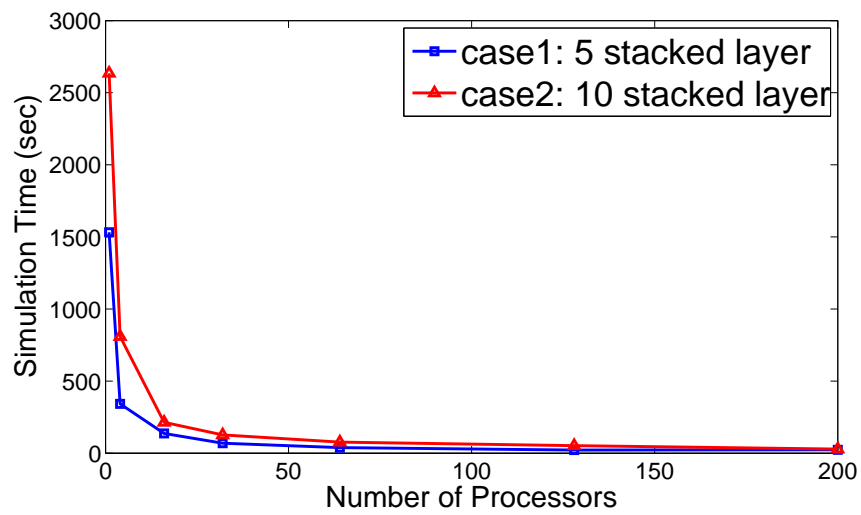


Figure 3.10: Simulation CPU time as a function of number of processors for 3D stacking.

Analysis of 3D PDN in Time Domain

Higher via density will result in lower voltage noise as Huang *et al.* concluded in [22]. However, this is not always the case. In Figure 3.5, we can see that for a certain via placement strategy, higher density may lead to lower impedance peak below MHz . In addition, it may also result in higher resonance peak at higher frequencies. This is because more vias in parallel reduce effective resistance, and thus increase the quality factor Q of the RLC tank which makes anti-resonance more pronounced. As a result, the output voltage noise depends on the frequency spectrum of the input current. Figure 3.11 shows the output voltage noise for different via densities. From Figure 3.11, we can observe that higher via density may result in larger voltage noise.

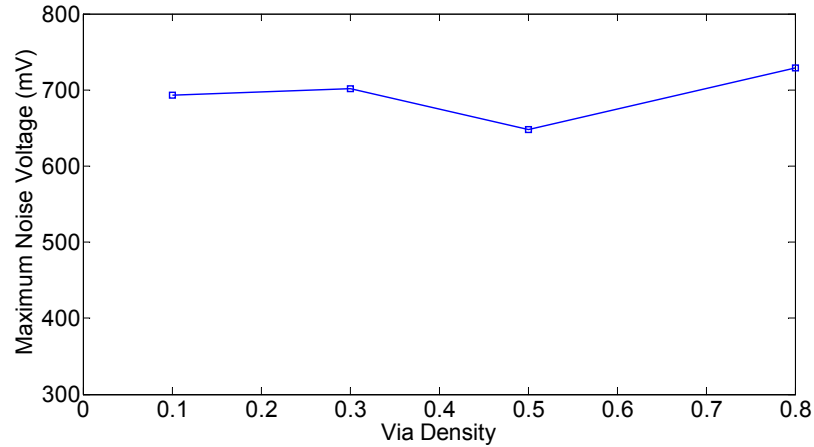


Figure 3.11: Time domain noise magnitude of 3D stacking versus via density.

3.4 Reliability-aware 3D PDN Model Considering Substrate Coupling

3.4.1 Substrate Coupling Model in 3D PDN

Substrate coupling in 3D stacking has dominant impact on logic performance. We present our modeling scheme for through silicon via in presence of the substrate coupling. The key technology that enables stacking of multiple dies in *system in package* (SiP) is through silicon via. Heterogeneous layers with different functionalities and high operation frequency can be stacked densely together. Silicon substrate is usually of low resistivity in digital SOC. Thus, silicon substrate propagate noise during system switching and result in huge noise.

None of the previous works on 3D power distribution [22, 51] to the best of our knowledge, model and consider the impact of substrate coupling in power noise. Three dimensional stacking parasitic interacts through the shared silicon substrate and among TSVs. Each time transition and switching occurs in the digital circuits, they inject noise into three dimensional stacked packages and among TSV. Silicon substrate is conductive and power noise easily propagates to analogue circuits destructively. This impact the performance due to deficient immunity of digital circuits [52]. Through silicon via model and extraction method are described in the next Section 3.4.2.

3.4.2 Frequency Dependent Through Silicon Via Model

At high frequencies current does not flow uniformly across the cross section of the conductor. Structure of through silicon via is shown in Figure 3.12. Instead because of skin effect current becomes increasingly concentrated near the edge of surface. Current density varies with respect to frequency. As a result, the effective cross section area is reduced with the onset of skin effect, causing the effective resistance at high frequencies to increase. Previous works [39, 4] addressed modeling of TSV for fix frequency and without consideration of the substrate coupling effect. Our proposed modeling takes into account frequency dependent

parasites and substrate coupling more accurately. TSV is electrically modeled and extracted using High Frequency Structure Simulation (HFSS [1]) full-wave solver from Ansoft that considers substrate loss of TSV.

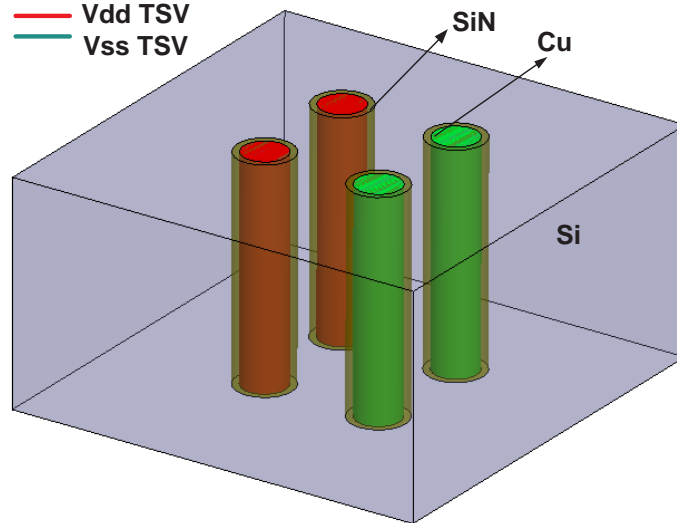


Figure 3.12: Through silicon via modeling in HFSS [2].

The following properties are assumed for the TSV structure modeling,:

1. TSV is modeled as an array of V_{dd} and V_{ss} pair in a $200 \times 200 \mu m^2$ grid for each tier. TSV has a silicon shared substrate through wafer via hole with insulation barrier formed on the substrate and via sidewall.
2. Substrate is a low resistive ($10\Omega \cdot cm$) silicon with small loss tangent (less than 0.0005). Substrate loss could be mitigated by using high resistivity or thick dielectric as an expensive packaging alternative solution [19, 37] based on design constraint.
3. TSV height is $100\mu m$.
4. The thin dielectric layer (silicon nitride) surrounding TSV is $0.2\mu m$ thick with $\epsilon_r = 7$.
5. TSV and substrate model are extracted over a broad range of frequency from $50MHz$ to $10GHz$.

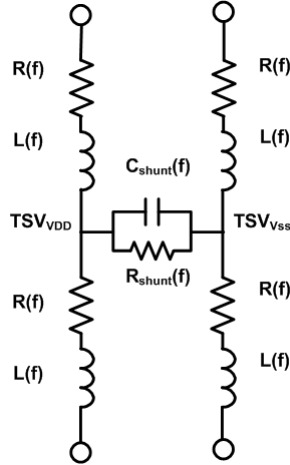


Figure 3.13: TSV *RLGC* equivalent *RLCG(f)* model with substrate coupling.

We look into sensitivity of noise to dimension of through silicon via and the density of the TSV pitch varied. We analyze the impact of these parameters to reduce the power noise and to maximize reliability. *S*-parameters model of the V_{dd} and V_{ss} pairs are obtained from the extraction. *S*-parameter of the TSV is converted to equivalent *RLGC* electrical model as described in Section 3.4.3.

3.4.3 Frequency Dependent Substrate Coupling Model

The extracted *S*-parameters model for TSV is converted to equivalent fitted electrical model as shown in Figure 3.13. We follow [12] technique where γ is the propagation constant and Z is the TSV transmission line based characteristic impedance:

$$\gamma = \sqrt{(R + j\omega L) \cdot (G + j\omega C)} = \alpha + j\beta \quad (3.2)$$

$$Z = \sqrt{\frac{(R + j\omega L)}{(G + j\omega C)}} \quad (3.3)$$

From Z and γ , equivalent fitted electrical *RLGC* model is derived. The model is suitable for frequency domain analysis flow as follows:

$$R(f) = \text{Re}\{\gamma \cdot Z\} \quad (3.4)$$

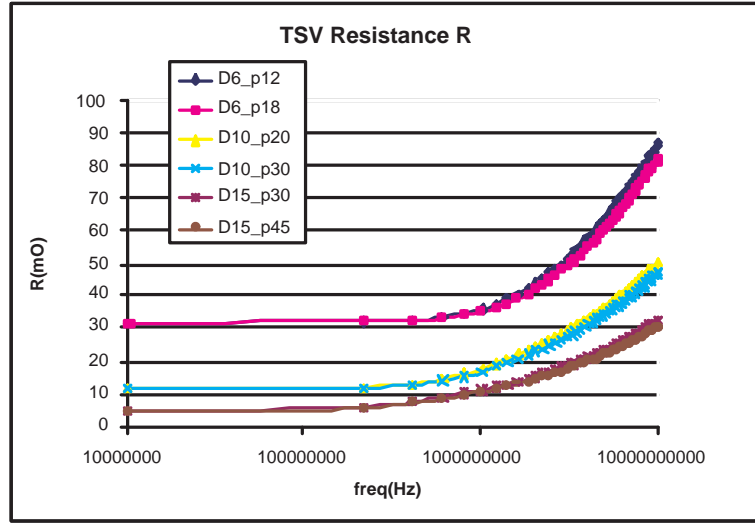


Figure 3.14: TSV equivalent resistance R (D =diameter, P =pitch).

$$L(f) = \text{Im}\{\gamma \cdot Z\} / \omega \quad (3.5)$$

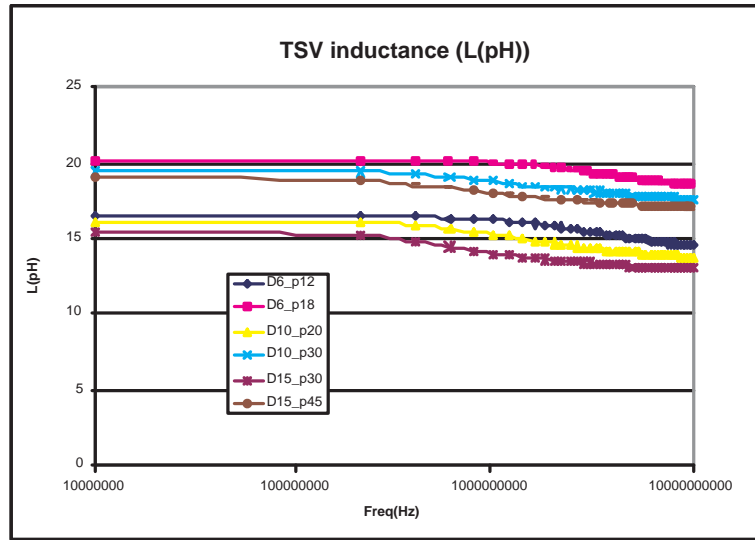


Figure 3.15: TSV equivalent inductance L (D =diameter, P =pitch).

$$R_{\text{substrate}}(f) = R_{\text{shunt}}(f) = \text{Re}\{\gamma/Z\}^{-1} \quad (3.6)$$

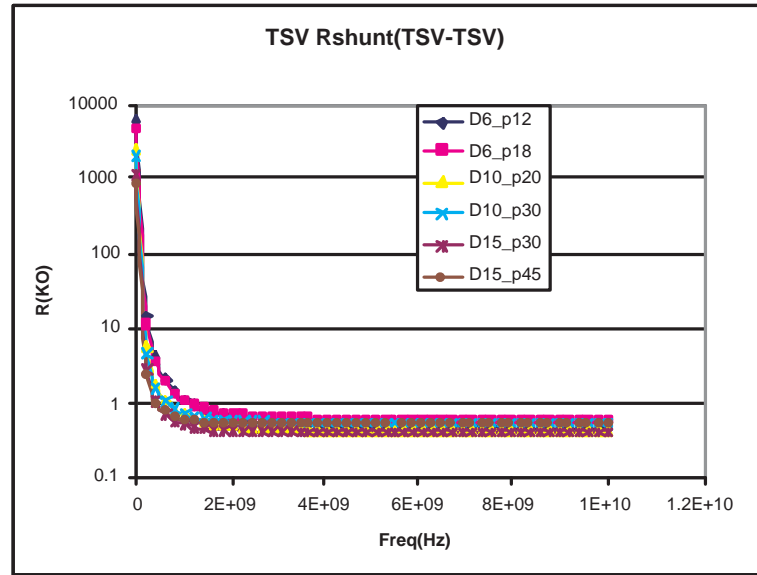


Figure 3.16: Substrate equivalent shunt resistance (D=diameter, P=pitch).

$$C_{substrate}(f) = C_{shunt}(f) = Im\{\gamma/Z\}/\omega \quad (3.7)$$

Figure 3.14 to Figure 3.17 depict complete set of TSV extractions for different combination of TSV geometry and density. TSV RLC parasitic components are then concatenated to represent stacked die. In our testcase a $1 \times 1mm^2$ die is partitioned into $200 \times 200\mu m^2$ uniform sections.

3.4.4 Frequency Dependent Power Grid Model in Each Tier

Power metal parasitics in each layer of stacked PDN is extracted for a $200 \times 200\mu m^2$ geometry. Each metal layer is concatenated as a lattice grid electrical model front to back and left to right. Q3D extraction simulator [2] is adopted to extract the frequency dependent parasitic of each metal layer as shown in Figure 3.18. Two dies are stacked together with TSV in between.

Figure 3.15 illustrates that inductance sensitivity depends on the pitch and diameter aspect ratio ($pitch/diameter$). TSV resistivity is function of its diameter

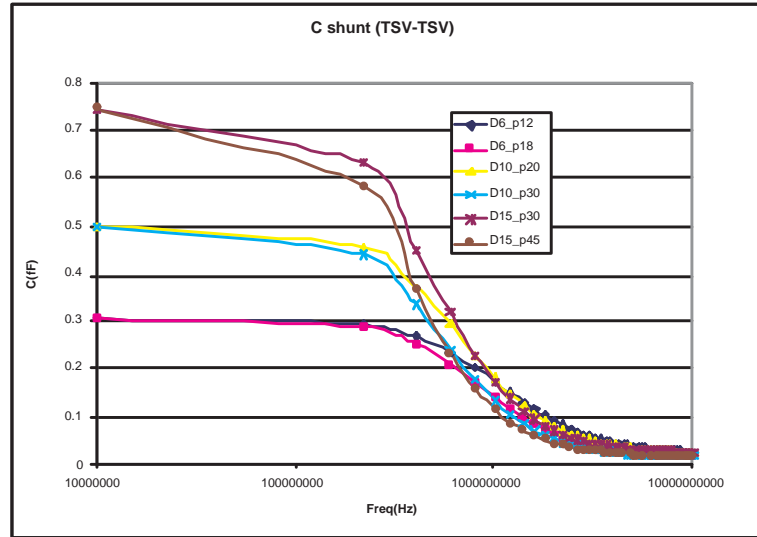


Figure 3.17: Substrate equivalent shunt coupling (D=diameter, P=pitch).

because of the skin effect. Figure 3.14 to Figure 3.17 show that in high frequency (over $1GHz$) substrate coupling decreases significantly.

3.5 Reliability of Stacked Power Grid

The *reliability* of a system is defined as the probability function $R(t)$, over the interval $[0, \infty]$ that the system operates without any failure. The reliability is defined as function of failure rate, $\lambda_f(t)$ or alternatively with *Mean Time To Failure (MTTF)* where $MTTF = 1/\lambda_f$. We use constant failure rate reliability model in our analysis. Reliability component is represented here by using exponential distribution [13] with a failure rate, λ_f , according to Equation(3.8):

$$R(t) = e^{-\lambda_f(t)} \quad (3.8)$$

3.5.1 Electromigration Constraint

In each power grid node i we have: $\Delta V_i(t) \leq \Delta V_{max}$ where ΔV_{max} is about 5% of nominal V_{dd} according to signoff corners in below $45nm$ technologies. Elec-

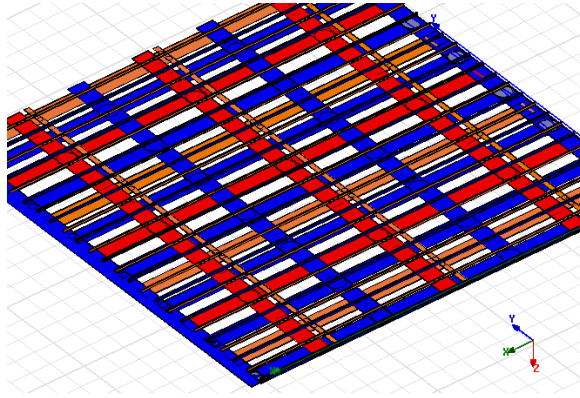


Figure 3.18: Power grid model in each tier.

tromigration is a major cause of momentum transfer from electrons to ions that makes interconnect lattice. In 3D IC, electromigration due to higher current variation and density, will lead to some major defects: shortening between adjacent metal layers, TSV and bumps; opening of metal lines and TSV contacts and increased resistance of metal lines and TSV contacts. Mean time to failure due the EM is modeled by the Black's model according to [41]:

$$MTTF_{TSV\ EM} = k_{bonding\ TSV} A_o (J - J_{crit})^{-n} \frac{E_a}{kT} \quad (3.9)$$

where A_o is an empirically determined constant. J is the current density in the interconnect. J_{crit} represents the threshold current density and K is the Boltzmann's constant, E_a is the activation energy and n is the scaling factor (usually n set to 2). EM in stacked wire segment sets an upper bound on the average current density. Thus, EM imposes a minimum wire width constraint. For a fixed thickness t_k of a layer k and given maximal current density J_{crit} , this constraint for a wire segment between node i and j can be expressed as:

$$\left| \frac{1}{T} \int_0^T (V_i(t) - V_j(t)) dx \right| \leq J_{crit} \cdot \rho l t \quad (3.10)$$

where ρ is the sheet resistance, l is the length of the wire segment, $V_i(t)$ is the voltage at node i , and $V_j(t)$ is voltage at node j and T represents period. Therefore the EM rule imposes either a minimum width of the wire with given power pitch or a minimum pitch with given wire width.

Black model for EM constraint is extended to take into account TSV exclusive EM factor due to the bonding. Stacked TSV has a bonded interface that may be Cu-Cu, or Cu-tin alloy bond. In addition to material differences, bond quality would also affect EM. $k_{bondingTSV}$ bonding coefficient in Equation (3.9) is the TSV empirical bonding coefficient and is determined by the fabrication process.

3.5.2 Thermo-mechanical Reliability of 3D Stacking

Thermo-mechanical modeling and analysis is a critical step in the design of a reliable 3D system. Global thermal mismatch between the copper filled TSV substrate and silicon chip is significant. Micro bumps would break under thermal stress conditions. Global thermal mismatch depends on both TSV size and pitch. Copper *coefficient of thermal expansion* (CTE), $\sim 17.5\mu/^\circ C$, is much higher than CTE of silicon, $\sim 2.5\mu/^\circ C$ [45]. Stacked die has delimitation potential due to thermal cycling and local thermal expansion mismatch between the copper and silicon substrate. Failure could potentially happen in the interface of TSV, dielectric, substrate, and grid [5]. λ_{TM} , effective stress and thermo-mechanical failure rate of TSV increases with TSV diameter or density reduction [45]. Selvanayagam *et al.* [45], performed a comprehensive study on thermo-mechanical stress and strain for multiple TSV dimensions. Thermal stress will degrade if density of the TSV increases because of reduction in thermal resistance. Figure 3.19 depicts normalized thermo-mechanical failure for various TSV densities and dimensions. TSV density(%) in Figure 3.19 shows TSV pitch percentage from maximum allocable TSV nodes.

3.5.3 Reliability-aware Experimental Results and Analysis

Proposed problem formulation is described in this section to derive the optimal reliable design parameters, i.e., TSV diameter and pitch of 3D PDN. The main objective of the optimization is to minimize the power noise while maximizing reliability. In the design of 3D PDN, two main constraints needs to be satisfied:

1. *Electromigration Current Density*: Current density should be bounded by

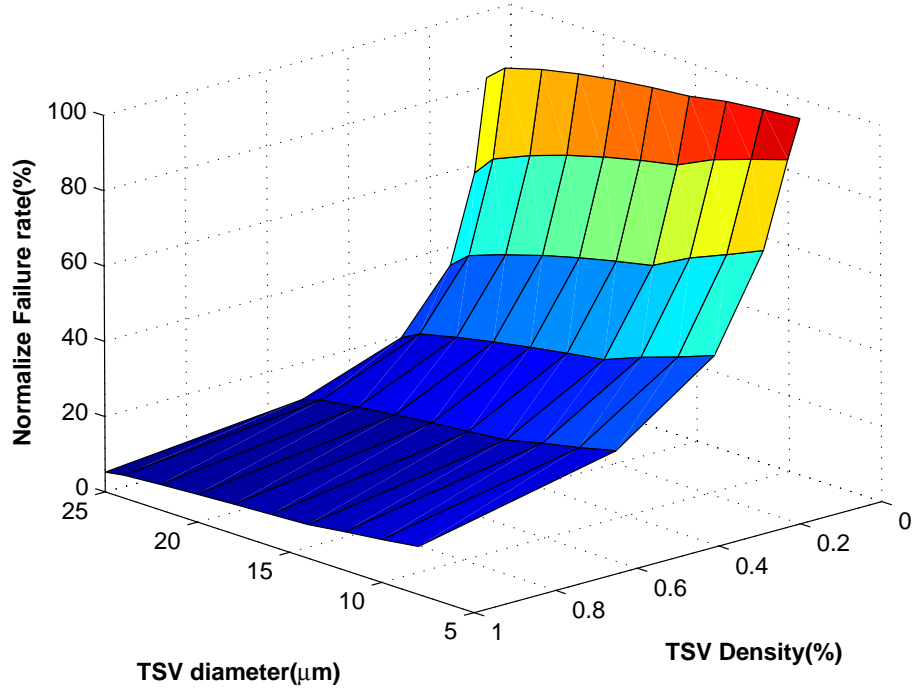


Figure 3.19: Normalized TSV thermo-mechanical failure rate(%).

the electromigration maximum current density as described in Section 3.5.1.

2. *Maximum Routable Area:* According to *Design Rules Checks*, surrounding TSV there is a block out region where no hard macro could be placed. The blocked area should be limited by the maximum available area dedicated to hard macro placement.

In addition, we define decoupling capacitor allocable area as $A_{decap} = k_{decap} \times A_{routable}$, where k_{decap} is the percentage of the routable area ($A_{routable}$). Therefore; there is a tradeoff between having dense TSVs with large size versus allocating more decoupling capacitors among tiers in the unblocked area as shown in Figure 3.20.

Reliability-aware Problem Formulation

The optimization problem is formulated as follows:

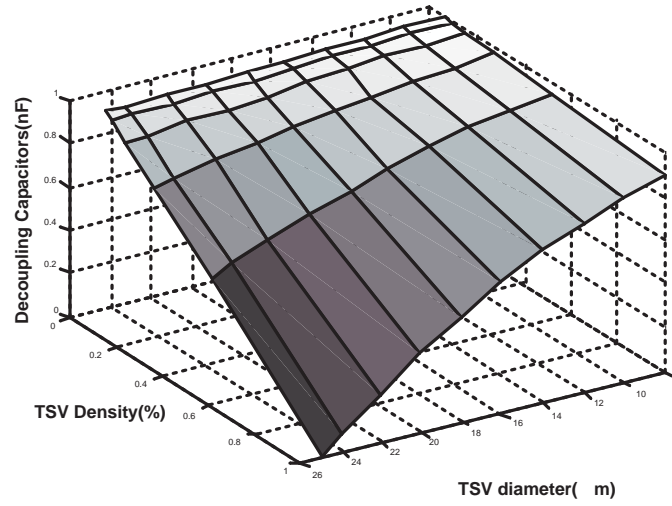


Figure 3.20: Decoupling capacitor allocation tradeoff versus TSV.

$$\text{Min} : \text{Failure} \times (\Sigma_i \Delta V_{max_{TSV}} + \Sigma_j \Delta V_{max_{tier}}) \quad (3.11)$$

Such that

$$(I) A_{blockout} \leq \min(A_{routable})$$

$$(II) \left| \frac{1}{T} \int_0^T (V_i(t) - V_j(t)) dx \right| \leq J_{crit} \cdot \rho l_{p,q} t \quad p, q \in 3DPDN$$

$$\text{Failure}(TSV_{diameter}, TSV_{density}) = 1 - R(t)$$

The objective is to minimize total power noise in TSV which are $\Delta V_{i_{TSV}}$ and in stacked layers $\Delta V_{j_{tier}}$. The keep out area in constraint (I) should be less than minimum routable area. Constraint (II) satisfies the EM maximum current density as discussed in Section 3.5.1. $R(t)$ is the reliability function as described in Section 3.5.

In the experiments, we model $1 \times 1 \text{ mm}^2$ of 5 layers of stacked dies with the tier-to-tier height of $100 \mu\text{m}$. Multiple diameters of the through silicon via are extracted and modeled. The current sources in each layer represent the active blocks

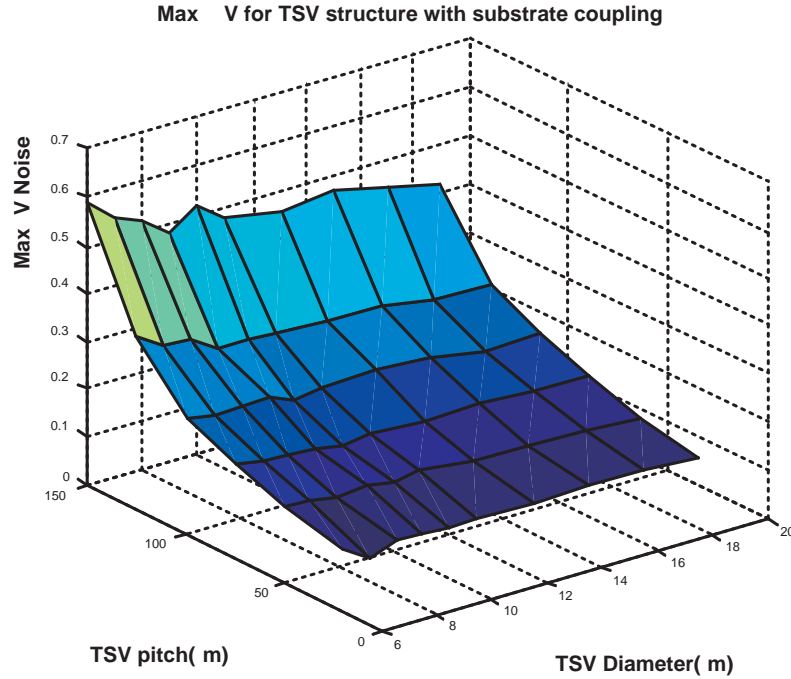


Figure 3.21: Power noise of 3D stacking considering substrate coupling model.

switching profile. The amplitude of the triangular current sources is calculated by maximum power divided by the total number of supply and nodes. Based on our flow described in Section 3.3.1, the maximum power noise is obtained as shown in Figure 3.21. We approximate the maximum voltage noise and reliability with a polynomial expression function of TSV diameter and pitch to solve the TSV linear optimization.

Finally, Figure 3.22 illustrates proposed cost function from optimization where the optimum configuration is derived. Figure 3.22 illustrates that based on the TSV density and diameter the optimal point is where the TSV density is the highest and the TSV diameter is in the middle range to minimize noise and maximize reliability. The noise value and reliability value are fitted into a polynomial expressions to derive the optimal cost function.

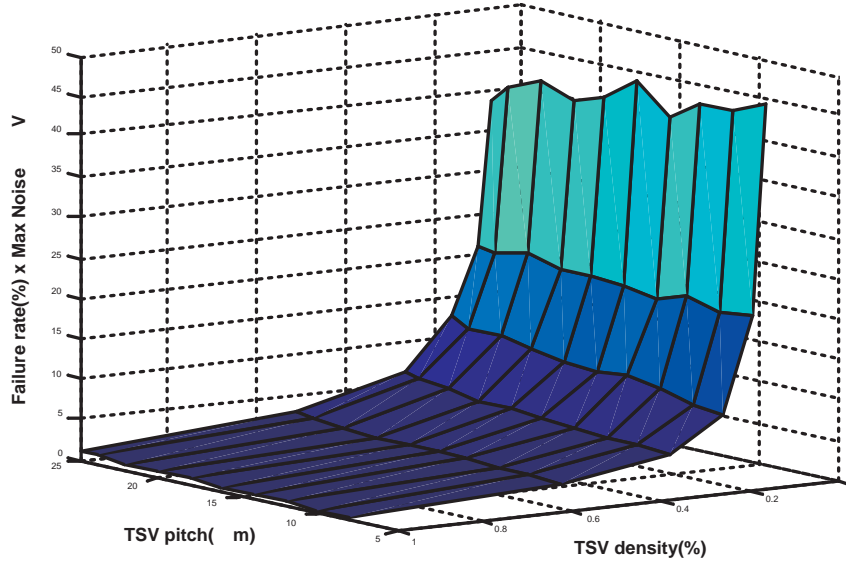


Figure 3.22: Optimization cost function($Failure\ rate_{TSV} \times max\ \Delta V$).

3.6 Summary

In this chapter, we developed an efficient parallel flow for the analysis of the 3D PDN. We explored different design parameters including the impact of the on-chip inductance, on-chip regulator and inductive through silicon via in the 3D PDN. In the second part of this chapter, a comprehensive 3D PDN stacked model is introduced. The model specifically focus on the substrate coupling of the 3D power grid and TSV with frequency dependent parasitic. The S -parameter model is converted into $RLGC(f)$ electrical model so that we can incorporate it into the frequency domain analysis flow. Electromigration current density constraint and reliability issues were discussed in this chapter. An optimization problem formulation is introduced with the objective of minimizing power noise under reliability constraints. The experimental results demonstrate efficacy of the proposed framework.

3.7 Acknowledgments

Chapter 3 is based on the following publications:

- A. Shayan, X. Hu, H. Peng, W. Zhang, M. Popovich, L. Chua-Eoan, C.K. Cheng, “Power Distribution Co-design for Nanoscale Stacked Silicon ICs”, *IEEE Conference on Electrical Performance of Electronic Packaging (EPEP)*, 2008.
- A. Shayan, X. Hu, M. Popovich, A.E. Engin, C.K. Cheng, “3D Stacked Power Distribution Considering Substrate Coupling”, *International Conference on Computer Design (ICCD)*, 2009.
- A. Shayan, X. Hu, H. Peng, W. Yu, T. Toms, M. Popovich, X. Chen, C.K. Cheng, “Reliability-aware Through Silicon Via Planning for Nanoscale Stacked Silicon ICs”, *Design Automation and Test in Europe (DATE)*, 2009.

The dissertation author was the primary researcher and author, and the co-authors involved in the above publications directed, supervised, and assisted in the research which forms the basis for that material.

Chapter 4

What Would be a Worst Current Scenario?

Worst case current stimulus that leads to maximum voltage noise is the focus of this chapter. We outline a frequency and time domain co-design flow that uses frequency domain results to construct time domain input vectors. The vectors are then adopted in a resonance-aware time domain analysis flow and can highlight low and mid frequency issues dominated by board and package. In the second section of this chapter, we propose a rogue wave-based vector generation algorithm that obtains realist worst-case current load.

4.1 Resonance-aware Methodology for System Level Power Distribution Network Co-design

4.1.1 Introduction

In the past decade, a number of EDA tools came into the market to assist designers optimizing PDN while minimizing its footprint in the overall system. These tools typically are divided into frequency and time domains. The frequency domain tools typically employ specialized fast electromagnetic solvers that take advantage of layered dielectric structures in package and board. Due to much

greater complexity and finer feature set of on-die power grids, silicon is usually modeled as a lattice of lumped RLC elements. For the frequency domain tools, the emphasis is therefore on analyzing the low and mid frequency PDN system responses [48].

One shortcoming of frequency domain tools is that the results are not directly expressed in millivolts of voltage drop seen by transistors. Supply voltage drop and ground bounce behavior could not be distinguished from the frequency domain electromagnetic solver results. Thus, it is difficult to reflect modifications needed to enhance the power delivery. They also do not analyze on-die power/ground grid structure in minute detail to assist silicon designers detect missing vias and shorts.

In the last few years, time domain PDN analysis tools grew out of the static IR drop tools that model full chip power/ground grid structures [54]. In the new time domain tools, die level interconnect is modeled in fine detail, along with on-die decoupling capacitors, power gating transistors, and switching elements. Through proprietary algorithms, time domain tools are now capable of running transient simulations to hundreds of nanoseconds. The results of time domain simulations can be very useful because they contain both current demand and instantaneous voltage information. Current demand, for example, is a critical piece of information when trying to set frequency domain impedance spec. Localized voltage drop hotspots can guide silicon designers to improve power/ground mesh in weak regions.

One limitation of time domain analysis is the necessity to stimulate a highly complex silicon design that consists of millions of transistors. Often an *activity factor* (AF) based input vector is used to toggle a percentage set of gates in the design at operational frequencies of individual functional blocks. In our work, both frequency domain and time domain PDN analysis flows are adopted to design the highly integrated mobile SOC chipsets. In this chapter, we will document a flow improvement to guide time domain analysis with PDN system resonances obtained through frequency domain analysis. In this fashion, the time domain analysis can better highlight PDN system behavior when it is perturbed at or near its resonance

frequencies.

The rest of this chapter is organized into five sections. Section 4.1.2 consists of a high level review of theoretical background. Section 4.1.3 focuses on frequency co-design of PDN system. Section 4.1.4 outlines a resonance-aware time domain methodology. Section 4.1.5 contains simulation results, and Section 4.3 summarizes our findings and concludes the chapter.

4.1.2 Theoretical Background: Resonance-aware Modulation

The relative importance in creating the right stimulus for PDN analysis can be seen in time domain results of Section 4.1.3. In this section, a simple lumped PDN model having a single resonant frequency at $75MHz$ is stimulated. The switching is from a high speed clock ($\sim 1GHz$) whose frequency is many times that of the PDN resonance. For the initial portion of the simulation the circuit draws a fixed charge per instruction cycle which sets up a steady DC current. The average of current has no energy near resonance. As a result there is little if any affect on noise since resonance is not disturbed. This is the period prior to $300ns$.

At a later time in the test the charge demand per cycle (top pane after $300ns$) is made to vary at a rate coincident with PDN resonance frequency. Now an AC current component in the package (middle pane) is developed. It perturbs the target frequency and a corresponding low frequency voltage component (lower pane) manifests itself in the power domain. It is of interest to note how the magnitude of this lower frequency component can rival or possibly exceed the localized high speed droop occurring at the clock rate. This result shows how time domain stimulus needs to target certain frequencies identified through frequency domain analysis even though the clock rate is well above PDN resonance.

Load Current Modulation

The notion of load modulation is very applicable for pipelined processors. The executed instructions can vary in both power magnitude as well as completion

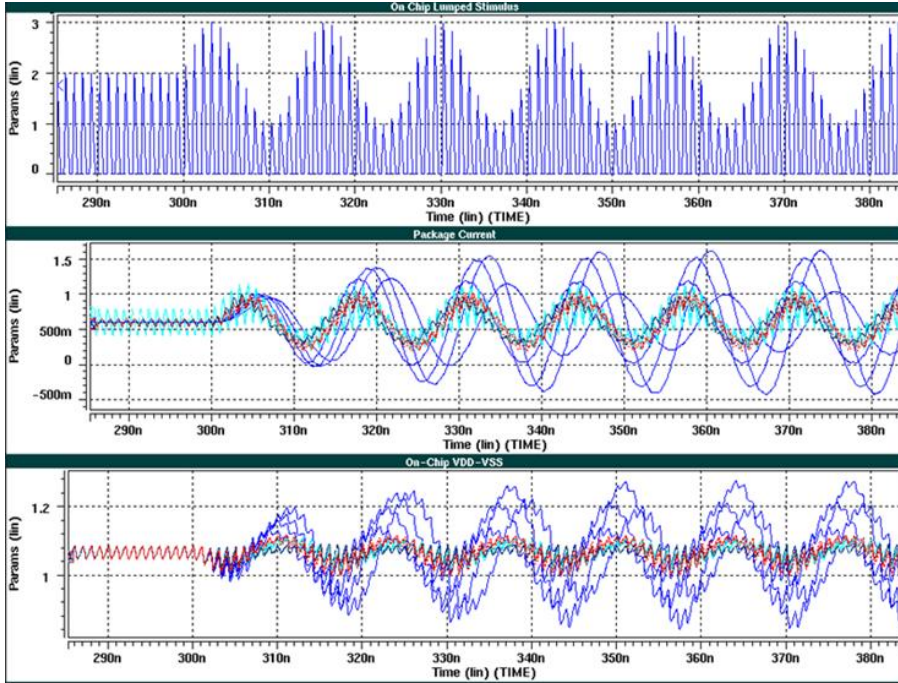


Figure 4.1: Resonance-aware load current modulation impact on voltage variation (Top: modulated current, Middle: package current, Bottom: on-chip voltage variation).

time relative to issue. The mathematics behind the load modulation properties for both time and frequency domain analysis lend themselves from basic AM radio modulation theory [53]. The equation for a generalized sinusoidal carrier can be expressed as follows:

$$y(t) = a(t) \cdot \cos(\omega_c t) \quad (4.1)$$

where we assume $a(t)$ varies slowly compared to the *carrier* frequency. The term $a(t)$ shows the envelope of the carrier frequency. For PDN analysis we would replace the single cosine term representing a *carrier* with a more appropriate Fourier sum. This way we represent the cycle-to-cycle current demand of the on-die circuits whose fundamental clock rate is known. Thus, the modulated current demand would take on the following form:

$$y(t) = a(t) \cdot \left[\frac{\alpha_0}{2} + \sum_{n=1}^N (\alpha_n \cos(n\omega_c t) + \beta_n \sin(n\omega_c t)) \right] \quad (4.2)$$

Even though the modulating carrier $a(t)$ is expressed as a general aperiodic signal, but the PDN issue occurs when enough periodicity exists such that resonance is only momentarily perturbed. This could be for only a few cycles over an extended length of time. In the frequency domain, we expect the ‘carrier’ or high speed current demand which occurs at high frequency to appear as a fundamental clock spur with its associated harmonics implied through Fourier analysis. In addition, we expect to see the modulated energy positioned relative to DC and also positioned around the carrier much like side-band energy with an AM carrier as shown in Figures 4.8(a) and 4.8(b) [53]. For a typical PDN system and resulting resonant frequencies, we only need to concern ourselves with the energy in the 50MHz region for standard cores and somewhat higher frequencies for I/O interfaces.

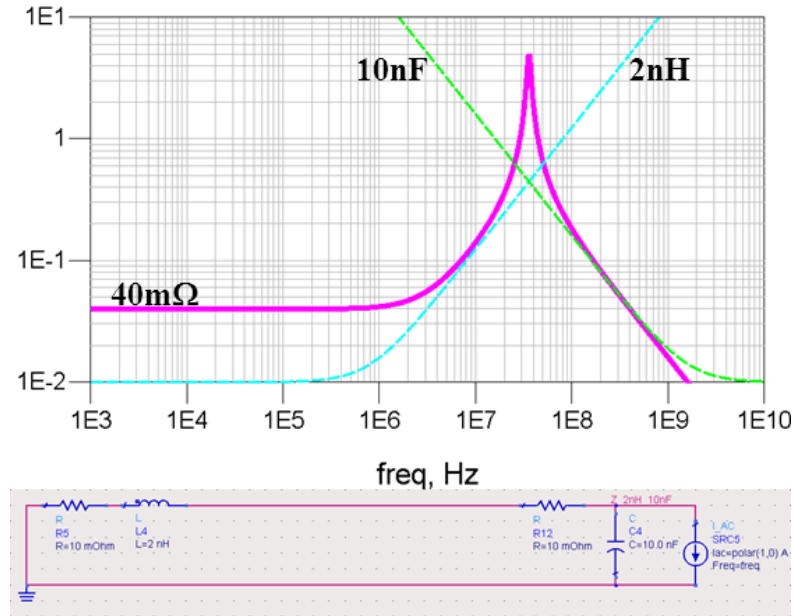


Figure 4.2: RLC resonance with $R=40m$, $L=2nH$, and $C=10nF$.

4.1.3 Broadband Frequency Domain Co-design

PDN design for highly integrated mobile chipsets spans from board, package, and to silicon. PDN modeling is complicated by the following unique properties:

1. Large problem size. Board PDN network is electrically large with power shapes ranging in 10's of millimeters.
2. Wide range of feature sizes. Feature sizes found in a PDN network varies from 10's of millimeter on board, to 10's of microns on package, to micros on silicon.
3. Fully coupled system. Power delivery shapes and lines are fully coupled in a single system. Accurate modeling of PDN at board/package and package/die resonances is critical to ensure sound design.
4. Highly non-uniform die-level capacitive and current loading.

PDN design essentially behaves as a RLC passive circuit. The resistance R is the accumulated board/package routing parasitic and on-die P/G mesh parasitics. The inductance is the accumulation of board and package routing. The capacitance comes in from silicon intrinsic capacitance and any silicon, package, and board decoupling capacitors.

Figure 4.2 shows a lumped circuit with $40m\Omega$ of resistance, $2nH$ of inductance, and $10nF$ of capacitance. This circuit has a resonance with magnitude of roughly $5m\Omega$ at about $30MHz$. If the current source on the right hand side perturbs this resonance, a large voltage drop is expected to develop across the inductor. The design of PDN tries to minimize this resonance and, if possible, push it outside of current spectral bandwidth expected from the current source. One of the methods for controlling the resonance peak and to shift it to higher frequencies is to break the inductance into segments by inserting shunt capacitances. In Figure 4.3, a $10\mu F$ capacitor and a $100nF$ capacitor are inserted to break the $2nH$ into $1nH$, $0.95nH$, and $50pH$ segments. At about $2.5MHz$, the effectiveness of the $10\mu F$ is maximized. The $0.95nH$ inductor after this capacitor starts to dominate

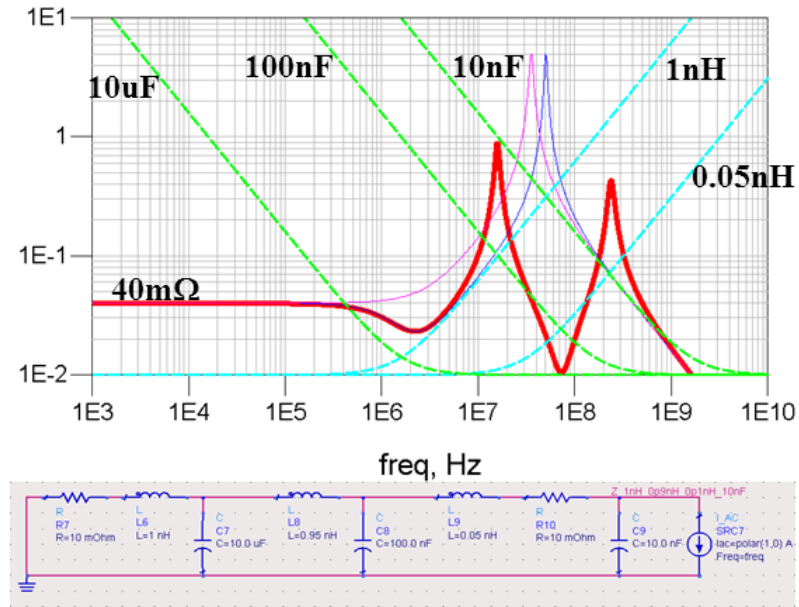


Figure 4.3: RLC resonance after the $2nH$ inductance is segmented with shunt capacitances.

the system impedance. The system reaches a resonance at about $20MHz$, which is a product of the $0.95nH$ inductor and $100nF$ capacitor. After this frequency, the $100nF$ capacitor dominates the impedance profile until around $100MHz$ when the $50pH$ inductor starts to dominate. The system reaches another resonance peak at around $200MHz$ and after this point the $10nF$ capacitor becomes the dominating factor. This simple RLC example serves as an example of how PDN system impedance can be controlled via insertion of decoupling capacitors. The large capacitors are inserted as close to the silicon package as possible.

4.1.4 Resonance-aware Time Domain Methodology

We have developed a method based on current demand shaping which can be used in time domain simulators to create a stimulus. We generate a current demand during a VCD analysis run, with important energy content close to the resonant frequency of the PDN as presented in Figure 4.3. One of the major challenges of PDN signoff is the ability to generate realistic design stimulus. The

proper stimulus should excite the power delivery network close to the resonant frequency. In other words, the ability to create a current stimulus which is design specific and has enough energy contents close to the resonant frequency of the system PDN. We have developed a new method based on current demand shaping which can be used in time domain simulators to create a designer shaped stimulus. This stimulus will generate a current demand during a VCD analysis which has important energy content close to the resonant frequency of the PDN as presented in Figure 4.3. Current modulation has following key parameters:

1. Modulation functions (sine, cosine, pulse, etc.).
2. Modulation frequency (functional block defines).
3. Modulation depth (user defined).

The resonant frequency of the PDN is fixed and is determined by the physical implementation of the die, package and PCB. Usually the resonant frequency is close to $20\text{-}50\text{MHz}$ as presented in 4.3. We have extracted the impedance profile for the complete PDN V_{dd} domains die, complete layout package and layout PCB. We can create a VCD stimulus or different vectorless runs by knowing the resonant frequency of the PDN which will stimulate the PDN close to the resonant frequency.

$$I_{comb}(t) = \begin{cases} I_1 & \text{if } (t \leq T_1) \\ 0 & \text{if } (T_1 \leq t \leq T_1 + T_3) \\ I_2(t) & (T_1 + T_3 < t \leq T_1 + T_3 + T_4) \\ I_3(t) & T_1 + T_3 + T_4 + n \times (T_5 + T_2) < t \leq T_1 + T_2 + T_3 + T_4 + n \times (T_5 + T_2) \\ 0 & T_1 + T_2 + T_3 + T_4 + n \times (T_5 + T_2) < t \leq T_1 + T_3 + T_4 + (n + 1) \times (T_5 + T_2) \\ & \text{for } 0 \leq n \end{cases}$$

The process of generating the desired current demand is controlled by a configuration which defines the duration of individual current wavelets. From multiple vector runs, we can create in time domain a new VCD run with the desired designed current profile as presented in Figure 4.4. The combined current shape can be validated using an FFT analysis to ensure that we have maximized

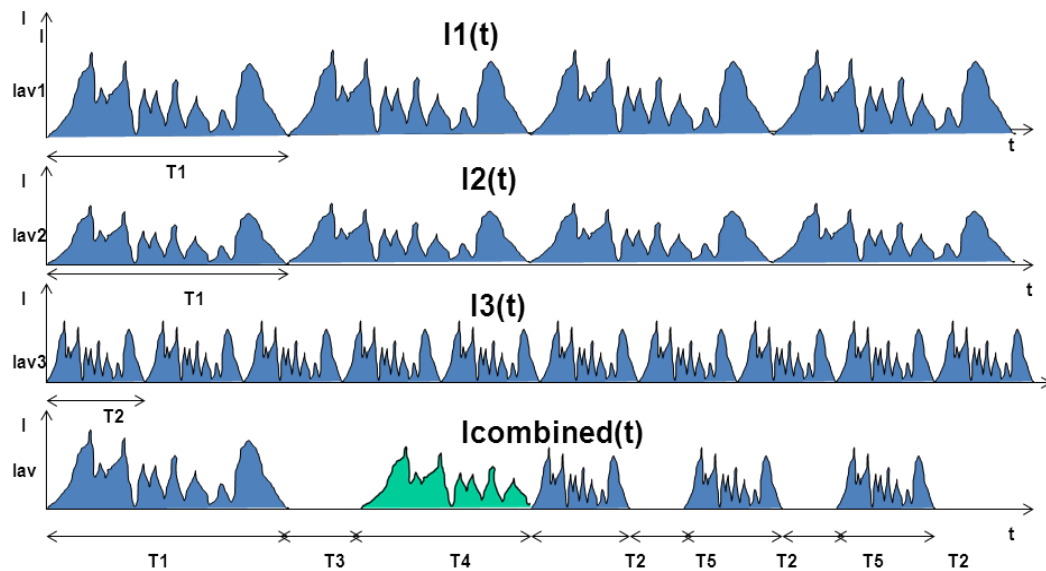


Figure 4.4: Current waveform generated from individual wavelets.

the spectral current components in frequency domain close to PDN resonance. The difference in spectral content for the current demand between a regular vector runs and the *designed* VCD run is presented in Figure ???. The basic function for the final designed current demand is presented below in Figure 4.4:

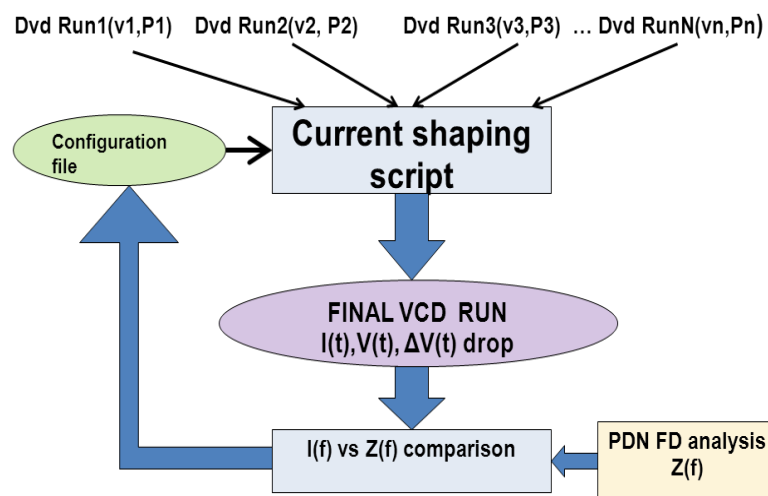


Figure 4.5: Flow chart of the time domain worst case Z-aware PDN analysis.

Architectural Perspective of Current Load

Table 4.1 shows how a processor core can stimulate any given resonant frequency. This goal is achieved by means of executing finite code loops at various clock frequencies. Table 4.1 is a setup guideline for targeted stimulus such that assembly code for either test purposes or VCD generation for dynamic IR drop analysis can be facilitated.

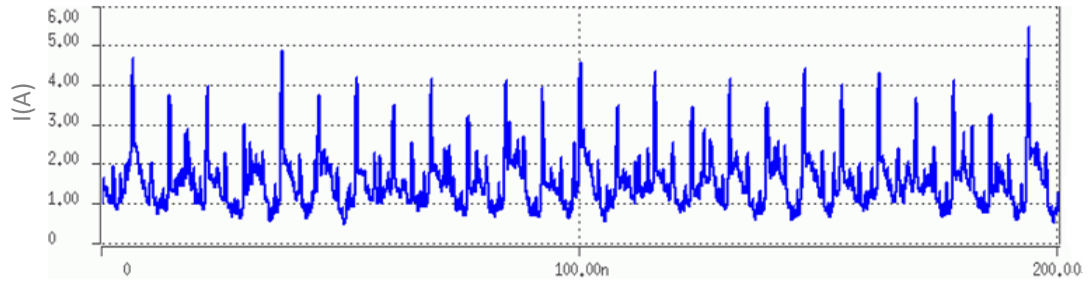
Table 4.1: Setup table for current modulation of the processor.

Cycle	Clock Frequency (MHz)								
	550	575	600	625	650	675	700	725	750
6	91.67	95.83	100.00	104.17	108.33	112.50	116.67	120.83	125.00
8	68.75	71.88	75.00	78.13	81.25	84.38	87.50	90.63	93.75
10	55.00	57.50	60.00	62.50	65.00	67.50	70.00	72.50	75.00
12	45.83	47.92	50.00	52.08	54.17	56.25	58.33	60.42	62.50
14	39.29	41.07	42.86	44.64	46.43	48.21	50.00	51.79	53.57
16	34.38	35.94	37.50	39.06	40.63	42.19	43.75	45.31	46.88
18	30.56	31.94	33.33	34.72	36.11	37.50	38.89	40.28	41.67
20	27.50	28.75	30.00	31.25	32.50	33.75	35.00	36.25	37.50

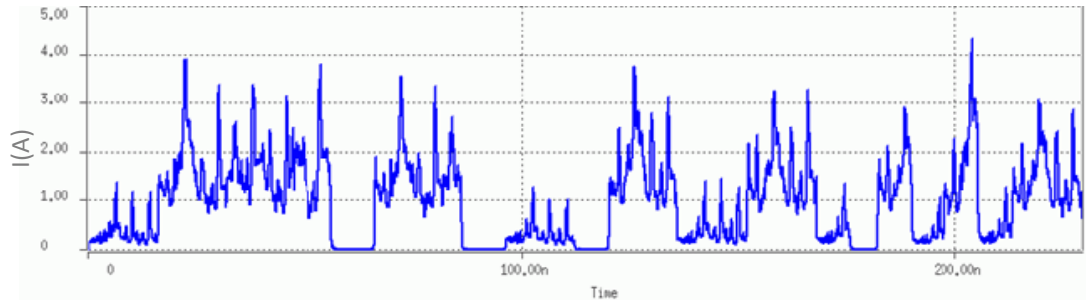
In this context, the cycle count is the implied code loop length which is repeated. Horizontal scale is the clock frequency in *MHz*. The purpose of any code loop is to generate a variation in current demand such that the perceived min and max possible current magnitudes are explored. Generating such stimulus is the essence of exploring what dynamic range of current variation for any given system component is, during a specific resonance cycle. Only real code running on a design for simulation purposes will yield the true bounds so that effective modulation indexes can then be established for further use.

4.1.5 Experimental Results

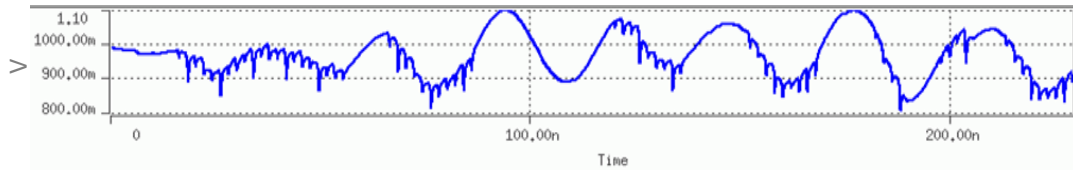
We generate the worst realistic di/dt using proposed method, which results to the worst voltage noise for the complete system. Based on the method presented above using time domain simulators, we are creating a system level simulation test bench with package and PCB. The goal is to produce the worst realistic voltage



(a) Original on-chip current demand from a vector run.



(b) Modulated current from the proposed method.

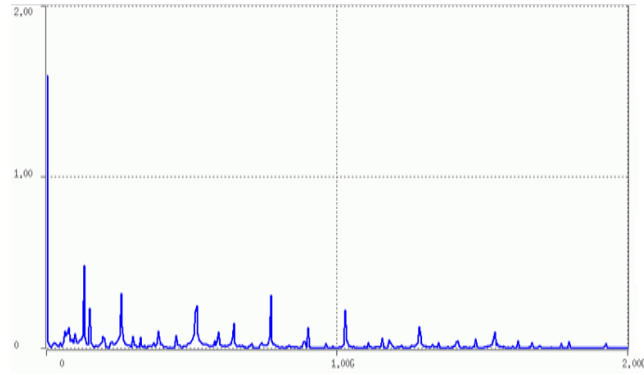


(c) Voltage waveform for an instance with modulation where the resonance and low frequency content is more dominant.

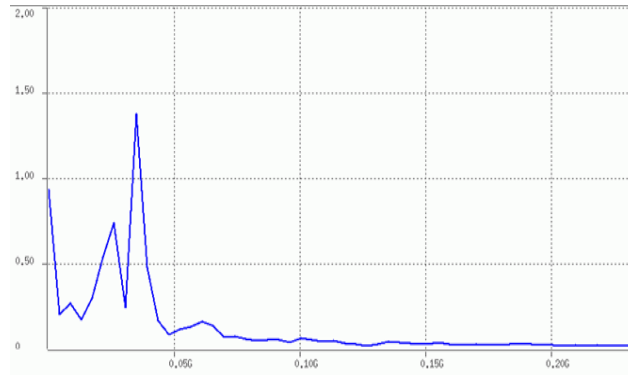
Figure 4.6: Synthetic resonance-aware current generation.

noise. This designer driven simulation scenario can be used for the signoff of the PDN, to validate the voltage corners used during time domain analysis and to validate the results of the frequency domain PDN optimization and analysis. As an example of multi resonances PDN design Figure 4.8 demonstrates time and frequency domain waveforms of Figure 4.2 with dual resonance of a digital system clocked at $1.0GHz$ but having two tone modulations at $20MHz$ and $200MHz$. The first pane is the time domain waveform while the later is the frequency domain representation of this modulation scenario. In this example, piece-wise linear

current waveform amplitude is $0.5A$. Rise time and fall time is about $0.5ns$. Dual resonance frequency ($f_1 = 20MHz$ and $f_2 = 200MHz$) modulation is performed based on: $a_1 \cdot \sin(2\pi f_1 t) + a_2 \cdot \sin(2\pi f_2 t)$, where $a_1 = a_2 = 0.25A$ and the modulation index is 0.5.



(a) FFT comparison between the two current demands.

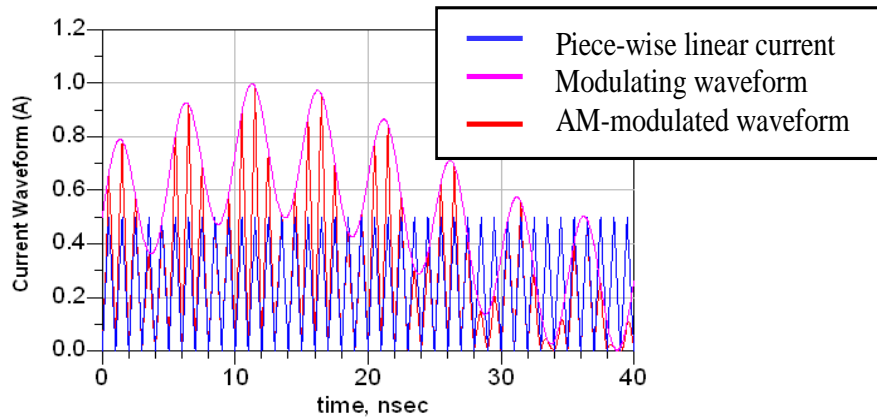


(b) The modulated current has a frequency content in the range of resonance (about $30-40MHz$).

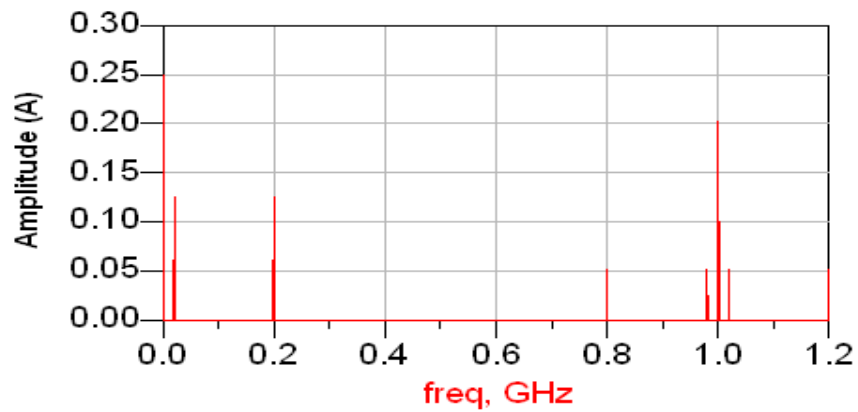
Figure 4.7: Frequency energy content of the modulated vector current.

In the provided waveforms one can observe a portion of the $20MHz$ envelope while multiple cycles of the $200MHz$ envelope are visible in the time domain pane. The frequency domain pane on the right side clearly shows that energy exists at $20MHz$ and $200MHz$. Although this example is contrived, it is the case that pipelined processors which iterate code loops of fixed length can focus energy into a

specific frequency while that energy level varies with time. In general, the problem is more aperiodic in nature, but we need to be aware of what can or will happen when resonance is perturbed and try to determine the low frequency voltage noise magnitude.



(a) FFT comparison between the two current demands.



(b) Dual resonance current modulation time domain waveform.

Figure 4.8: Multiple modes resonance-aware modulation.

4.2 Vector-based Rogue Wave Synthesis Algorithm Using Realistic Current Activity

The goal of the second section of this chapter is to outline our algorithm for combining realistic vectors and theoretical worst-case load. This final result would be a synthetic current that represent worst current called rogue wave. With a given current load pattern I , which is sampled from VCD vector or SPEC benchmarks. We can synthesize the realistic worst-case current to induce the worst-case voltage noise. We synthesize the worst-case current by adopting realistic vectors current pattern.

The current pattern keeps the temporal correlation between two consecutive current peaks. That means it is likelihood that one current peak will follows its nearby current peak. Therefore, when synthesizing the worst-case current stimulus, we use a sliding window to preserve the temporal correlation information.

Algorithm 1 shows our pseudo code for pre-processing of current and sorting. We use the sorted current to do computation for the worst-case synthesis. The computation involves a heuristic that the descending sorted list is chosen in the positive first elements, and vice versa. The heuristic cannot guarantee the elements follow the positive first sections are all positive as it depends on the sliding window size. Under this heuristic, most of the windows can find the maximum voltage noise as soon as possible. The enhanced pseudo code is listed below.

We partition the impulse response into equal sliding windows. We use convolution operation to find which interval of current pattern will cause the maximum voltage noise. Because window can preserve the temporal correlation of the current pattern, we can have more realistic voltage noise. After we find the position of the maximum voltage noise, we will use the current peaks in the corresponding interval of current pattern as our worst-case current stimulus of the window.

4.2.1 Complexity of Rogue Wave Synthesis Algorithm

The complexity analysis of the proposed algorithm is as follows: We assume there are K intervals for the time spanning with window size of m . The complexity

Algorithm 1 Algorithm for vector-based rogue wave current generation.

```

1:  $M \leftarrow$  is the size of current pattern
2:  $N \leftarrow$  is the size of impulse_response
3: for  $i = 1 \rightarrow N - window\_size$  do
4:   sum each current peak of current pattern( $i, i + window\_size - 1$ )
5: end for
6:  $sorted\_list\_des =$  sorting the sum of the intervals of current peak descending
7:  $sorted\_list\_asc =$  sorting the sum of the intervals of current peak ascending
8: for  $i = 1 \rightarrow N - window\_size$  step  $window\_size$  do
9:   sum each current peak of current pattern( $i, i + window\_size - 1$ )
10:  if  $impulse\_response(i) > 0$  then
11:     $current\_list = sorted\_list\_des$ 
12:  else
13:     $current\_list = sorted\_list\_asc$ 
14:  end if
15:  for  $j = 0 \rightarrow M - window\_size + 1$  do
16:     $idx\_current = current\_list(j)$ 
17:     $tmp\_val = convolution\ of\ impulse\_response(i, i + window\_size - 1)$ 
18:     $current\_pattern(idx\_crnt, idx\_crnt + window\_size - 1)$ 
19:    if  $tmp\_val > max\_val$  then
20:       $max\_val = tmp\_val$ 
21:       $max\_current = current\_pattern(idx\_crnt, idx\_crnt + win\_size - 1)$ 
22:    else
23:      break
24:    end if
25:  end for
26: end for

```

for pre-processing sorting of current is $O(K \log K)$. For the for-loop of synthesizing worst-case current signature, the complexity of convolution is $O(m \log m)$. The loops repeats convolution for $K \times (N - m)$ times. Therefore, the overall complexity of synthesizing worst-case is $O(K \times (N - m) \times m \log m)$. Then, the complexity of the

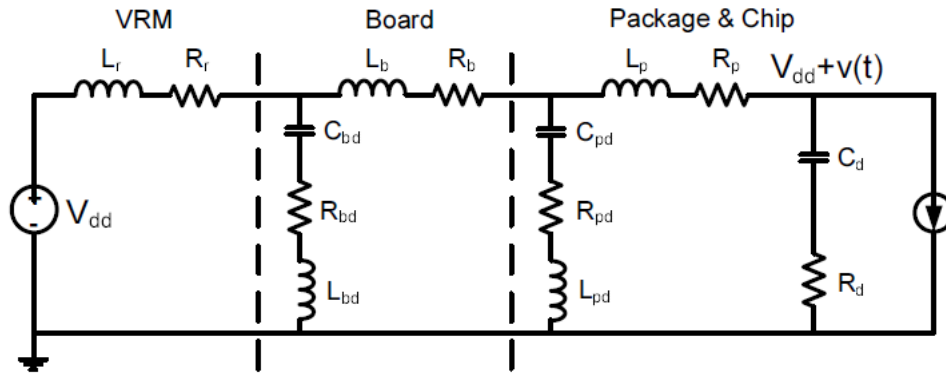
Table 4.2: Sensitivity of noise to Q -factor setting of synthesized rogue wave.

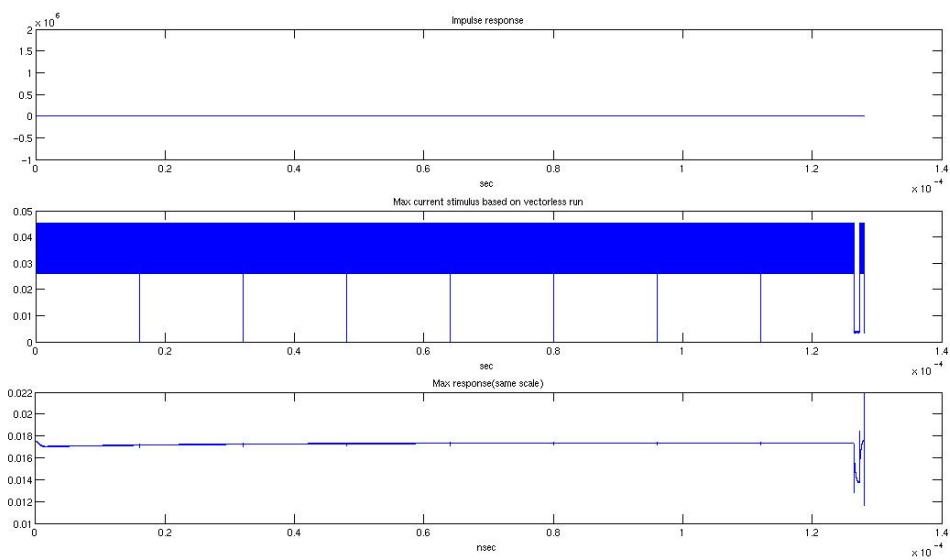
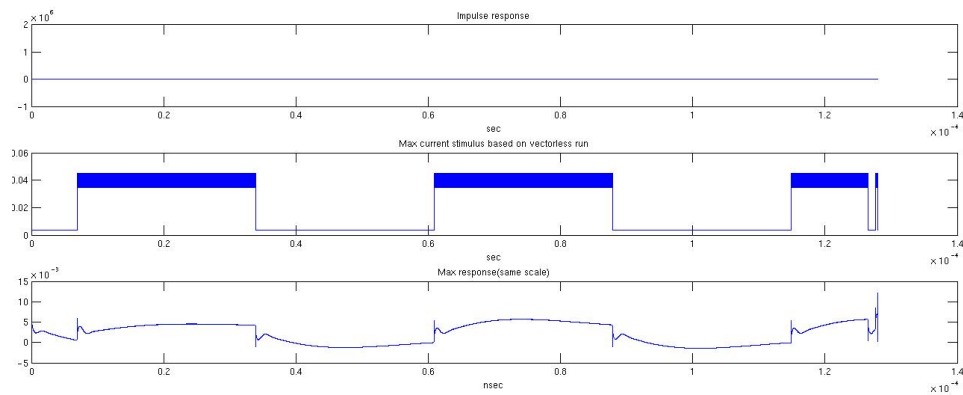
Lr	Cbd	Rr	Rbd	Q
10^{-9}	5×10^{-3}	5×10^{-3}	1×10^{-3}	0.2357
10^{-9}	5×10^{-3}	1×10^{-3}	0.5×10^{-3}	0.9428
50×10^{-9}	0.5×10^{-3}	1×10^{-3}	0.5×10^{-3}	6.6667
50×10^{-9}	0.5×10^{-3}	1×10^{-3}	0.5×10^{-3}	2.1082
50×10^{-9}	0.5×10^{-3}	0.5×10^{-3}	0.5×10^{-3}	3.1623
50×10^{-9}	4×10^{-3}	0.5×10^{-3}	0.5×10^{-3}	3.5355

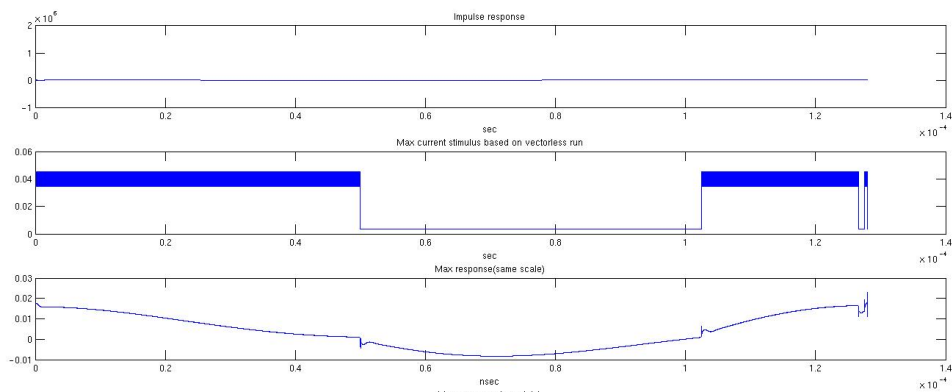
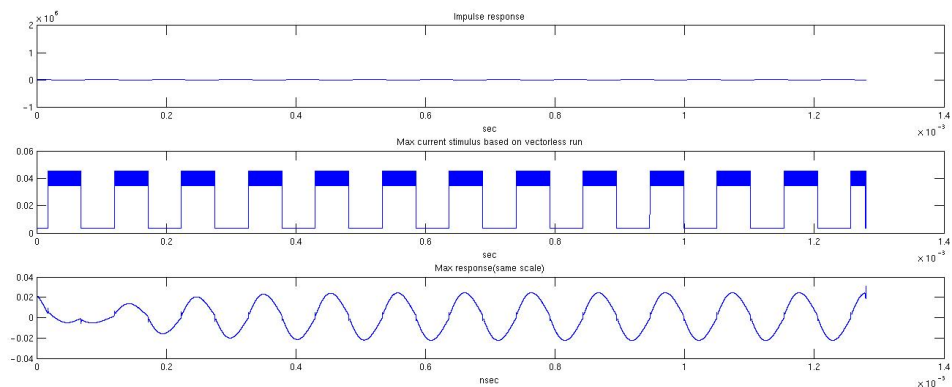
algorithm is $O(K \times (N - m) \times m \log m + K \log K)$. We know that $K = \text{ceiling}(N/m)$. The total complexity of the proposed algorithm can be expressed as $O(N^2 \cdot \log m)$.

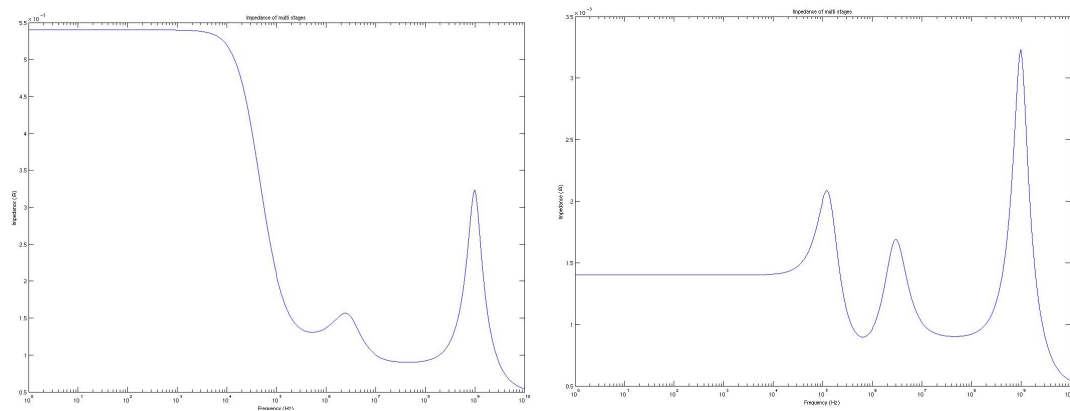
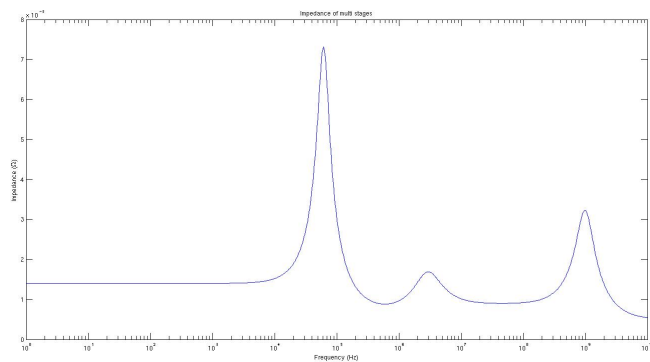
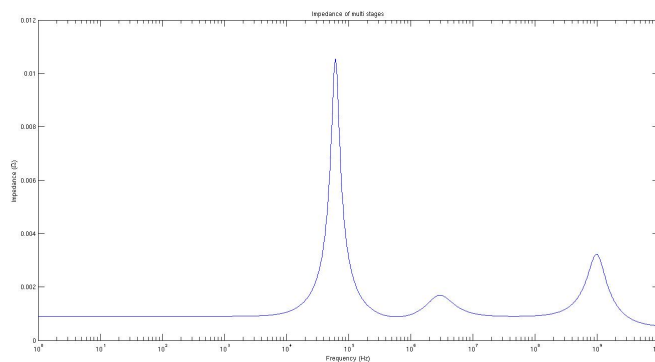
4.2.2 Sensitivity of Voltage Noise to Different Q -factors

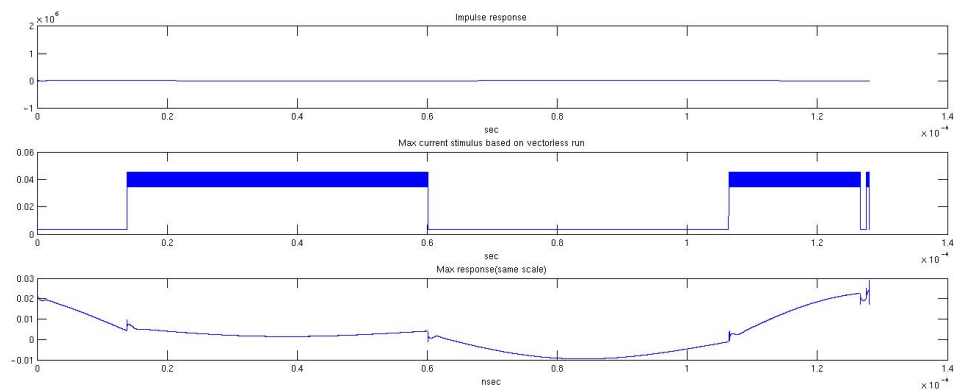
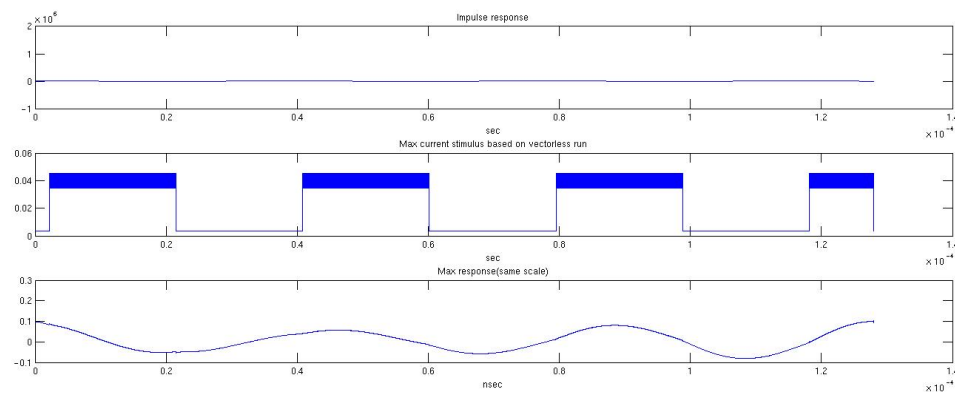
The impact of the rogue wave stimuli is demonstrated through a three stage network model. Table 4.2 and Figures 4.11 to 4.14 illustrate our experimental results for different Q -factors.

**Figure 4.9:** Lumped model of multiple stage power distribution.

(a) $Q = 0.2357$ (b) $Q = 0.9428$ **Figure 4.10:** Synthesized rogue wave stimuli for multiple Q -factors set 1.

(a) $Q = 2.1082$ (b) $Q = 3.1623$ **Figure 4.11:** Synthesized rogue wave stimuli for multiple Q -factors set 2.

(a) $Q = 0.2357$ (b) $Q = 0.9428$ (c) $Q = 2.1082$ (d) $Q = 3.1623$ **Figure 4.12:** PDN Impedance profile change for different Q -factors.

(a) $Q = 3.5355$ (b) $Q = 6.6667$ **Figure 4.13:** Synthesize of rogue wave stimuli for multiple Q -factors set 3.

4.3 Summary

We proposed an enhanced flow for the analysis of PDN under worst-case realistic load current. The proposed methodology is based on resonance-aware modulation of core power activity to perturb the PDN under resonance. This way we mimic core activity realistically under worst-case operation mode to guide design of a robust PDN.

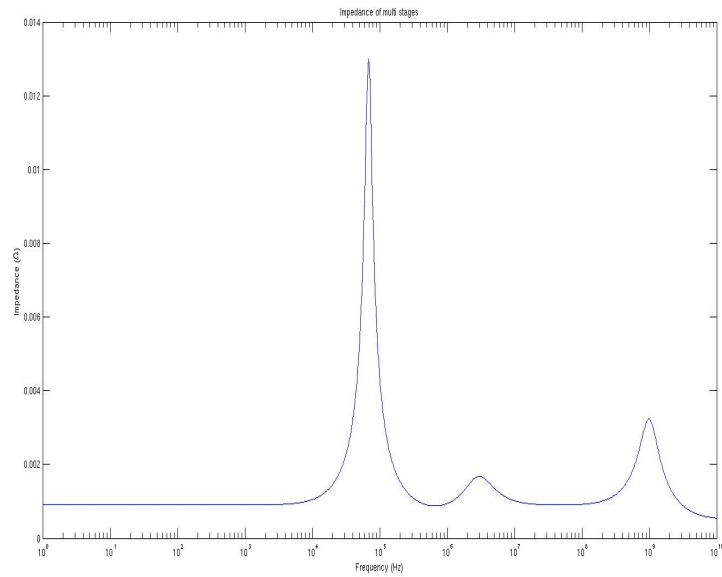
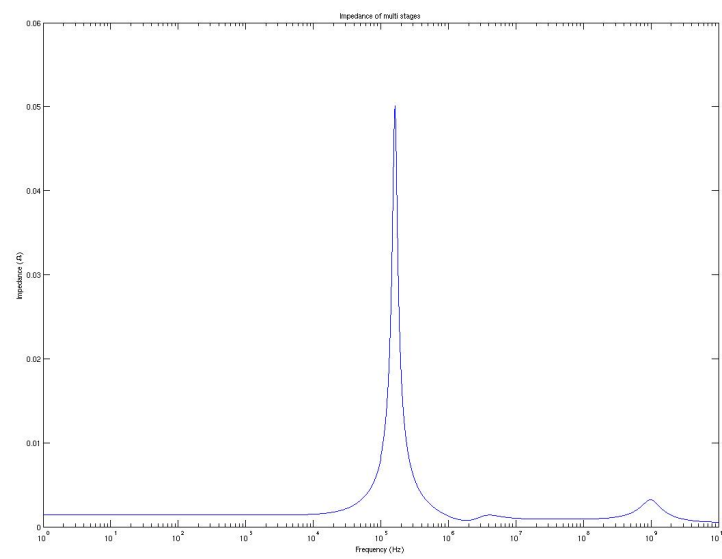
In second section of this chapter, we proposed an algorithm for generating a realist current for worst voltage drop. We outlined our vector based rogue wave current generation algorithm in this chapter. We preserved the temporal characteristics of the vectors by using a flexible sliding window of observation in the current generation process. The final current load will provide designers a better understanding of the worst voltage drop.

4.4 Acknowledgments

Chapter 4 is based on the following publication:

- A. Shayan, K. Bowles, S. Dobre, M. Popovich, X. Chen, C. Pan , “Resonance-aware Modulation Methodology for System Level Power Distribution Co-design”, *IEEE Conference on Electrical Performance of Electronic Packaging (EPEP)*, 2009.

The dissertation author was the primary researcher and author, and the co-authors involved in the above publications directed, supervised, and assisted in the research which forms the basis for that material.

(a) $Q = 3.5355$ (b) $Q = 6.6667$ **Figure 4.14:** Impedance profile in different Q -factor PDN system.

Chapter 5

Low Power Distribution Network Regulations

On-die regulation is a promising technique both for power saving and noise reduction. We introduce a regulation design methodology based on on-chip *Linear Dropout Regulators* (LDO). We then look into the regulator impact on power noise. An efficient parallel flow for the design of the complete power distribution system is proposed. The analysis focus on the impact of the voltage regulator model in both frequency and time domain response. Based on the experimental results, it is observed that including the voltage regulator model in the PDN model increases the transient voltage drop. The increase in noise cannot be ignored in current processing node with below 1V of supply. The flow runtime is optimized using parallel processing to speedup slow response simulation time of the off-chip voltage regulator. The study highlights the power integrity issues related to voltage regulator in broadband frequency ranges. The experimental results show speedup of up to 22× times with single processor and more than 430× times using up to 200 processors compared with HSPICE and other commercial simulators. Simulation time of PDN is reduced from hours to less than a minute.

5.1 On-die LDO-based Design and Optimization of PDN under Worst Loading

We introduce a methodology to design and to optimize power distribution with on-die Linear Dropout Regulators. An analytical formula for worst voltage drop is derived that takes into account LDO-PDN system poles and zeros. LDO power and decoupling capacitor optimization flow is proposed to meet the system voltage noise requirements.

5.1.1 Introduction

Power distribution network (PDN) design continues to be a major challenge because of the increased power demand and limited chip resource. Many functional blocks compete for the expensive chip area. Thus, conventional methods for power distribution such as decoupling capacitors and pure passive network, are not sufficient for robust power delivery. On-die Linear Dropout Regulator (LDO) circuit improves some of the shortcomings of pure passive power network [28, 32, 26].

Power routing in multiple voltage domains system is non-uniform and scattered. As a result, the passive network has high inductance and high resistance impedance. Adoption of on-die LDO relaxes the off-chip impedance. Thus, multiple voltage domains share same power rail in the input side of the LDOs with less parasitics to reach a relaxed off-chip impedance. On the output side, LDOs regulate each domain individually. Multiple voltage domains will also benefit from the optimal local LDO regulation. Fast response time of the on-die regulator enables real time adoption to the load changes and hence make finer grain power management feasible for each block.

In this paper, we focus on the design of the power distribution with on-die linear dropout regulators. In the first step, we derive the impedance profile. Main dominant poles and zeros of the LDO and PDN system are identified. We introduce an analytical formulation to obtain the worst case voltage drop in presence of the LDO. The analytical results will be adopted as means to design and to optimize power distribution. We analyze LDO power consumption and decoupling capacitor

area tradeoff for minimizing worst voltage drop. Experimental results show that we are able to efficiently reduce the worst voltage drop to the target value.

The rest of this chapter is organized as follows. In Section 5.1.2, we describe the worst case noise optimization problem formulation in the LDO based PDN system. Section 5.1.3 provides description of the LDO design and the main poles and zeros in frequency domain. We derive the impedance of the power distribution. Next, we calculate the analytical formulation for the worst voltage drop in Section 5.1.4. LDO-PDN system tradeoff and experimental analysis are discussed in Section 5.1.6.

For the second part, we discuss the VRM based analysis flow: In Section 5.2.2, we will explain the details of power distribution model, voltage regulator model and 2D partitioned model that we use in the flow. In Section 5.2.3, we will discuss the methodology and parallel flow. Section 5.2.6, demonstrates the complete PDN analysis and our experimental results. Finally, we conclude the summary of this chapter in Section 5.3.

5.1.2 Problem Formulation

On-die LDO regulates the supply rail of the functional blocks and decouples the power network from noisy off-chip path. The LDO block provides specific supply voltage and current, based on the functional block load demand. The LDO-based power distribution model in this work is comprised of both on-chip decap and off-chip decap to represent a generic scenario. The off-chip decap path includes the parasitics of the C4 bumps and routing path seen by the load as shown in Figure 5.1. The contribution of this section is a methodology that:

- Derive the analytical worst voltage drop in presence of the LDO and power distribution system poles and zeros.
- Split the resource between decap and LDO (area and power) to:
 - Minimize the worst voltage drop.
 - Maintain the stability of the LDO feedback system.

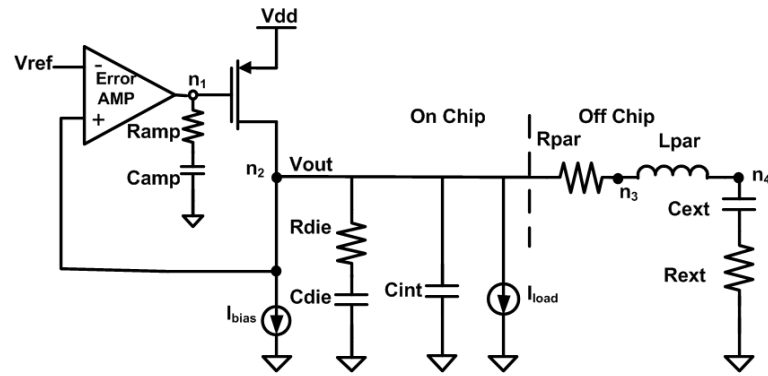


Figure 5.1: Model of LDO-based power distribution.

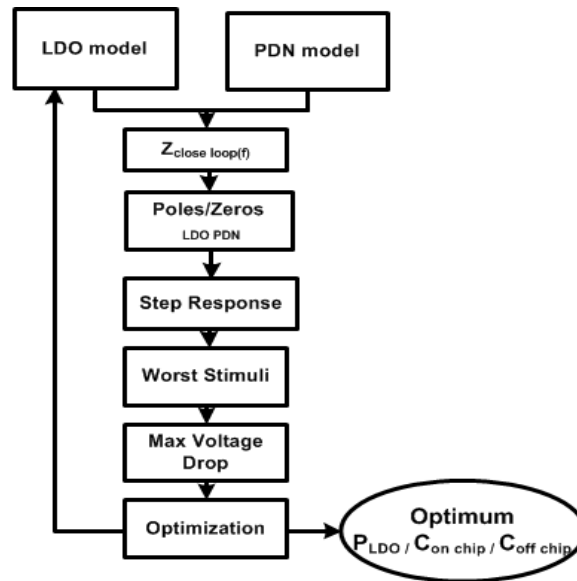


Figure 5.2: LDO-PDN system worst-case noise optimization flow.

Our goal is to minimize the worst voltage drop (V_{max}), such that we meet our total noise, LDO power and area budget. Problem formulation thus becomes Equation (5.1):

$$\begin{aligned}
& \underset{(P_{LDO}, C_{on\ chip}, C_{off\ chip})}{\text{Minimize}} && V_{max}(t) = I_{load\ step}(t) * Z_{LDO-PDN} \\
& \text{subject to} && P_{LDO} < P_0 \\
& && C_0 \leq C_{on\ chip} < C_1 \\
& && C_{off\ chip} < C_2 \\
& && \| I_{load\ step}(t) \| \leq I_{peak}
\end{aligned} \tag{5.1}$$

where $Z_{LDO-PDN}$ is the impedance of the system seen by the load and varies for different frequencies. We derive dominant poles and zeros of the impedance and employ them to calculate the step response of the network. Worst case step current, $I_{load\ step}(t)$ is generated by maximizing the current in the positive sections of the impulse response and minimizing the negative impulse sections. Finally, worst voltage drop, V_{max} is obtained from convolution of the $I_{load\ step}(t)$ with $Z_{LDO-PDN}$ in time domain as illustrated in Figure 5.2.

Magnitude of the step current, $\| I_{load\ step}(t) \|$ is bounded by the peak load of the functional block, I_{peak} . LDO power P_{LDO} , is consumed mainly in the form of bias current and is bounded by maximum power budget P_0 . Decoupling capacitors are constraint with minimum C_0 to make feedback loop stable and C_1 and C_2 as area constraints.

5.1.3 LDO Frequency Domain Characteristics

Feedback loop of the on-die LDO and PDN system have three main poles (P_i) and a zero (Z_1) introduced by on-chip and off-chip decoupling, load, power MOSFET parasitics and error amplifier parasitics as shown in Figure 5.4:

$$\begin{aligned}
P_1 &= \frac{1}{2\pi(R_{on\ LDO} + R_{die})C_{die}} \\
P_2 &= \frac{1}{2\pi(R_{die}C_{int})} \\
P_3 &= \frac{1}{2\pi(R_{amp}C_{amp})} \\
Z_1 &= \frac{1}{2\pi(R_{die}C_{die})}
\end{aligned} \tag{5.2}$$

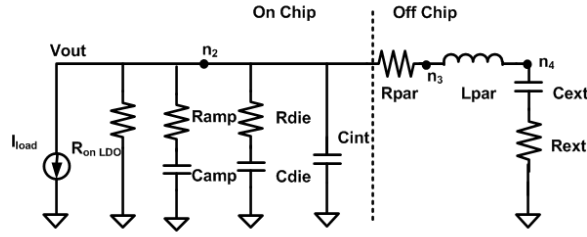


Figure 5.3: Lumped approximation of the LDO-PDN model.

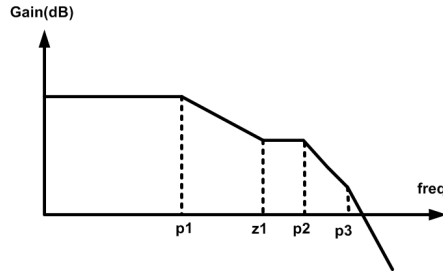


Figure 5.4: LDO poles and zeros impact on the system loop gain.

Description of subscripts are as follows. *die* represents the on-die decoupling capacitor, *int* is the intrinsic load decap, and the LDO ON resistance is shown as *on LDO*. Finally, *amp* shows the error amplifier parasitics as depicted in Figures 5.1 and 5.3.

Impedance Profile of LDO

Output impedance of the LDO is a function of the off-chip and on-chip network and LDO as detailed in Equation (5.3):

$$\begin{aligned}
 Z_o &= [R_{ds} \parallel (\frac{1}{sC_{int}}) \parallel (R_{die} + \frac{1}{sC_{die}}) \\
 &\parallel (R_{ext} + R_{par} + \frac{1}{sC_{ext}} + sL_{ext})] = \\
 &\approx \frac{sC_{die}R_{ds}R_{die} + R_{ds}}{sC_{die}R_{ds} + s^2C_{die}C_{int}R_{ds}R_{die} + sC_{int}R_{ds} + 1}
 \end{aligned} \tag{5.3}$$

where R_{ds} is the drain-source ON impedance of the LDO power MOSFET in the regulation mode. Impedance of the power MOSFET is a function of LDO bias current, I_{ds} and directly contribute to the LDO total power. In Equation (5.4), λ is the channel-length modulation parameter of the power device. Therefore, the

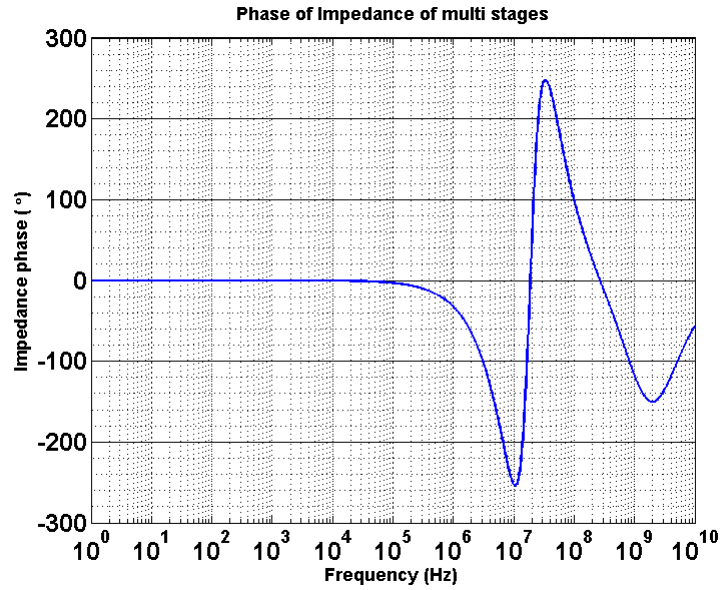


Figure 5.5: Phase of the power distribution impedance $Z(f)$.

more power we burn in the form of LDO bias current, the less we can bring the impedance. As a result the voltage noise is reduced with the cost of more power consumption in the LDO block. Here, λ_1 is a constant co-efficient for power.

$$R_{ds} = \frac{1}{\lambda I_{ds}} \sim \frac{1}{\lambda_1 P_{LDO}} \quad (5.4)$$

Figures 5.5 and 5.6 illustrate the output impedance phase and amplitude of the feedback system in presence of the LDO. Bandwidth of the LDO is roughly early decades of the MHz range. Thus, the error amp can only regulate the output voltage up to MHz . Beyond 10s of MHz range, we can only rely on the power MOSFET and decaps for noise compensation. The impedance profile of the feedback loop is reduced by error amp loop gain ($A_{opamp}(f)$) factor as shown in Equation (5.5).

$$Z_{LDO \text{ closed-loop}}(f) = \frac{Z_{LDO \text{ open-loop}}}{1 + A_{opamp}(f)} \quad (5.5)$$

5.1.4 Analytical Worst-case Voltage Drop

We derive worst case voltage drop from the step response of the system. Impulse response is then obtained from the derivative of the step response. We

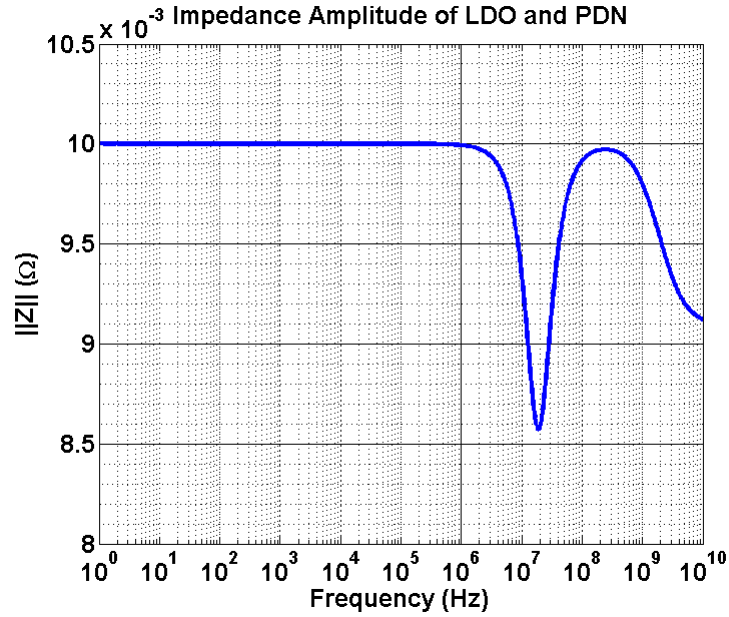


Figure 5.6: Amplitude of the power distribution impedance $Z(f)$.

partition the impulse response into positive and negative time sections. In the positive sections of the impulse response, we maximize the voltage convolution and worst drop by increasing current to the peak. Similarly, in the negative sections of the impulse response, we minimize the peak current. Superposition of the maximum and minimum current, leads to the worst voltage drop of the system. From impedance profile, we identify poles and zeros of the system. We focus on dominant poles and zeros and derive the worst analytical voltage drop. We then apply the flow on the testcase to study LDO power and decap tradeoff.

Several published works propose different methods to find the worst case power supply noise. Ghani *et al.* investigated the upper bound in a vectorless methodology that avoids detailed simulation [16]. Shayan *et al.* focused on the system resonance and introduced resonance-aware current generation [46]. However, the voltage drop obtained in this way might not be the upper bound of noise. Drabkin and Hu *et al.* [10, 21], obtained an aperiodic stimulus from the network's impulse response, which is then used in a time-domain simulation to find the worst-case voltage swings. Most of these works consider the conventional passive power networks. The focus of this work is the interaction between LDO

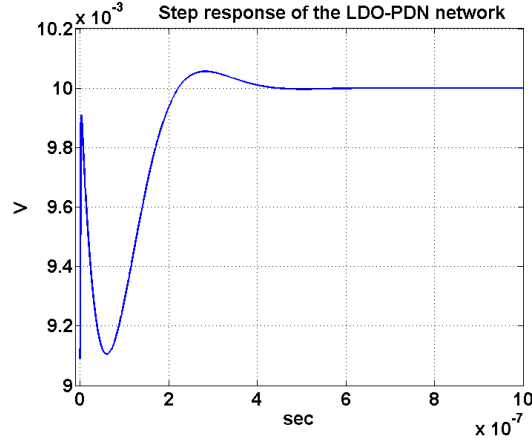


Figure 5.7: Step response of the LDO based power distribution.

and PDN. We derive the exact voltage drop of system. The impedance of the network is represented by three poles and zeros in Equation (5.6):

$$Z_{LDO-PDN}(s) = \frac{(s - z_1)(s - z_2)(s - z_3)}{(s - p_1)(s - p_2)(s - p_3)} \quad (5.6)$$

$$\begin{aligned} z_1(s) &= -2.011 \times 10^9 \\ z_2(s) &= -0.0107 + 0.0156i \times 10^9 \\ z_3(s) &= -0.0107 - 0.0156i \times 10^9 \\ p_1(s) &= -1.8177 \times 10^9 \\ p_2(s) &= -0.0125 + 0.0142i \times 10^9 \\ p_3(s) &= -0.0125 - 0.0142i \times 10^9 \end{aligned} \quad (5.7)$$

We apply the step current to the system impedance in frequency domain as shown in Figure 5.7:

$$\begin{aligned} v(s) &= \frac{1}{s} Z(s) = \frac{k_1}{s} + \frac{k_2}{s-p_1} + \frac{k_3}{s-p_2} + \frac{k_4}{s-p_3} \\ k_1 &= \frac{z_1 z_2 z_3}{p_1 p_2 p_3} = 1.1008 \\ k_2 &= \frac{(p_1 - z_1)(p_1 - z_2)(p_1 - z_3)}{p_1(p_1 - p_2)(p_1 - p_3)} = -0.1011 \\ k_3 &= \frac{(p_2 - z_1)(p_2 - z_2)(p_2 - z_3)}{p_2(p_2 - p_1)(p_2 - p_3)} = A + Bi = 0.0002 + 0.1396i \\ k_4 &= \frac{(p_3 - z_1)(p_3 - z_2)(p_3 - z_3)}{p_3(p_3 - p_1)(p_3 - p_2)} = A - Bi = 0.0002 - 0.1396i \end{aligned} \quad (5.8)$$

Step response is converted back to time domain:

$$\begin{aligned}
v(t) &= k_1 + k_2 e^{-p_1 t} + 2e^{-\alpha t} [A \cos(\beta t) - B \sin(\beta t)] \\
&= 1.1 - 1.1 e^{-1.8 \times 10^9 t} + 2e^{-0.0125 \times 10^9 t} \\
&\times [0.0017 \cos(\beta t) - 0.002 \sin(\beta t)] \\
\alpha &= 0.0125 \times 10^9 \\
\beta &= 0.0142 \times 10^9
\end{aligned} \tag{5.9}$$

5.1.5 Exact analytical value of LDO peak voltage drop

We derive the exact solution for the worst case voltage drop in time domain in Equation (5.10):

$$\begin{aligned}
V_{worst\ case} &= k_1 + k_2 e^{-p_1 t_0} + \\
&+ 2A e^{-\alpha t_0} + 2[A \cos(\beta t_k) - B \sin(\beta t_k) e^{-\frac{\alpha}{\beta} \arctan(2)}] \cdot \frac{1}{1 - e^{-\alpha \pi \beta}} \\
&= 1.2048V
\end{aligned} \tag{5.10}$$

We calculate two different extreme time corners from the peak voltage to obtain t_0 and t_k . First, when t is small we have:

$$\left. \frac{\partial v}{\partial t} \right|_{t\ is\ small} \approx 0 \Rightarrow t_0 = 2.1203 \times 10^{-9} \tag{5.11}$$

Next, when t is reaching infinite we have:

$$\left. \frac{\partial v}{\partial t} \right|_{t\ is\ large} \approx 0 \Rightarrow t_k = \frac{1}{\beta} \left[\arctan\left(\frac{0.003}{0.0017}\right) + k\pi \right] \tag{5.12}$$

5.1.6 LDO-based Experimental Results and Tradeoff Analysis

In the experimental section, we first derive the worst case voltage drop of the system. We then try to minimize worst voltage drop by optimizing LDO power and decap area tradeoff. Increasing the bias current of the LDO (power of the LDO) will reduce the regulation mode impedance of the power MOSFET. We first calculate the worst voltage drop, based on the step response. Figure 5.8-(a) shows the step response of the system. Impulse response is obtained from derivative of step response and is employed to generate the worst load current. Maximum

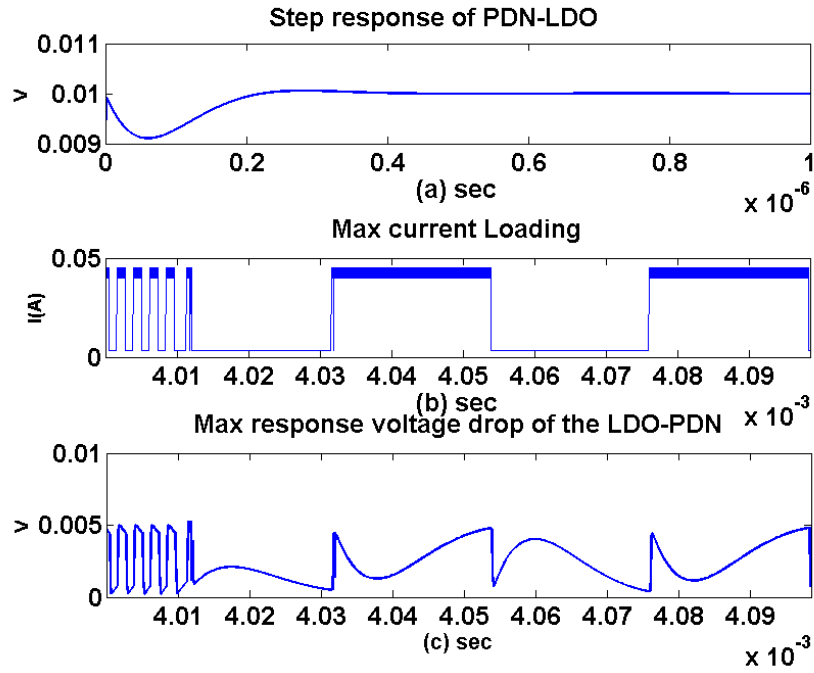


Figure 5.8: Worst-case voltage drop based on LDO step response

load current, $I_{load\ step}(t)$ is generated based on the impulse response as shown in Figure 5.8-(b) and represents dominant frequency peaks in time domain.

Worst voltage drop $V_{max}(t)$ is obtained from convolution of the impulse response and $I_{load\ step}(t)$ (Figure 5.8-(c)). Figures 5.9 and 5.10 illustrate the tradeoff of the LDO power and decoupling capacitors area for voltage drop minimization. Here, voltage overshoot affects device reliability and timing and voltage undershoot results in functional failure. We assume the nominal supply voltage value is $V_{dd} = 1V$ in our analysis. Consuming more power in LDO will reduce the power MOSFET ON resistance R_{ds} , and as a result the voltage drop overshoot and undershoot become smaller. In addition, decap area competes with the LDO power for reducing worst drop as shown here. The main effect of the decoupling capacitors are on middle and high frequencies noise compensation ($> MHz$). While system rely on the LDO feedback loop and power MOSFET for low frequency compensation. Overall as shown in Figures 5.9 and 5.10, we meet the target drop.

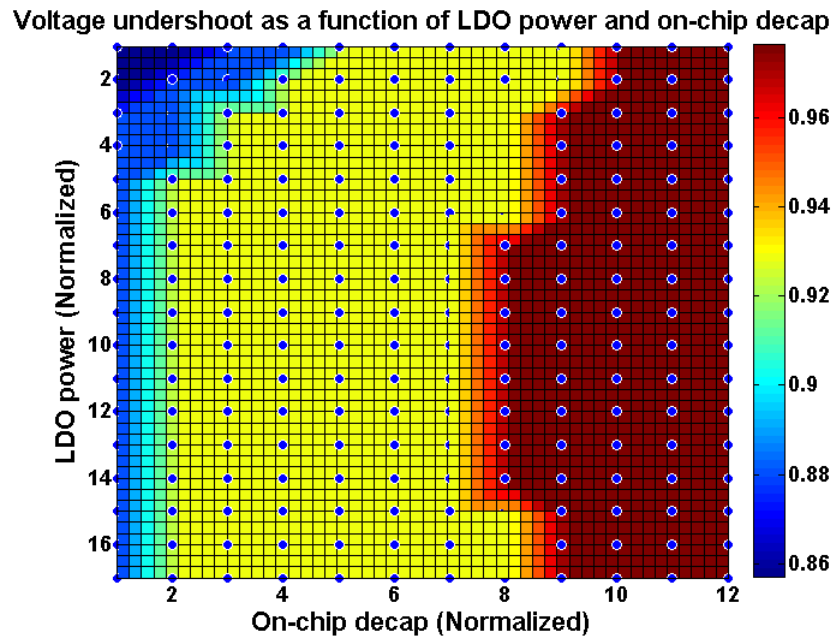


Figure 5.9: Peak voltage drop undershoot as a function of LDO power and on-chip decap.

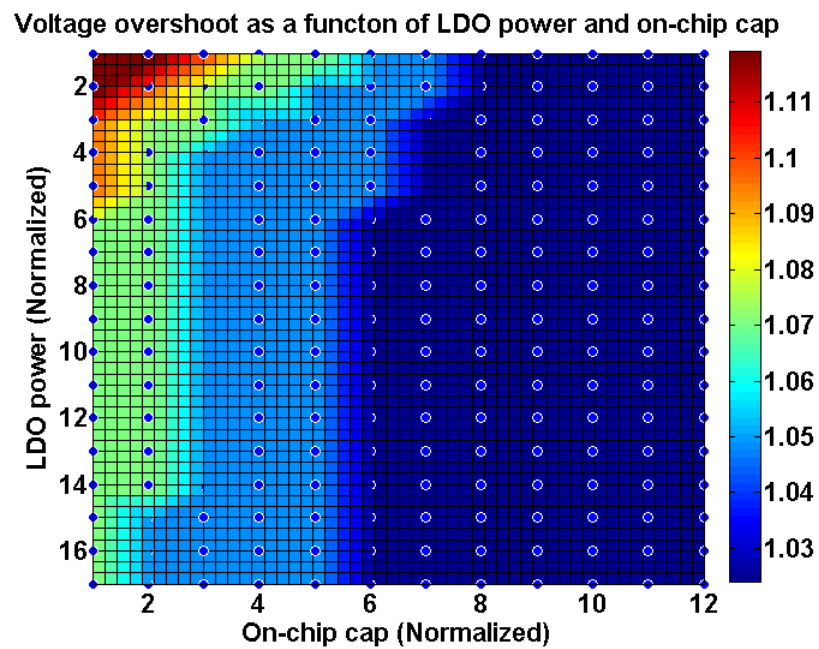


Figure 5.10: Peak voltage drop overshoot as a function of LDO power and on-chip decap.

5.2 Parallel Flow to Analyze the Impact of the Off-chip Voltage Regulators

5.2.1 Overview of Voltage Regulator Impact

With technology scaling and increase in the complexity of VLSI chips, design and analysis of power distribution network have become challenging. Noise margin becomes narrower in successive technology generations which is resulted from decrease in the V_{dd} voltage [6]. A poorly designed power network can become the cause for different types of issues such as loss of circuit performance, temporarily failures, electromigration and chip failures [34, 62].

Inclusion of voltage regulator in the power delivery analysis will increase the voltage drop up to 10% of the nominal V_{dd} based on our experimental results. Due to the low frequency response time of regulator, we need up to micro seconds transient simulation to highlight the power noise. In this work, a custom flow is developed which speeds up slow transient simulation time for longer simulation window. Time domain simulations are pattern dependent where the response of the PDN to a specific stimulus is obtained. From frequency domain simulation the anti-resonance peaks are identified. Therefore for the design of a non-conservative reliable PDN, we need to iterate between both time and frequency domains.

Conventional conservative methodology for designing a good PDN is to define a target impedance for the network that should be met over a broadband frequency [50]. This type of design methodology is significantly expensive and might lead to over conservative designs.

Much research work in academia and industry tries to come up with analysis solutions [35, 44, 49, 62, 64] to address tremendous number of variables for the design of PDN. In this work, we developed an efficient parallel processing analysis flow for the full power distribution network in both frequency and time domain as shown in Figure 5.11. The regulator-based flow enables iteration between both domains simultaneously. We run the flow on FWgrid infrastructure [15] using from one up to 200 processors. The PDN simulations time is reduced from hours down to hundreds of seconds with speedup of up to 430 times. Frequency domain analysis

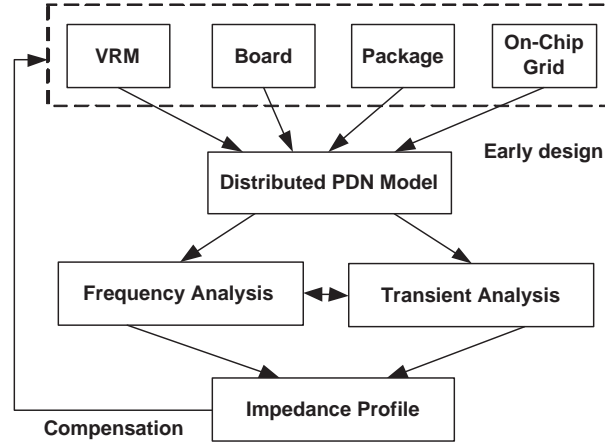


Figure 5.11: Early analysis flow for the complete PDN network.

is performed in a broadband range from DC up to GHz to guide the designers to identify the weak hotspot locations and meet the target impedance for the entire range. The time domain waveform is recovered from the frequency domain. That helps to examine the voltage drop of the PDN during chip performance. Our experimental results show speedup over HSPICE and Cadence Spectre.

In this chapter, we highlight the effect of the non-linear VRM for power delivery noise. VRM has low frequency response in the range of hundreds of KHz. Thus, transient analysis of the complete PDN is challenging and requires many clock cycles. We demonstrate that including the voltage regulator model will increase the voltage drop. VRM effect could not be ignored in the nanoscale technology as V_{dd} is lowered. The power distribution network includes a complete path from voltage regulator module, board, and package to on-chip power grids.

To identify the power noise hotspot spatial distribution, we partition 2D package and on-chip grid model into $m \times n$ sections with high resolution. We demonstrate that this model will help the designers identify the exact location of the resonance peak of PDN much better than uniform model.

5.2.2 Complete Distributed PDN Model

The power distribution network in our analysis is a complete path including non-ideal voltage regulator module, board, and package and die as shown in

Figure 5.12.

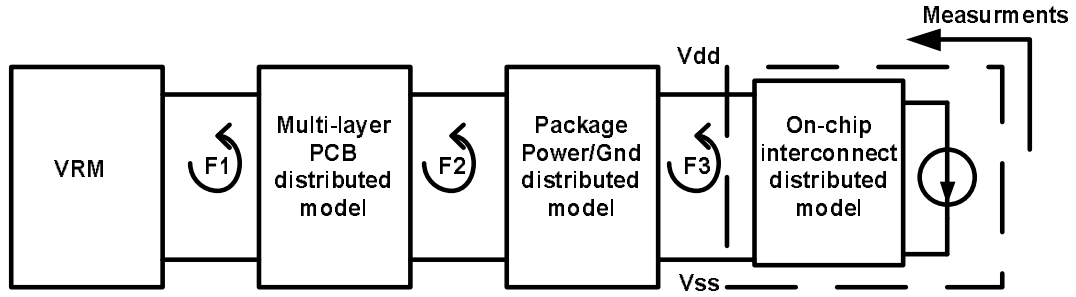


Figure 5.12: Complete power distribution network model.

Current available commercial PDN extraction tools for board, package and on-chip PDN include controlled current and voltage sources in addition to RLC components in their extracted PDN models. Therefore, we extend the flow to modify nodal analysis instead of nodal analysis to model the complete PDN components in frequency domain. We use PowerSI from Sigriety [3] to extract the $RLGC$ model of the board and package.

Voltage Regulator Model

In this section, we discuss details of *switching-mode power supply* (SMPS) model that we adopt in our PDN analysis. A switching mode power supply is a device which transforms one level of voltage from AC power line or unregulated DC power line to another level that is required by the logic circuits.

A switching mode power supply is a power supply that provides the supply voltage using loss-less components such as capacitors(C), inductors(L), and transformers and the switches that have two stages: *ON* and *OFF* mode. Voltage regulator module is required to maintain the target voltage level within a specified tolerance such as $1.0V$ with $\pm 2\%$ ripple.

In today's VRM design the bandwidth is in the range of hundreds of KHz . Thus, VRM is considered as a low frequency module in the PDN [36]. In order to use the VRM in the PDN analysis, we need to derive the closed-loop output impedance (Z_{oc}) of the VRM. Here L is the open loop inductance of the VRM, C

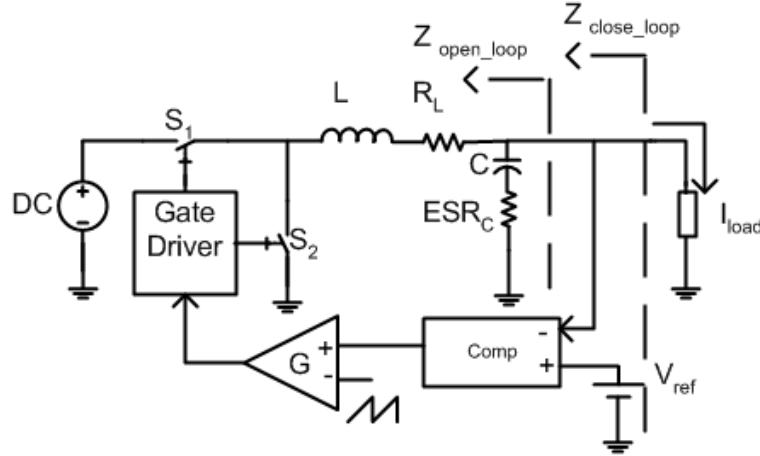


Figure 5.13: Voltage regulator model.

is the output capacitor, R_L and ESR_C correspond to the equivalent series resistors of the L and C as depicted in Figure 5.13. We use small-signal analysis method to calculate the open loop output impedance Z_o as well as the closed-loop output impedance Z_{oc} [58]:

$$Z_o(s) = R_L \frac{(1 + s/\omega_c)(1 + s/\omega_L)}{1 + s/(Q\omega_o) + s^2/\omega_o} \quad (5.13)$$

$$Z_{oc}(s) = \frac{Z_o(s)}{1 + T(s)} \quad (5.14)$$

$$\begin{aligned} \omega_c &= \frac{1}{1 + T(s)}, & \omega_L &= \frac{L}{R_L} \\ \omega_o &\approx \frac{1}{\sqrt{LC}}, & Q &\approx \frac{\sqrt{L/C}}{R_L + ESR_C} \end{aligned} \quad (5.15)$$

In Equation (5.15), R_L includes the DC resistance of the inductor L , conduction resistance of the VRM switches S_1 and S_2 and the parasitic resistance of the traces. The ESR_C is the equivalent series resistance of the VRM output capacitor C . The ω_o is the power stage double poles and $T(s)$ is the closed-loop gain.

In high frequency Z_{oc} is approximately equal to ESR_C . The impedance of the VRM for high frequencies is the same as the open loop impedance in high frequency. Impedance is independent from the feedback transfer function $T(s)$ above VRM bandwidth. The VRM feedback controller only attenuates output

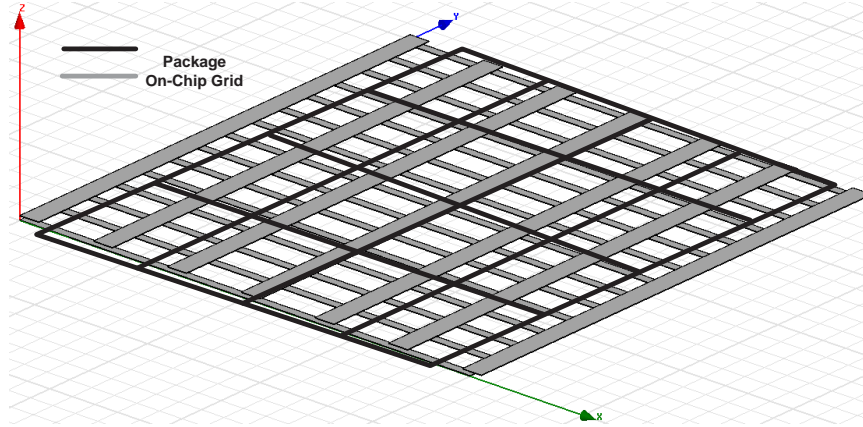


Figure 5.14: 2D partitioned model for package and on-chip grid.

impedance in low frequency range. We use the closed-loop impedance of the VRM model in the analysis flow.

Package and On-chip PDN Partitioned Model

We introduce our PDN partitioned model for the package and die in this section. During PDN model extraction, we partition both package and chip PDN into $n \times m$ sections based on the physical floorplan resolution. Partitioning will reflect the spatial distribution of noise and will help designers in PDN optimization. Each partition represents finer spatial distribution of the PDN properties. Based on the experiment results, this is required especially in the low power chip in order to optimize the design of the PDN. Power noise sensitivity varies as a function of area. In our testcases based on the package and chip topology n and m are 4×4 as illustrated in Figure 5.14. In the next section, the details of the MNA-based analysis methodology and flow are explained.

5.2.3 MNA-based Early PDN Analysis Flow

We implement parallel processing using *Message Passing Interface* (MPI) to speedup the matrix solving. In the proposed early analysis, PDN is modeled in frequency domain applying modified nodal analysis. To recover the time domain voltage response, vector fitting method is applied. Simulation time is reduced

compared with HSPICE transient analysis and Cadence Spectre. We ran multiple clock cycles which is the bottleneck when the PDN includes low frequency modules such as VRM and PCB.

Zhang *et al.* [60] proposed an analysis method for the on-chip *RLC* power ground models for multiple clock domains based on the frequency domain for single processor. They applied Laplace transform to the current sources. The log scale sampling points are solved with nodal analysis for any arbitrary RLC circuit in the frequency domain.

We adopt their technique, and enhanced and optimized it for the extracted PDN model. The extracted model includes controlled current and voltage sources. Thus, we need to extend the flow to modified nodal analysis instead of nodal analysis. Also, to reduce the CPU time, we implemented the flow using parallel processing. In each of the processing threads, the solution from previous run is adopted as the initial guess for the GMRES iterative solver to reduce the iteration convergence time. In the following section we describe the methodology of our MNA analysis.

5.2.4 Modified Nodal Analysis-based Frequency Domain Analysis

The goal of our analysis flow is to identify the resonance frequencies and the power grid noise due to IR drop, ground bounce and di/dt . Power grid networks consist of the resistance, capacitance and inductance of the power grids which introduce multiple anti-resonance frequencies [34]. We derive the impedance profile for the entire power grid network which includes the closed-loop path from VRM to die grids.

First, we convert the current sources which represent the die transistors switching behavior from time domain to frequency domain representation. We assume that the input stimuli $I(t)$ is described in a *piecewise linear* (PWL) form. To convert the time domain stimulus into frequency domain, $I(s)$, Laplace transform is applied. Complex modified nodal analysis matrix is then formed. We solve the matrix iteratively to compute the frequency domain voltage response. Finally, we

convert the voltage $V(s)$ using vector fitting to approximate the discrete samples in partial fractional expression to time domain representation $V(t)$.

Modified nodal analysis is an efficient method to process voltage sources, and current and voltage dependent circuit elements. The voltage $V(s)$ of the PDN nodes is calculated as:

$$\begin{bmatrix} Y_R & B^T \\ C & D \end{bmatrix} \begin{bmatrix} V(s) \\ J(s) \end{bmatrix} = \begin{bmatrix} I(s) \\ E(s) \end{bmatrix} \quad (5.16)$$

where Y_R is a reduced form of the nodal matrix excluding the contributions of the voltage sources, current controlling elements, etc. B^T contains the stamps from the current equations with respect to the additional current variables. Thus, it contains $+/-1$ for the elements whose branch relation is introduced. The branch constitutive relations, differentiated with respect to unknown vectors are presented in by matrix C and D . It is noted that the zero-nonzero pattern of C is basically the same as that of B^T except for some nonreciprocal elements [55]. The vector $I(s)$ and $E(s)$ are the excitations from input stimuli. $V(s)$ is the node voltage vector and $J(s)$ is the input current source vector. We apply the iterative solver *Generalized Minimal Residual Method* (GMRES) method with *Incomplete LU* (ILU) as pre-conditioner to speedup complex matrix solving. We chose GMRES with ILU preconditioning because of the faster convergence compared to the other available iterative solvers and direct solver for the *Modified Nodal Analysis* (MNA) complex matrix.

The logarithmic scale is used to sample the output voltage response in frequency domain. To reduce convergence time, we use solution of the previous iteration as the initial guess for the GMRES next iteration.

In the next flow step, the time domain voltage drop of each node in the PDN is recovered.

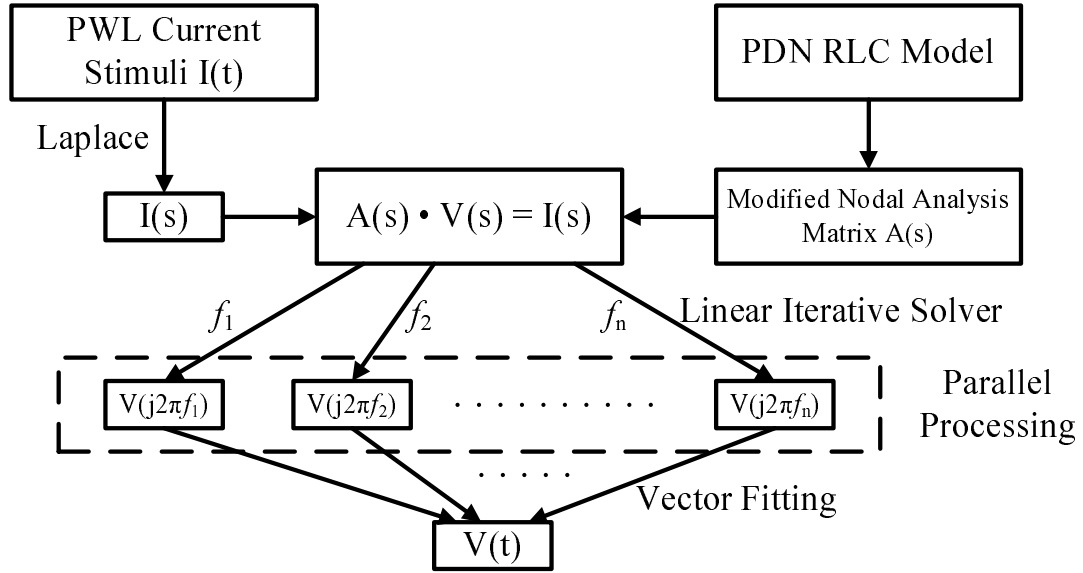


Figure 5.15: Framework for calculating the voltage response of the P/G networks.

5.2.5 Time Domain Signal Recovering Using Vector Fitting

Because the clock frequency of the ICs is reaching GHz range and the transistor current load includes many high frequency components, current signature is complex. Thus, conventional vector fitting method [9] cannot fit high frequency components well and the time domain recovered result has a large error (ΔV_{VF}):

$$\Delta V_{VF} = V(f) - V_{VF}(f) \quad (5.17)$$

where ΔV_{VF} is the deviation of the fitted frequency domain approximation and the original signal. We enhanced the vector fitting process and use the remainder of the each vector fitting ΔV_{VF} iteration as the next iteration input and perform the vector fitting process iteratively until ΔV_{VF} reaches the acceptable error boundary, i.e., less than $1e - 15$.

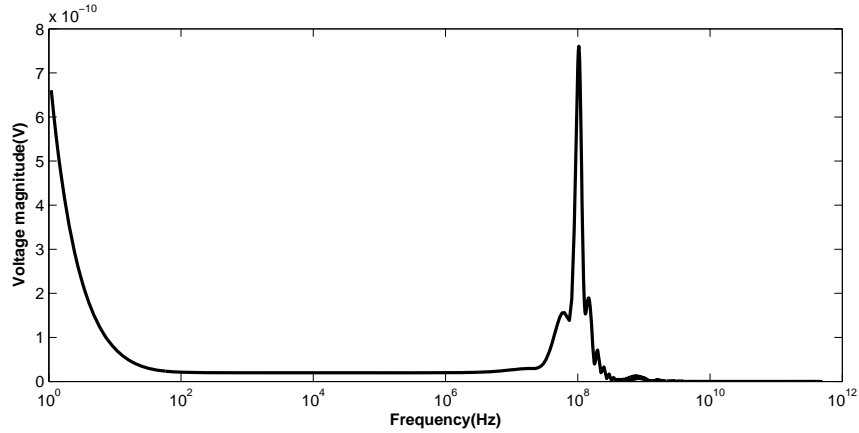


Figure 5.16: Magnitude of voltage response of PDN in frequency domain without VRM.

MNA-based Impedance Profile of PDN

To obtain the impedance profile of a given node for the broadband frequency range, we apply $I(s) = 1$ current to the node of interest. In this case, from modified nodal analysis equation $Z(s) \cdot I(s) = V(S)$, we can see that the node voltage $V(s)$ corresponds to the impedance of the node of interest in that frequency range.

Parallel Processing Flow

We implemented the matrix solver using message passing interface (MPI) to take advantage of the speedup of parallel computing. We use FWgrid infrastructure [15] to run our parallel flow. We run the flow by using from 1 up to 200 processors. In FWgrid, there are 94 dual 1.6GHz Opteron nodes with 2GB RAM, 2 dual 2.2GHz Opteron nodes with 16GB RAM and 224 dual Xeon (64× 2.8GHz & 160× 3.2GHz) nodes with 4GB RAM. Totally it provides 640 processors, 1.15TB of Memory and 160TB of storage. Currently, multiple core machines are also available in market for a low cost.

5.2.6 VRM-based PDN Results and Analysis

In this section, we will discuss the results from running our flow on the power distribution network model. We first show the results for the package and power

grid model connected to an ideal voltage regulator. Then for next section, we add the details of the board and VRM impedance Z_{oc} as we discussed in Section 5.2.2. We observe the low frequency anti-resonance peak (in KHz range) resulted from the regulator model and PCB. The IR drop is worsened after adding the model for the VRM and the transient result takes longer time compared to ideal VRM to reach the steady state mode. Therefore this would be a more accurate and conservative early design analysis as we consider the worst-case voltage drop. This analysis is very important particularly for the low power platforms power grid design with low V_{dd} and narrow noise margin. The measurement and observation points are from the on-chip power grid to VRM.

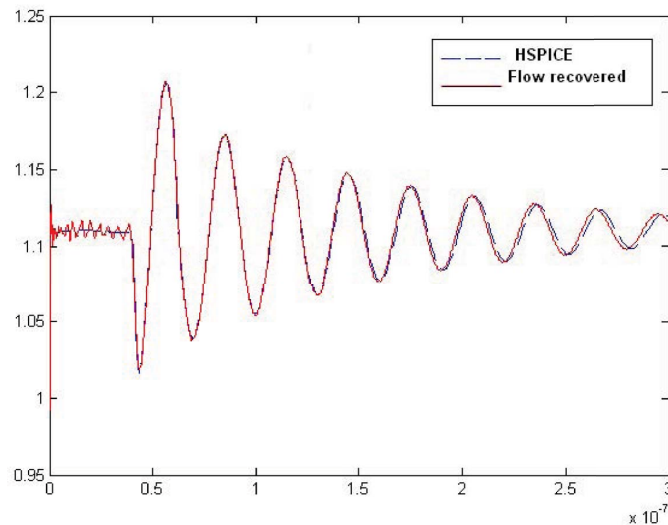


Figure 5.17: Recovered time domain voltage response comparison with HSPICE.

Ideal Power Distribution Network Model

We first use the ideal regulator with package and on-chip power grid in our analysis. The current stimulus that we apply to the PDN has duration of $20nsec$ starting from $40nsec$. Current represents the switching current of the transistors in the chip operation mode and includes a wide range of frequencies. The frequency domain results show that the ideal PDN network has a resonance peak in the range

of $100MHz$ because of package as shown in Figure 5.16.

Figure 5.17, shows the comparison between our time domain recovered signal using vector fitting and HSPICE transient simulation. Table 5.1 and 5.2 are the transient simulation time of our flow for $300nsec$ compared with HSPICE and Cadence Spectre simulations. In our method each frequency sample solution takes about $16sec$ for P/G case 1 with GMRES and ILU as the preconditioner. Analysis requires a total of 200 frequency sampling points to recover the time domain wave form. We ran the flow on both single processor and multiple processor Linux machines as shown in Tables 5.1 and 5.2.

Table 5.1: Simulation time and speedup of proposed frequency flow versus HSPICE and Cadence Spectre for $300nsec$ using 1 and 4 cores.

Test	HSPICE	Spectre	1 proc.		4 proc.	
	<i>sec</i>	<i>sec</i>	<i>sec</i>	Speed Up	<i>sec</i>	Speed Up
case1	44302	1601	2509	17.66	673.98	65.73
case2	55117	3120	2856	19.30	884.88	62.29
case3	62588	3218	2962.6	21.13	899.26	69.60
case4	71643	10918	3226.26	22.21	1030	69.56

Table 5.2: Simulation time and speedup of proposed frequency flow versus HSPICE and Cadence Spectre for $300nsec$ using 8 and 16 cores.

Test	HSPICE	Spectre	8 proc.		16 proc.	
	<i>sec</i>	<i>sec</i>	<i>sec</i>	Speed Up	<i>sec</i>	Speed Up
case1	44302	1601	243.15	182.2	115.62	383.17
case2	55117	3120	525.47	104.89	291.3	189.21
case3	62588	3218	520.14	120.32	268.7	232.39
case4	71643	10918	530	135.18	292.5	244.93

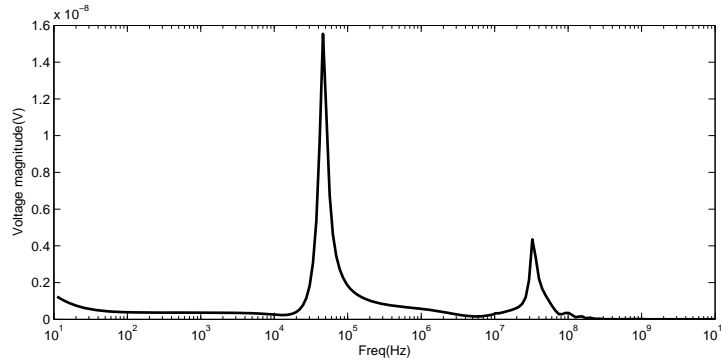


Figure 5.18: Frequency domain resonance peaks of the PDN with VRM.

5.2.7 Voltage Regulator Module Impact on Noise

In this section, we add the details of the voltage regulator model to the PDN. The regulator bandwidth is in the range of hundreds of KHz . From the frequency domain simulation results, we observe additional low frequency resonance peak for the voltage frequency response as illustrated in Figure 5.18. This is due to the VRM-board and package-board anti-resonance phenomena [34]. In Figure 5.18, the KHz peak, the ten MHz peak, and hundreds MHz peak are respectively from the VRM and PCB, PCB and package, and package and on-chip grid anti-resonant loop.

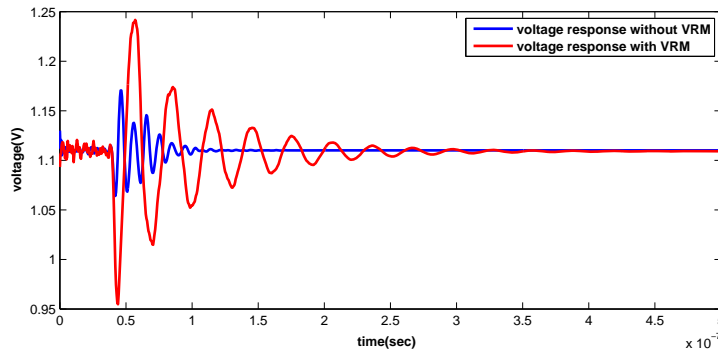


Figure 5.19: Comparison of time domain voltage response with and without VRM.

The IR drop is worsened after adding VRM model in the flow. Time domain transient result takes longer time which is about $300nsec$ compared to the ideal

case which used to be $100nsec$ as shown in Figure 5.19. To reduce the IR drop within the tolerable noise margin, we allocate decoupling capacitors as discussed next.

Broadband Impedance Profile

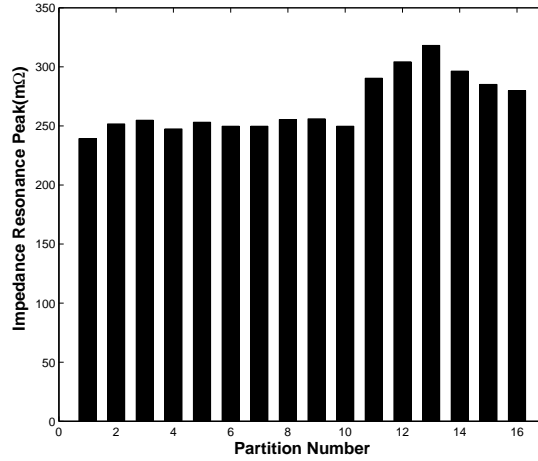


Figure 5.20: Impedance profile for the 4×4 2D PDN partitions.

The broadband impedance profile helps us to identify the location in the PDN which needs compensation. Figure 5.20 is the impedance profile resonance peak for the 16 partitions in the PDN based on the 4×4 model in Section 5.2.2. From the profile, the minimum impedance is about $240m\Omega$ and the maximum is $317m\Omega$ with about $77m\Omega$ difference. Based on the current consumption and impedance profile of the hard macro in each partition, we calculate the amount of decoupling capacitor of that region. To calculate the amount of necessary decoupling capacitances in order to maintain a limited voltage drop, one can use the following:

$$P = C_T \cdot V_{dd}^2 \cdot f \cdot p_{0-1} \quad (5.18)$$

where P is the total chip power consumption, V_{dd} denotes the supply voltage, f is the clock frequency, C_T is the effective chip capacitance, and p_{0-1} is the probability that a $0 - 1$ transition occurs. The original design has a low frequency impedance

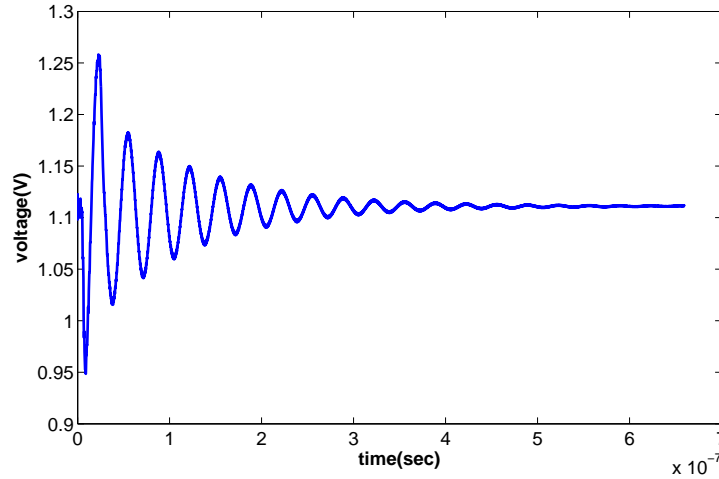


Figure 5.21: PDN voltage response over chip operation current profile.

peak of 1.17Ω at about $158KHz$. This was resulted from the VRM and board low frequency loop resonance. We reduce the peak impedance from 1.17Ω down to less than $250m\Omega$ by using bulk decoupling capacitors with minimum size, numbers and ESL value.

In our analysis, first we analyze the PDN under chip switching current profile. Then we repeat the current stimulus with the resonance frequency period of the PDN to test the PDN design under worst-case current profiles.

Figure 5.21 is an example where the V_{dd} node of the on chip PDN takes about $660nsec$ to respond to the chip load transition with single stimulus of $20nsec$. V_{dd} node on-chip has a resonance peak at $30MHz$. We repeat the current load with $30MHz$ to see the PDN transient response under the switching current with resonance frequency. The peak IR drop increases from $0.14V$ to $0.25V$ compare to the original single current stimulus as shown in Table 5.1. Up to micro second simulation period is required to complete the run.

Simulation Time

We run the flow on the 3D power grid test structure with the dimension of 32×119 nodes with 5 to 10 stacked metal layers. Figure 3.10 shows simulation time versus number of processors for two cases on FWgrid machine with up to 200

processors. Simulations of testcase 1 and 2 in HSPICE take 9911sec and 110479sec correspondingly for 300nsec simulation. Simulation time is reduced to hundreds sec using parallel flow as illustrated in Figure 3.10.

5.3 Summary

In Chapter 5, we introduced an LDO-based design of the power distributions with main focus on the impedance and LDO power. We propose a methodology to define the worst voltage drop in presence of the LDO using step response. The exact analytical worst voltage drop is outlined. The flow is adopted for optimization of voltage noise of the PDN. The voltage drop minimization is shown as a convex problem.

In the second part of this chapter, a parallel flow for analyzing the impact of the off-chip VRM is introduced. The parallel flow reduces the CPU time significantly from hours to less than hundred seconds. We run the parallel flow on the clustered Linux machine using up to 200 processors to reduce simulation time down orders of magnitude. From the parallel flow, we obtain the speedup of up to 430 times over SPICE simulation. The simulation results are compared with commercial simulators such as Cadence Spectre. We discussed the impact of the voltage regulator model in the complete power delivery path which leads to increase of the voltage drop and longer settling time. Regulator design analysis is crucial especially in nanoscale low power PDN design with low V_{dd} values.

5.4 Acknowledgments

Chapter 5 is based on the following publications:

- A. Shayan, X. Hu, H. Peng, W. Yu, W. Zhang, C.K. Cheng, M. Popovich, X. Chen, L. Chua-Eaon, Xiaohua Kong, “Parallel Flow to Analyze the Impact of the Voltage Regulator Model in Nanoscale Power Distribution Network”, *The International Symposium on Quality Electronic Design (ISQED)*, 2009.

- A. Shayan, X. Hu, C.K. Cheng, W. Yu, C. Pan, “Linear Dropout Regulator based Power Distribution Design under Worst Loading”, *IEEE International Conference on ASIC*, 2011.

The dissertation author was the primary researcher and author, and the co-authors involved in the above publications directed, supervised, and assisted in the research which forms the basis for that material.

Chapter 6

Power Distribution Impact on Performance

Final goal for power distribution design is to have a predictable performance. In this chapter, we compare the pre-silicon analysis with post-silicon measurements. The sensitivity of F_{max} to quality of network is discussed. We introduce our model for predicting worst-case performance under voltage and temperature variation. The model enables designers to have an assessment of the PDN design quality.

Integrity of power delivery network has long been one of the key design concerns for high power high performance CPUs operating in multi-GHz with tens of amperes. Arabi *et al.* documented the sensitivity of CPU performance to voltage noise in different frequency ranges [38]. Waizman developed a method to accurately measure CPU PDN impedance using FFT and clock gating [53]. Kantorovich and Houghton outlined maximum tolerable power supply noise for high performance CPUs [24]. Because of the low average current consumption that is often less than 1A, low power processors are not expected to have the same sensitivities to PDN impedance as high performance processors in GHz range. In Chapter 6, we quantify performance sensitivity of the low power application processor based on the quality of its power delivery network.

Rest of this chapter is organized as follows. In the first section of this chapter, we discuss our results from estimation of power integrity impact to low power

processor performance through *Pre-Silicon* simulation and *Post-Silicon Measurements*. Section 6.1.1 describes an improved frequency domain simulation flow that was used to assist silicon, package, and board designs. Section 6.1.2 outlines measurement techniques and the correlation between low and middle frequency power integrity with maximum performance.

In the second part of this chapter, we present our performance predictive model: In Section 6.2.2, we review and contrast prior related work. Section 6.2.3 describes our implementation flow and the scope of our study. Section 6.2.4 describes our modeling methodology using machine learning-based regression techniques. In Section 6.2.6 we describe the impact of different parameters on gate delay and output slew, and present our proposed worst-case performance model. In Section 6.2.8 we validate our proposed models against SPICE simulations. Finally, Section 6.3 concludes this chapter.

6.1 Estimation of Power Integrity Impact to Low Power Processor Performance through Pre-Silicon Simulation and Post-Silicon Measurements

We proposed the enhanced pre-silicon simulation to accurately capture the power delivery quality. The post-silicon performance measurements correlate well with pre-silicon analysis and demonstrate that power integrity could impact performance up to 15% in low and high frequencies.

6.1.1 Power Delivery Network Pre-Silicon Analysis

Accurate PDN simulation [40, 3] is a key to assist making design decisions by estimating impact to processor performance. Furthermore, in [53] Waizman developed a methodology for measuring PDN impedance, and showed correlation between measurement and one commercial tool. We improved the frequency do-

main PDN simulation flow in the course of designing the processor chipset.

One weakness in the frequency domain analysis flow was die modeling. For high performance CPUs, the need for detailed die model is less because on-die power grid is designed as uniform and robust as possible. Low power processors, on the other hand, are often designed as part of a larger SOC and can be quite non-uniform. In modeling of application processor, we enhanced the die model by allowing current to flow top to bottom through vias, left to right through even metal layers, and front to back through odd metal layers, as shown in Figure 6.1. We also greatly increased tile granularity to better capture on-die power grid and capacitance non-uniformity.

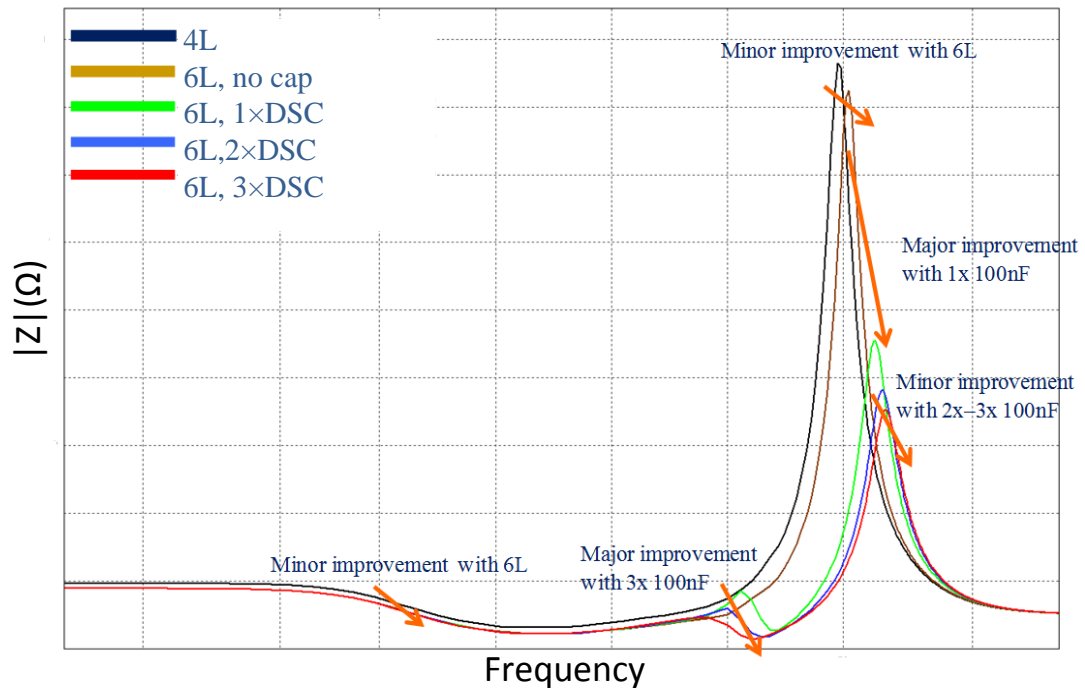


Figure 6.1: System impedance sensitivity due to package layer scaling and die side cap scaling¹.

For low power processors, package and silicon level power grid often are not designed as robust as high performance processors. Greater grid resistance can

¹Over 6500 tiles can be used to model a die. These tiles were used to model the application processor power grid. We simulated the complete PDN model using a 4-layer package and a 6-layer package. We also studied the impact of including 1, 2, and 3 on-package die-side caps.

lead to local impedance hotspots. To validate this hypothesis, we compared our solution to a commercially available solution. The commercially available solution partitions the silicon to coarse $M \times N$ tiles. An electrically equivalent circuit model is then built for each tile with matching transient response. The model we proposed is a physical representation of the on-die power grid. To show the need for fine resolution, we first simulated a die with 3×6 tiles using the commercially available equivalent circuit tile solution and then simulated the design again with over 6500 tiles using our physical model solution. As shown in the first graph on left side in Figure 6.2, the impedance obtained from the two die models did not show any correlation. We found the reason for the discrepancy was that with only 18 equivalent circuit tiles, the impedance at one of the port was heavily dependent on whether the tile clipped a large memory region. To verify this, we regrouped some of the 6500 tiles in the physical model differently with higher circuit granularity. As shown in Figure 6.2, impedance correlation improved dramatically between the two solutions.

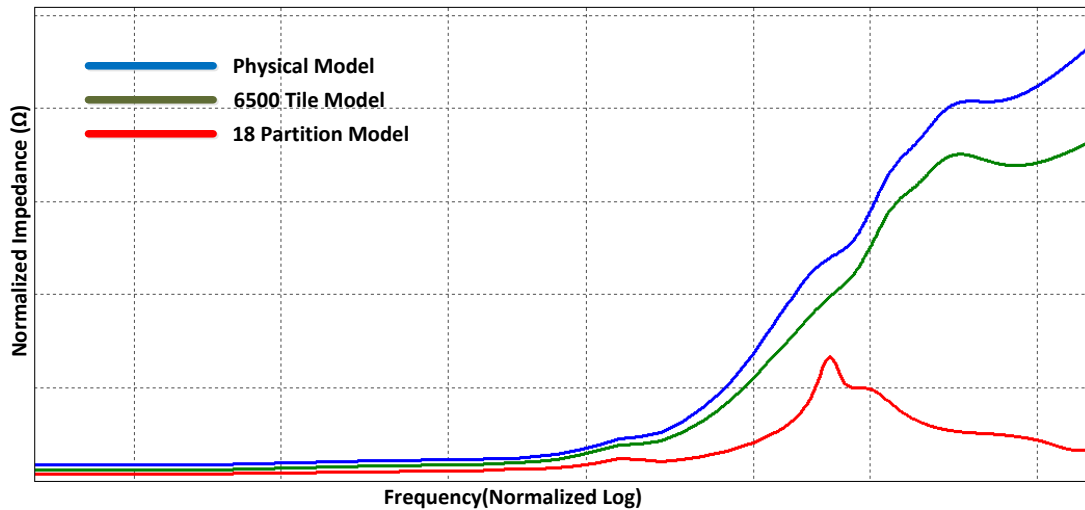


Figure 6.2: Miscorrelation between commercially available equivalent circuit model and physical model we developed².

²The equivalent circuit model used 18 tiles to partition the die, while the physical model used 6500.

6.1.2 Silicon Measurement for Impact of Power Integrity on Processor Performance

Post-silicon measurement is critical for correlation. We will look at the impact of reduction of impedance on F_{max} for both low and high frequency ranges. Our post-silicon measurement consists of two main steps:

1. Direct silicon measurement of the processor impedance using impedance measurement technique described in [53]. Impedance measurement reconstruction with three main harmonics correlates the simulation data as shown in Figure 6.3.
2. Measurement of the processor F_{max} while running different types of application code with different power activities.

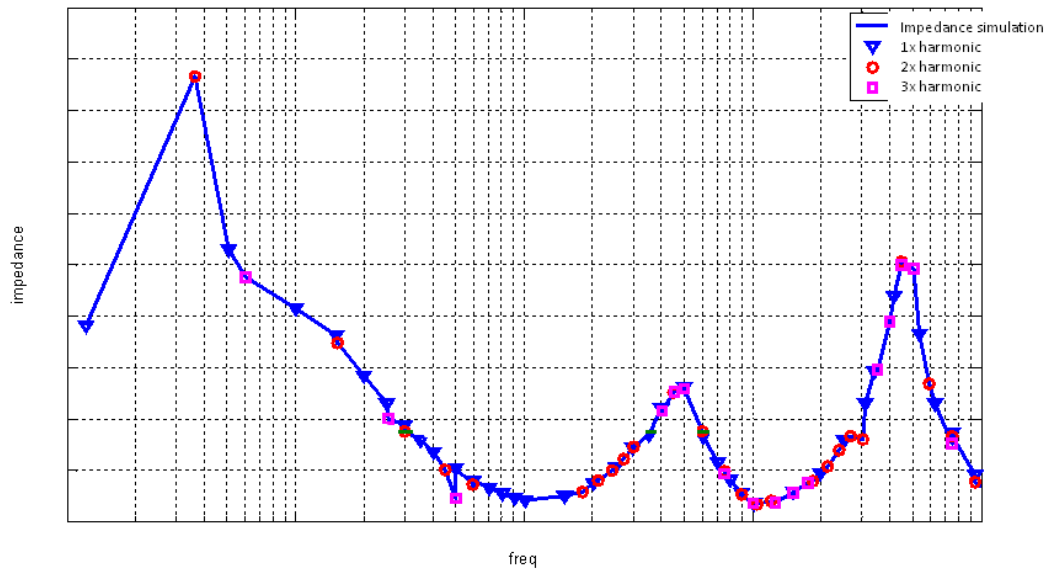


Figure 6.3: Impedance measurement correlation with pre-silicon simulation using three main harmonics.

6.1.3 Mid-High Frequency Impedance Resonance Impact on Performance

Middle frequency resonance is mainly influenced by package and board PDN quality. Mid-frequency impedance could be improved with improving package stacking and adding proper decoupling capacitor on die and package. Figure 6.1 depicts the systematic removal of die side capacitor can increase the impedance peak in the mid-frequency significantly. We systematically reduce the number of DSC from $3\times$ down to $1\times$ while measuring the F_{max} . Table 6.1 is the table showing the impact of DSC cap removal on improving the impedance and frequency. Single DSC can significantly reduce impedance mid-frequency peak. F_{max} reduction is from 3.6% to 12% in different supply voltage as we remove the DSCs.

According to Table 6.1, marginal impact to F_{max} reduction is noted until all on-package decaps are removed. This finding correlates to the simulated PDN impedance in Figure 6.1, where adding the first on-package decap led to the largest high frequency impedance drop.

6.1.4 Low Frequency Impedance Impact on Performance

Time domain simulation of low frequency voltage droop is challenging, and is usually ignored due to large problem size (full die) and long simulation time (up to $100ms$). However, we found that low voltage drop could be as important and can influence the performance of the processor. Pre-silicon frequency domain impedance simulation can highlight a potential low frequency voltage droop. A test code can be used post-silicon to verify the low frequency droop and document the impact on F_{max} .

Figure 6.4 and Table 6.2 depict addition of bulk decaps can improve the impedance profile in low frequency range. Increasing the number of bulk decap systematically from $1\times$ to $15\times$ could translate to maximum -15% frequency degradation for $1\times$ DSC and up to 5% F_{max} improvement for $13\times$ DSC. Figure 6.4 and Table 6.2 summarize the improvement of the F_{max} due to the enhancement of impedance with board bulk decap across fast, nominal and slow process. F_{max}

Table 6.1: Measured F_{max} sensitivity to systematical removal of package capacitors.

Performance degradation(%)			
VDD(V)	3-DSC to 2-DSC	2-DSC to 1-DSC	1-DSC to 0-DSC
0.850	-9.09%	-9.09%	0.00%
0.875	0.00%	0.00%	0.00%
0.900	-7.14%	-7.14%	-7.14%
0.925	0.00%	0.00%	-6.67%
0.950	0.00%	0.00%	-6.25%
0.975	0.00%	0.00%	-5.88%
1.000	-5.56%	-5.56%	0.00%
1.025	-5.00%	-5.00%	-10.00%
1.050	0.00%	0.00%	-9.09%
1.075	0.00%	0.00%	-8.70%
1.100	0.00%	0.00%	-8.33%
1.125	0.00%	0.00%	-12.00%
1.150	-3.70%	-3.70%	-7.41%
1.175	0.00%	0.00%	-7.41%
1.200	0.00%	0.00%	-3.70%
1.225	-3.57%	-3.57%	-3.57%
1.250	0.00%	0.00%	-3.57%
1.275	0.00%	0.00%	-3.57%
1.300	0.00%	0.00%	-3.57%
1.325	0.00%	0.00%	0.00%
1.350	0.00%	0.00%	0.00%

impact for slow process is greater than other processes due to the larger variation space.

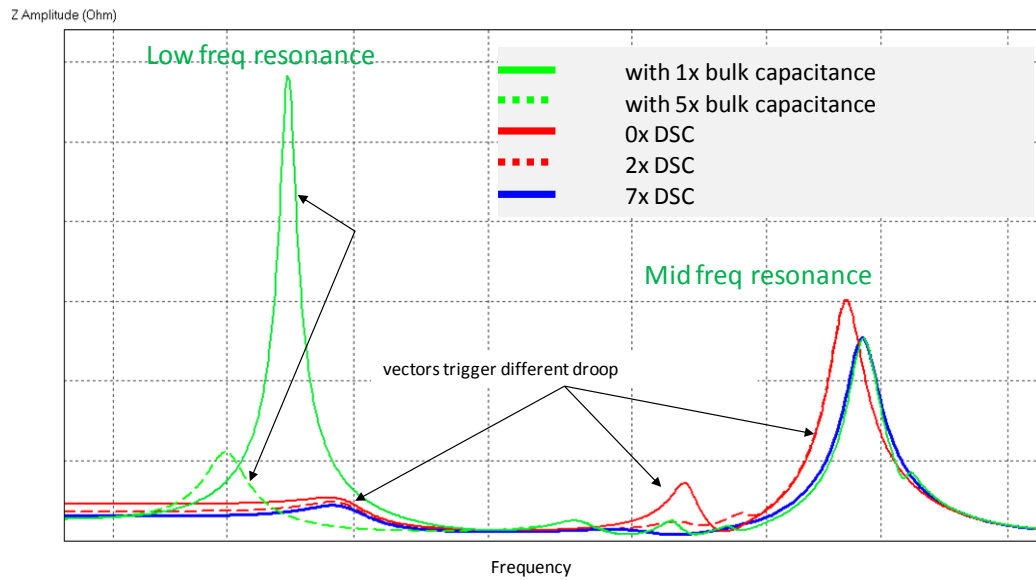
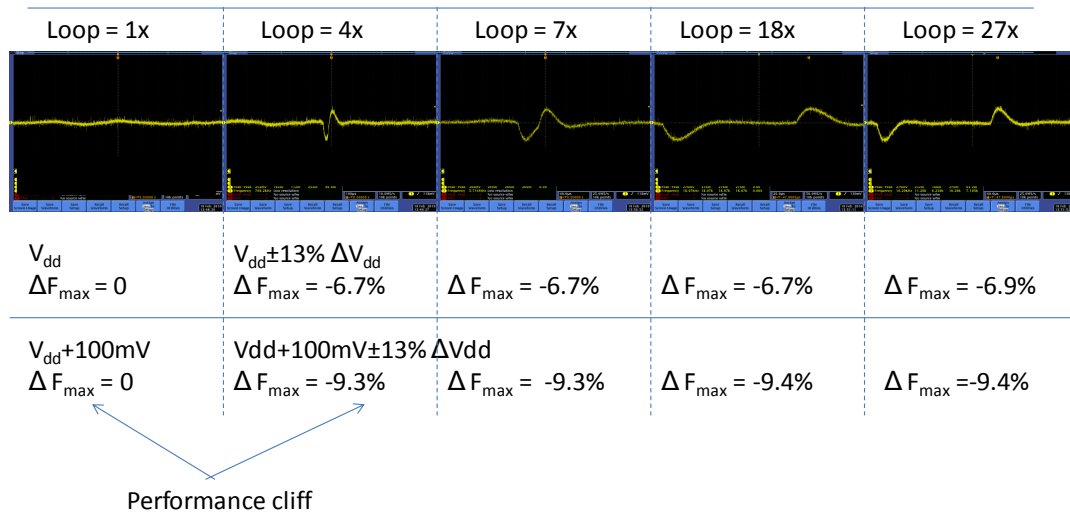


Figure 6.4: Impact of systematic increase of bulk cap on improving the low frequency impedance on F_{max} across different processes.

We ran multiple codes for many loop cycles to exercise the low frequency impedance resonance. Figure 6.5 illustrates that after running the code from $1\times$ to $27\times$ loops, we may observe up to -9.4% F_{max} degradation. This highlights the importance of maintaining low PDN impedance in the low frequency range.

Table 6.2: Bulk decap impact on the performance.

Process	F_{max} (Normalized)	F_{max} degradation from ideal Freq		
		1× bulk cap	5× bulk cap	15× bulk cap
Slow	0.88	-11.42%	-0.46%	3.20%
Slow	0.96	-15.48%	-5.44%	-2.09%
Slow	0.92	-15.21%	-4.72%	-1.22%
Nominal	1.00	-10.61%	-2.73%	0.64%
Nominal	1.00	-10.61%	-2.73%	0.64%
Nominal	1.04	-11.20%	-3.47%	-0.23%
Fast	1.12	-4.66%	-2.08%	4.66%
Fast	1.16	-7.96%	-1.45%	1.04%

**Figure 6.5:** Measurement results for *Low* frequency voltage droop impact on F_{max} for multiple clock cycles and loops.

6.2 Worst-case Performance Prediction under Supply Voltage and Temperature Variation

Power delivery network is a major consumer of interconnect resources in deep-submicron designs (i.e., more than 30% of the entire routing area) [57]. Hence, efficient early-stage PDN optimization enables the designers to ensure a desired power-performance envelope. On the other hand as technology scales, gate delays become more sensitive to power supply variation. In addition, emerging 3D designs are more prone to supply voltage and temperature variation due to increased power density. In this chapter, we develop accurate inverter cell delay and output slew models under supply voltage and temperature variation. Our models are within 6% of SPICE simulations on average. We use our single-cell delay and output slew models to estimate the delay of a path (i.e., an inverter chain, etc.). We also present a methodology to find the worst-case input configuration (i.e., input slew, output load, cell size, noise magnitude, noise slew, noise offset and temperature) that causes the delay of the given path is maximized. We believe that our models can efficiently drive accurate worst-case performance-driven PDN optimization.

6.2.1 Overview of Dynamic Supply Noise

In sub-65nm designs, power/ground voltage level fluctuations (PG noise) has become a primary concern for power integrity as circuit timing becomes more susceptible to supply voltage noise. Thus, designers must take into consideration the impact of supply voltage noise to ensure successful chip design [33]. Rising supply voltage variation has become a challenge for *power distribution system* (PDS) verification. Typically, PDS verification is based on simulation; however, all possible current waveforms and load circuits are not known early in the design cycle. Hence, it is important to develop methods of accurately predicting worst-case supply voltage noise to ensure that the design timing is met.

Existing works [11, 61, 59] on supply voltage noise and its implications on power distribution network optimization or PDS verification are oblivious to the timing impacts of supply voltage noise. In this work, we develop early-stage closed-

form performance models under supply voltage and temperature variations that aid designers to assess the impact of their PDN design choices on the performance of the design. Timing degradation due to PG noise is often estimated by considering voltage drops through static IR-drop analysis. However, these analyses fail to capture the dynamic behavior of the supply voltage noise.

On the other hand, temperature variation affects transistor characteristics including threshold voltage, drive current, drive resistance, and off-current. Hence, it is important to accurately model the impact of temperature on circuit performance. Existing literature [23, 17] propose closed-form expressions that consider the impact of temperature on cell delay; however, in this work we consider the combined effect of supply voltage and temperature variation on circuit performance.

In addition, emerging 3D designs are more prone to supply voltage noise due to increase in power/current demand and variations among tiers. Compensation of the supply voltage variation requires a fair amount of the silicon real estate (e.g., decoupling capacitance allocation, etc.), routing resources, and increased packaging cost. Increased power density in 3D designs also requires close attention to the impact of temperature on circuit performance. Hence, to guarantee a given performance envelope, designers need to characterize the impact of supply voltage and temperature variation on circuit timing. Furthermore, [40] points out to a number of problems caused by dynamic effects of supply voltage noise. These effects include (1) change in maximum frequency of a critical path, (2) degradation of the clock network performance, etc. Thus, designers must consider the dynamic effect of supply voltage noise early in the design cycle.

Finally, the PDN is a major consumer of resources (e.g., more than 30% of the entire routing area) in wire-limited deep-submicron designs [57]. Conventionally, the PDN is designed to satisfy power integrity constraints, but without understanding the true implications of supply noise on delay, correct optimization of PDN is impossible. To close this gap, our present work gives a methodology for closed-form modeling of the delay impact of supply voltage noise (characterized by noise slew, offset, and magnitude). We believe our models can efficiently drive accurate worst-case performance-driven PDN optimization, as shown in Figure 6.6.

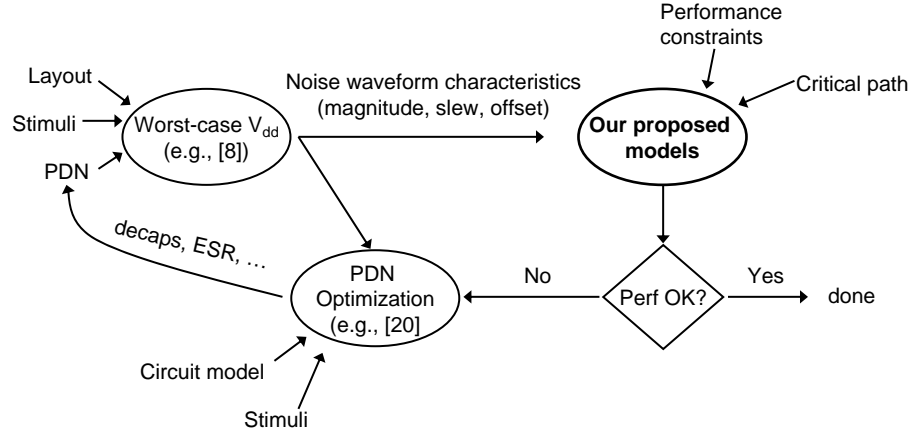


Figure 6.6: Accurate worst-case performance-driven power distribution network optimization flow.

In this section of chapter, we propose a new modeling paradigm in which we use *machine learning-based nonparametric* regression techniques to develop accurate early-stage performance models under dynamic supply voltage and temperature variations. The contributions of our work are as follows.

- We propose a framework for gate delay modeling under supply voltage and temperature variations, using machine learning-based nonparametric regression methods. Further, we introduce a reproducible flow to aid *automatic* generation of accurate performance estimation models (e.g., using generic critical paths).
- We develop early-stage performance models using our basic gate delay models, to enable worst-case performance prediction that can efficiently drive PDN optimization.
- We validate our models against SPICE simulations using 65nm foundry SPICE models.

6.2.2 Related Work

Gate delay models under supply voltage noise can be classified as (1) static or (2) dynamic; with the former type, the dynamic behavior of the noise waveform

is ignored. The majority of the existing literature focuses on the former type [8, 18, 31, 43, 56]. Hashimoto *et al.* [18] propose to replace supply voltage noise with an equivalent power/ground voltage. However, this method assigns static voltage value (time-invariant) during the static timing analysis (STA), and cannot appropriately capture the dynamic behavior of the noise waveform. Martorell *et al.* [31] present a probabilistic approach to estimate supply voltage noise bound given performance criteria. However, they assume that all gates in a combinational path have the same supply voltage value; this assumption is incorrect due to the presence of dynamic supply voltage noise.

Chen *et al.* [8] propose closed-form equations to estimate the change in delay of buffers in the presence of supply voltage noise. However, the authors do not consider specific noise waveform characteristics (magnitude, offset, and slew) in their analysis. In another effort, Weng *et al.* [56] propose a methodology to improve the accuracy of gate delay calculation under supply voltage noise by taking into account (time-varying) IR drop waveforms. To capture the dynamic impact of supply voltage noise, the authors of [56] partitioned noise waveform to discrete sections and assign an equivalent DC value across different time intervals. The DC values are calculated as the average supply voltage values over the entire interval. This method still does not capture the ‘true’ dynamic behavior of the supply noise waveform. To assess the impact of supply voltage noise on circuit performance, [43] suggests that using average supply voltage, rather than dynamic behavior, can be well-correlated with measurements; however, the authors fail to demonstrate the limitations of timing analysis using static IR-drop analysis as noted in [33].

Recently, Okumura *et al.* [33] have proposed a gate delay calculation approach which considers the dynamic behavior of the supply voltage noise by considering noise waveform slew and magnitude. However, in their characterization setup they do not allow all the relevant parameters (i.e., input slew, noise slew, noise magnitude, etc.) to change simultaneously; this limits the applicability of their proposed model. In our present work, we develop new gate delay and output slew models under supply voltage and temperature variations, where all the

relevant parameters can interact with one another.

6.2.3 Implementation Flow

Figure 6.7 shows our implementation flow, which begins with SPICE simulations using foundry SPICE models and extracted or CDL SPICE netlists for each gate type. We measure the 50% delay and output slew of each gate with respect to a number of different parameters. In our experiments we have three main axes: (1) cell delay parameters, (2) supply voltage noise parameters, and (3) temperature. These parameters, and the values that they take on in our experiments, are explained below. Cell delay parameters include (1) input slew $slew_{in}$, (2) output load $load_{out}$, and (3) cell size $cell_{size}$. For supply voltage we use 0.9V as the nominal value, with noise waveform superimposed on it.

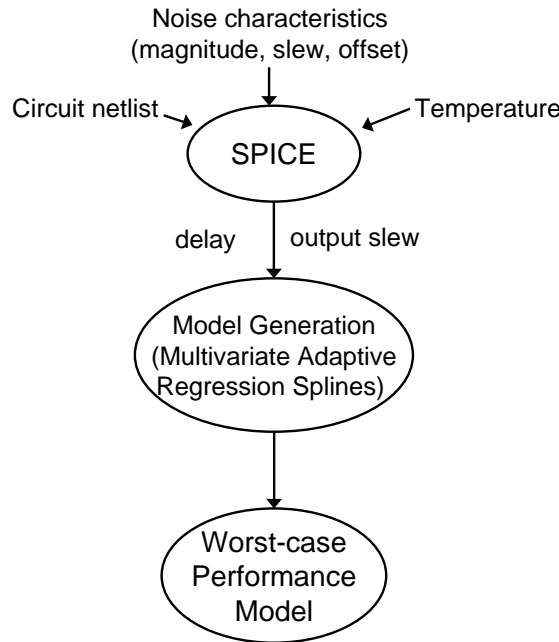


Figure 6.7: Implementation flow.

Supply voltage noise parameters include (1) noise amplitude amp_{noise} , (2) noise slew $slew_{noise}$, and (3) noise offset $offset_{noise}$. Noise offset denotes the noise transition time with respect to that of the input signal transition. Finally, temperature denotes the operating temperature of the transistors. In our studies, we use

two different cells (1) inverter, and (2) 2-input NAND to show the applicability of our modeling approach. For worst-case performance model we implement our basic cell delay and output slew models in C++. Using our basic delay and output slew models we construct path delay models with arbitrary number of stages and a mix of different cells. We run a total of 30720 SPICE simulations and gather delay and output slew values corresponding to different parameters (cf. Table 6.3).

We use *Synopsys HSPICE v.Y-2006.03* [20] for SPICE simulations using 65nm foundry SPICE models and netlists. We perform our experiment using typical corner and normal- V_{th} (NVT) transistors. We also use *MARS3.0* [30] to implement nonparametric regression techniques.

Table 6.3: List of parameters used for performance modeling.

Parameter	Values
$slew_{in}$	{0.00056, 0.00112, 0.0392, 0.1728, 0.56, 0.7088}ns
$load_{out}$	{0.0009, 0.0049, 0.0208, 0.0842}pF
$cell_{size}$	INV: {1, 4, 8, 20} ND2D: {1, 2, 4, 8}
amp_{noise}	{0, 0.054, 0.144, 0.27}V
$slew_{noise}$	{0.01, 0.04, 0.07, 0.09}ns
$offset_{noise}$	{-0.15, -0.05, 0, 0.05, 0.15}ns
$temp$	{-40, 25, 80, 125}°C

6.2.4 Modeling Methodology

Previous delay estimation techniques do not consider dynamic impact of supply voltage noise on cell delay [18, 31, 43]. By contrast, we propose to pursue a different modeling paradigm in which we use *machine learning-based nonparametric regression techniques* to capture the dynamic impact of supply voltage noise on cell delay. To illustrate the basic idea, consider the following baseline model generation flow:

- We begin with a parameterized SPICE netlist for a given inverter cell. We refer to this as a *configurable* inverter SPICE specification, which will be used to generate the representative inverter cell delay under different cell

and supply voltage noise parameters. For example, a given SPICE simulation setup can be configured with respect to (1) input slew, (2) output load, (3) inverter size, (4) supply voltage noise magnitude, (5) supply voltage noise width (i.e., frequency), (6) voltage noise offset (i.e., with respect to the input transition), and (7) temperature.

- Using a small subset of selected configurations for *training*, we run through each configuration in this training set through SPICE simulations, to obtain an accurate cell delay for each instance.
- Finally, we apply machine learning-based nonparametric techniques on the training set of delay to derive the corresponding cell delay estimation models.

In general, the modeling problem aims to approximate a function of several to many variables using only the dependent variable space. This generic formulation has applications in many disciplines. The goal is to model the dependence of a target variable y on several predictor variables x_1, \dots, x_n given R realizations $\{y_i, x_{1i}, \dots, x_{ni}\}_1^R$. The system that generates the data is presumed to be described by

$$y = f(x_1, \dots, x_n) + \epsilon \quad (6.1)$$

over some domain $(x_1, \dots, x_n) \in \mathcal{D} \subset \mathcal{R}^n$ containing the data [14]. Function f captures the joint predictive relationship of y on x_1, \dots, x_n , and the additive stochastic noise component ϵ usually reflects the dependence of y on quantities other than x_1, \dots, x_n that are neither controlled nor observed. Hence, the aim of the regression analysis is to construct a function $\hat{f}(x_1, \dots, x_n)$ that can accurately approximate $f(x_1, \dots, x_n)$ over the domain \mathcal{D} of interest. There are two main regression analysis methods: (1) global parametric, and (2) nonparametric. The former approach has limited flexibility, and can produce accurate approximations only if the assumed underlying function \hat{f} is close to f . In the latter approach, \hat{f} does not take a predetermined form, but is constructed according to information derived from the data. Multivariate adaptive regression splines (MARS) is a nonparametric regression technique which is an extension of linear models that

automatically models nonlinearities and interactions, and is used in our methodology. In this section, we use MARS-based approach to model the dynamic impact of supply voltage noise on cell delay.

6.2.5 Multivariate Adaptive Regression Splines

Given different cell and supply voltage noise parameters \mathcal{X} , we apply MARS to construct cell delay model, $d_{cell} = \hat{f}(x_1, \dots, x_n)$. Variables x_1, \dots, x_n denote cell and supply voltage noise parameters. The general MARS model can be represented as [63]

$$\hat{y} = c_0 + \sum_{i=1}^I c_i \prod_{j=1}^J b_{ij}(x_{ij}) \quad (6.2)$$

where \hat{y} is the target variable (i.e., inverter delay and output slew in our problem), c_0 is a constant, c_i are fitting coefficients, and $b_{ij}(x_{ij})$ is the truncated power basis function¹ with x_{ij} being the microarchitectural parameter used in the i^{th} term of the j^{th} product. I is the number of basis functions and J limits the order of interactions. In our experiments we set the number of basis functions to 100 and the order of interactions to 6, i.e., every parameter can interact with all the other parameters. The basis functions $b_{ij}(x_{ij})$ are defined as

$$\begin{aligned} b_{ij}^-(x^{\mu arch} - t_{ij}) &= [-(x^{\mu arch} - t_{ij})]_+^q & (6.3) \\ &= \begin{cases} (t_{ij} - x^{\mu arch})^q & x^{\mu arch} < t_{ij} \\ 0 & \text{otherwise} \end{cases} \end{aligned}$$

$$\begin{aligned} b_{ij}^+(x^{\mu arch} - t_{ij}) &= [(x^{\mu arch} - t_{ij})]_+^q & (6.4) \\ &= \begin{cases} (x^{\mu arch} - t_{ij})^q & x > t_{ij} \\ 0 & \text{otherwise} \end{cases} \end{aligned}$$

where q (≥ 0) is the power to which the splines are raised to adjust the degree of \hat{y} smoothness, and t_{ij} is called a knot. When $q = 0$ simple linear splines are applied.

¹ Each basis function can be a constant, a hinge function that is of form $\max(0, c - x)$ or $\max(0, x - c)$, or a product of two or more hinge functions.

The optimal MARS model is built in two passes. (1) Forward pass: MARS starts with just an intercept, and then repeatedly adds basis function in pairs to the model. Total number of basis functions is an input to the modeling. Backward pass: during the forward pass MARS usually builds an overfit model; to build a model with better generalization ability, the backward pass prunes the model using a generalized cross-validation (GCV) scheme

$$GCV(K) = \frac{1}{n} \frac{\sum_{k=1}^n (y_k - \hat{y})^2}{[1 - \frac{C(M)}{n}]^2} \quad (6.5)$$

where n is the number of observations in the data set, K is the number of non-constant terms, and $C(M)$ is a complexity penalty function to avoid overfitting.

6.2.6 Accurate Cell Delay Modeling

In this section, we discuss the impact of supply voltage noise and temperature variation on cell delay, and note that delay modeling under supply voltage and temperature variation is a nontrivial task. We show an example of our proposed delay and output slew models derived from machine learning-based nonparametric regression techniques. We also propose a methodology to find the worst-case input configuration that maximizes the delay of a given path.

Cell Delay and Output Slew Models

In the existing literature [18, 31], supply voltage variation is assumed to be constant (time-invariant). When the supply voltage varies slowly with respect to the clock period, this is reasonable. This assumption enables to predict the timing impact of the supply voltage noise: the worst-case delay corresponds to the worst-case noise that can occur when the target cell is switching. In other words, when the supply voltage varies slowly, the delay degradation is proportional to the peak of the noise [43]. However, to better capture the impact of time-varying supply voltage noise we must consider the noise waveform characteristics including (1) noise amplitude, (2) noise slew, and (3) noise offset. Figure 6.8 shows the impact of noise slew on cell inverter delay. We observe that noise slew affects cell delay

only when it is comparable to input slew. Hence, we must take into consideration the specific noise waveform characteristics to ensure more accurate delay modeling.

Existing PDN optimization frameworks [61, 59] use *fluctuation area*, i.e., the area under the noise waveform, as the metric to represent the supply voltage noise. However, it is easy to see that such an approach can incur significant error in the delay estimation. Consider two scenarios: (1) $slew_{noise}=0.2\text{ns}$, $amp_{noise}=0.2\text{V}$ and (2) $slew_{noise}=0.4\text{ns}$, $amp_{noise}=0.1\text{V}$. Using a triangular waveform to represent the supply noise, the two scenarios have different noise waveforms, yet have similar areas under the noise curve. When we evaluate gate delay under each of these scenarios, we observe 22% difference. (In this evaluation, we use a single inverter, with other parameters values being $slew_{in}=0.4\text{ns}$, $load_{out}=0.002\text{pF}$, $cell_{size}=1\text{X}$, $offset_{noise}=0\text{ns}$, and $temp=25^\circ\text{C}$.) We conclude that to accurately model the impact of supply voltage noise on cell delay, we must consider both noise slew and noise magnitude parameters, and not simply the area under the noise waveform.

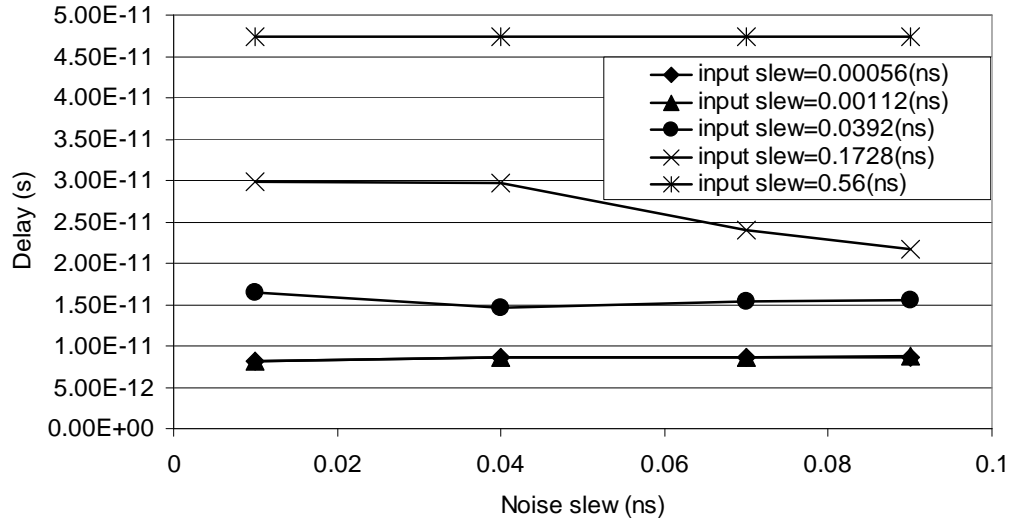


Figure 6.8: Delay of an inverter cell versus noise slew, for different input slew values.

The other important supply voltage noise characteristic is *noise offset*, which denotes the time of the voltage noise transition relative to the time of the input signal transition. We expect that as long as the supply voltage noise waveform is outside of the input signal transition window, it should not have any impact

on cell delay. However, when the noise waveform overlaps with the input signal transition, there will be an effect on cell delay. Figure 6.9 shows the impact of noise offset on cell delay. In our experiment, input slew and noise slew are $0.09ns$ and $0.1ns$, respectively. In our delay model, we explicitly consider noise offset as an input to the model.

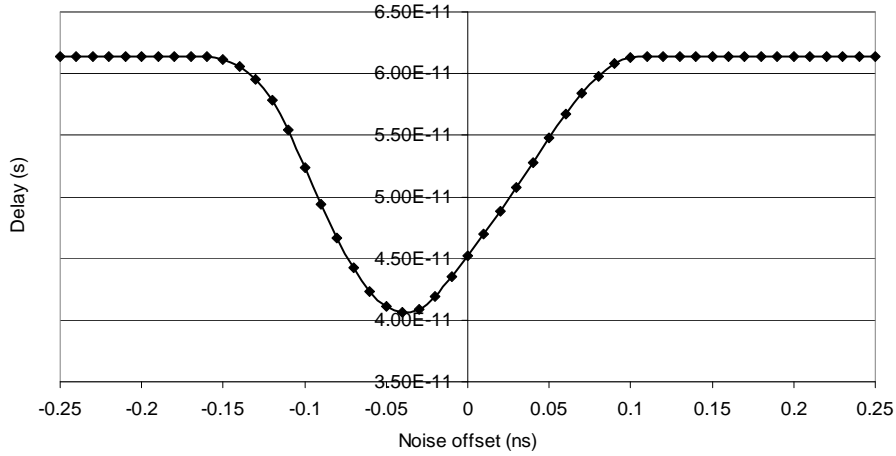


Figure 6.9: Impact of supply voltage noise offset on cell delay.

In addition, cell characteristics are influenced by temperature. Temperature impacts cell delay through voltage threshold, mobility, etc. parameters [17]. For example, as temperature decreases, both threshold voltage and mobility increase; the latter causes increased saturation current. However, the impact of temperature on cell delay depends on the gate voltage. The gate voltage at which the temperature shifts of threshold voltage and mobility exactly compensate each other's effects on delay is typically called zero-temperature-coefficient (ZTC) [29]. Hence, cell delay can increase or decrease with the increase in temperature. These complex relationships between cell delay and the aforementioned parameters make delay modeling a nontrivial task.

Finally, since our gate delay model depends on input slew, we must also model output slew of the previous stage of the critical path. Given the above discussion, we note that approximating CMOS gate delay is a nontrivial task with non-obvious implications, as seen from Figure 6.8. This has motivated us to explore machine learning-based nonparametric regression techniques to develop accurate

cell delay and output slew models. Figure 6.10 illustrates the form of resulting inverter delay and output slew models using 65nm foundry SPICE models.²

Delay Model
$b_1 = \max(0, load_{out} - 0.0208);$ $b_2 = \max(0, 0.0208 - load_{out}); \dots$ $b_{98} = \max(0, offset_{noise} - 0.05) \times b_{92};$ $b_{100} = \max(0, offset_{noise} + 2.4e-12) \times b_{37};$ $d_{cell} = 1.018e - 11 + 7.353e - 10 \times b_1 - 5.890e - 10 \times b_2$ $- 2.172e - 11 \times b_3 + \dots - 1.708e - 7 \times b_{96} +$ $2.431e - 7 \times b_{98} - 3.031e - 8 \times b_{100}$
Output Slew Model
$b_1 = \max(0, load_{out} - 0.0009);$ $b_2 = \max(0, cell_{size} - 4) \times b_1; \dots$ $b_{99} = \max(0, 0.05 - slew_{noise}) \times b_{55};$ $b_{100} = \max(0, offset_{noise} + 0.15) \times b_{94};$ $slew_{out} = 1.227e - 11 + 1.529 \times b_1 - 2.051e - 10 \times b_2$ $+ 2.050e - 9 \times b_3 + \dots - 1.081e - 8 \times b_{98}$ $- 4.327e - 9 \times b_{99} - 7.422e - 9 \times b_{100}$

Figure 6.10: Sample inverter delay and output slew models in 65nm.

6.2.7 Worst-case Performance Model

In this subsection we formalize the problem of finding the worst-case performance under dynamic supply voltage and temperature variations. We are interested in the specific configuration, i.e., set of seven parameters (7-tuple) described in Table 6.3 that causes the delay of a given path with arbitrary number of stages to be maximum.³ Note that we construct our path delay model using our basic cell delay and output slew models. Our proposed delay and output slew models are essentially mappings f and g , respectively, from the set of all 7-tuples Q (cf. Table 6.3) to the positive reals, i.e., $f : Q \rightarrow \mathcal{R}^+$ and $g : Q \rightarrow \mathcal{R}^+$, where $Q = slew_{in} \times load_{out} \times cell_{size} \times amp_{noise} \times slew_{noise} \times offset_{noise} \times temp$.

²Note that our methodology can be straightforwardly applied to future technologies, as long as necessary SPICE models and device-level netlists are available.

³In our experiments a path consists of (1) only inverter, (2) only 2-input NAND, and (3) a mix of inverter and 2-input NAND.

For a single stage the problem of finding the worst-case configuration seeks $\vec{q}^* \in Q$ such that $f(\vec{q}^*)$ is maximized. With more than one stage in a path, i.e., $k > 1$, the output slew of the previous stage becomes the input slew to current stage, and the noise offset must be adjusted accordingly. Then, we seek \vec{q}_1^* such that $f(\vec{q}_1^*) + \dots + f(\vec{q}_k^*)$ is maximized, where $\vec{q}_m^* = \vec{q}_1^*$ for all stages $1 < m < k$, except that the $slew_{in}$ component is replaced by $g(q_{m-1}^*)$ and the $offset_{noise}$ component is adjusted at the beginning of each stage. Note that the worst-case configuration is always going to be an element of the cross-product of the various sets of parameter values. In other words, it is one of $|slew_{in}| \times |load_{out}| \times |cell_{size}| \times |amp_{noise}| \times |slew_{noise}| \times |slew_{noise}| \times |offset_{noise}| \times |temp|$ configurations. In our studies, the worst-case configuration is out of 30720 different configurations.

6.2.8 Experimental Results and Validation

To generate our models, we randomly select 10% of our entire data set as training data; we then test the models on the other 90% of the data. To show that the selection of the training set does not substantially affect model accuracy, we randomly select 10% of the entire data set five times and show the corresponding models' maximum and average error values (Table 6.4).

Table 6.4: Model stability versus random selection of the training set.

Experiments	delay % diff		output slew % diff	
	max	avg	max	avg
Exp 1	56.993	5.660	55.117	6.012
Exp 2	53.342	5.458	56.896	5.976
Exp 3	53.661	5.401	56.237	5.526
Exp 4	55.419	5.552	54.883	5.311
Exp 5	55.015	5.609	55.614	5.672

To show the accuracy of our worst-case performance model, we compare our worst-case predictions with SPICE simulations. We construct three different paths with different number of stages, each consists of (1) only inverters, (2) only 2-input NAND, and (3) a mix of inverter and 2-input NAND gates. For (3), we construct the path starting with an inverter, and then alternating 2-input NAND

gates with inverter gates. In our experiments, one of the NAND gate inputs is connected to supply voltage (v_{dd}). We evaluate our predictions using two metrics: (1) correlation of our predictions against SPICE results, and (2) relative (%) difference in delays between our proposed model and SPICE. For (1) we rank our model predictions (total of 30720 data points) in descending order with respect to the delay of the given path. Each delay value corresponds to a set of parameters (i.e., 7-tuple including all the parameters shown in Table 6.3). Next, we compare our predicted worst-case configuration with SPICE, and find the rank ($rank_{SPICE}$) of our predicted worst-case configuration within SPICE results. For multi-stage paths with $k > 1$ stages, we need to adjust the noise offset for each stage. To perform this we need to identify the time at which the input to stage i , where $i = 1 \dots k$, makes the transition. This value can be estimated by calculating the delay up to stage $i - 1$, and subtracting $\frac{slew_{in}^i}{1.6}$ from it, where $slew_{in}^i$ is the input slew to stage i , and $\frac{slew_{in}^i}{1.6}$ determines the 50% output slew transition.⁴

Tables 6.5, 6.6, and 6.7 show the comparison our worst-case performance model with SPICE for a path consists of (1) only inverter, (2) only 2-input NAND, and (3) a mix of inverter and 2-input NAND gates, respectively. The second and third columns, represent our (2) and (1) comparison metrics, respectively. The fourth column shows where the SPICE worst-case configuration is ranked according to our proposed model ($rank_{MARS}$). We observe that our path delay models are within 4.3% of SPICE simulations. In addition, our predictions are always ranked in the top 3 (out of 30720 configurations) of the SPICE list ($rank_{SPICE}$). We note that the ability of our worst-case performance model to correctly predict worst-case configuration is beneficial for early-stage design and optimization of power distribution networks. Finally, the SPICE-computed worst-case performance value is always among top 5 predictions of our model.

⁴In our experiments, 10%-90% transition time is the slew value.

Table 6.5: Comparison of our proposed worst-case performance model and SPICE for an inverter chain. Rank values are out of 30720 configurations.

#Stage	delay % diff	$rank_{SPICE}$	$rank_{MARS}$
1	1.08	1	1
3	3.54	3	2
5	4.29	1	1
10	3.26	2	4
20	2.42	1	1
30	2.88	1	1

Table 6.6: Comparison of proposed worst-case performance model and SPICE for a 2-input NAND chain. Rank values are out of 30720 configurations.

#Stage	delay % diff	$rank_{SPICE}$	$rank_{MARS}$
1	1.34	1	1
3	3.21	1	1
5	3.69	2	3
10	3.11	1	1
20	3.43	2	3
30	2.37	2	2

Table 6.7: Comparison of proposed worst-case performance model and SPICE for a mixed inverter-NAND chain. Rank values are out of 30720 configurations.

#Stage	delay % diff	$rank_{SPICE}$	$rank_{MARS}$
1	1.08	1	1
3	2.73	2	4
5	3.24	3	5
10	3.36	1	1
20	3.93	2	4
30	2.85	1	1

6.3 Summary

In this chapter, we first outlined an improved frequency domain simulation method to analyze and assist silicon, package, and board PDN design. We then described a method to measure impedance of PDN. Finally, measurement re-

sults showed good correlation between processor performance sensitivity to PDN impedance. We quantified the impact of the middle and low frequency power integrity on the processor performance and F_{max} .

For the second part of this chapter, we have developed a methodology, based on nonparametric regression, to obtain accurate closed-form cell delay and output slew models under dynamic supply voltage and temperature variations. Our proposed models are within 6%, on average, of SPICE simulations. We show that our basic gate delay and output slew models can be used to construct delay estimates under supply noise for arbitrary critical paths. We also show that our models can accurately find the worst-case supply noise configuration that leads to worst-case delay performance. We believe that our proposed models can be beneficial in an accurate worst-case performance-driven power distribution network optimization, such as that shown in Figure 6.6.

6.4 Acknowledgments

Chapter 6 is based on the following publications:

- A. Shayan, C. Pan, M. Popovich, K. Bowles, “Estimation of Power Integrity Impact to Low Power Processor Performance through Pre-Silicon Simulation and Post-Silicon Measurements”, *AMSE InterPack*, 2011.
- Chung-Kuan Cheng, Andrew B. Kahng, Kambiz Samadi and Amirali Shayan, “Worst-case Performance Prediction under Supply Voltage and Temperature Variation”, *ACM/IEEE System-Level Interconnect Prediction (SLIP)*, 2010, pp. 91–96.

The dissertation author was the researcher and co-author of both papers. My co-authors (Dr. Christopher Pan, Dr. Mikhail Popovich, Kevin Bowles, Prof. Chung-Kuan Cheng, Prof. Andrew B. Kahng, and Dr. Kambiz Samadi) have all kindly approved the inclusion of the aforementioned publications in my thesis.

For the rest of the publications, my co-authors (Prof. Wenjian Yu, Prof. Ege Engin, Lew Chua-Eoan, Christopher Pan, Mikhail Popovich, Xiaoming Chen,

Xiaohua Kong, Wanping Zhang, Xiang Hu, He Peng, Sorin Dobre, Kevin Bowles, Thomas Toms, and Du Peng) have all kindly approved the inclusion of the aforementioned publications in my thesis.

Chapter 7

Conclusions

We entered new era for designing nanoscale low power and high performance processors where system level understanding of the whole complex power distribution is necessary. The success of an efficient power distribution will be translated into different aspects:

- A predictable performance is the main mission of a power distribution. The performance should be guaranteed under different processor and system on chip loading scenarios. Usually, designing a predictable goal could be achieved via cost effective PDN design.
- Maximum performance could be achieved if a fair enough cost-performance optimization is performed. With a solid power delivery system, architecture designers could push the processor frequency based on the PVT corners to the maximum upper bound.
- A robust power delivery sustains the performance of the system both in test and functional mode. Due to large current demand of the processors in test mode, performance might vary between functional mode and test mode. Proper design of PDN resource will achieve comparable results in test mode as well as the functional mode.

This dissertation, studied system level angels of designing a robust power distribution. The main contributions of each chapter are summarized as follows.

7.1 Thesis Summary

Chapter 2 provides an overview of different aspects of the whole power delivery system (PDS) including board, package and die design. An effective early stage modeling of power delivery system is outlined in this thesis that evolves through stages of design. A co-simulation time domain and frequency domain flow is proposed and developed to highlight the quality of the power distribution from different aspects.

Chapter 3 focus on study of die stacking in mobile chipsets as an alternative solution for scaling problem. A complete modeling of the 3D chip power distribution is detailed. The through silicon via and substrate coupling model are adopted in our Time-Frequency co-simulation flow. The failure mechanism of TSV and reliability aspects of the 3D stacking structure are discussed. The experimental results show that single TSV failure could increase voltage variation up to 70% in a local weak hotspot. A reliability-aware through silicon via and power delivery optimization flow is proposed to assist the complex design of the 3D stacked dies.

To provide a fair estimation of power integrity quality, PDS designers need to analyze the worst-case current load of network. We propose in Chapter 4 a vector-based resonance-aware current generation methodology to highlight the performance of the system under multiple impedance anti-resonance peaks. We propose an algorithm for generating synthetic vector based rogue wave current with complexity of $O(n^2 \cdot \log m)$. The proposed current waveform serves as a realist mean to assess the performance of the power delivery under worst loading.

Voltage regulation is one of the key pieces of the power distribution. Conventional PDS achieves this goal via bulk regulators and rely on the entire off-chip path to deliver current. Regulator path is most of the time ignored or is underestimated during analysis of the PDS. In Chapter 5, we highlight the problem of missing VRM in the PDS analysis and emphasize on the requirement of VRM inclusion in model. Furthers more, silicon measurement in Chapter 6 shows ignoring VRM in the loop will lead to unpredictable performance. An efficient parallel flow is introduced to assess the impact of the VRM noise of system.

On-die regulation is one of the main directions that addresses increase in

the complexity and power demand of the portable devices. On-die regulation is beneficial for the whole system because of the faster response and independency from unbalanced off-chip parasitic. On-die regulation enables optimal power mode configuration of different functional blocks which conventionally shared the rail. Each block tapping of the regulator could have its own optimum power configuration. In Chapter 5 an LDO-based power distribution design under worst loading is outlined and design tradeoff was discussed.

The final goal of power delivery design is that processor and functional blocks have a predictable performance. We provide a detailed analysis of the power integrity impact on low power processor performance in Chapter 6. The correlation of Time-Frequency domain analysis from the previous chapters with silicon measurement and processors performance were studied in Chapter 6. A model for predicting performance under worst voltage and temperature variation is proposed that achieves SPICE like accuracy.

Overall the framework of this thesis, tries to tackle different aspects of the power distribution design and provides means for designing a reliable mobile power delivery.

7.2 Future Research Directions

Looking to future, the complexity and number of applications and blocks integrated on mobile platforms are substantially increasing. On the other hand, as technology scales geometry and processor frequency are increasing. Hence, all these are hints toward a need for more advance power delivery techniques. We discuss in next section few possible research directions that author believes are viable or currently need more emphasis.

7.2.1 Architecture/Software/PDN Co-design (PDN-friendly Architectures)

Currently there is a considerable gap between architecture and power distribution design. More cost and area saving, performance increase and power saving

would happen if there is an in depth understanding of the architecture and PDN interactions. A strong motivation exists for developing software as well as the hardware architectures that perform in awareness of power noise. New system level policies for hardware and software are needed to adapt to dynamic variation of the voltage and temperature variation across the PDS system. Going forward this is one of the key directions for mobile platform power integrity.

7.2.2 Variation-aware Timing-Voltage Analysis Integration

We are reaching sub $16nm$ era where variation is even more pronounced. Operational voltage below $\sim 0.5V$ along with geometry scaling increase the variation of chip substantially. There is a strong push for understanding and modeling of temporal and spatial distribution of the voltage variation. Coupling this information with timing closure will lead to more reliable design. An integrated variation-aware *Timing-Power* closure is a direction for future research in mobile design.

7.2.3 Distributed On-die Regulations and Power Management

Distributed on-chip regulator circuits, policies and techniques are very appealing for the low power systems. Developing heterogeneous power distribution with localized regulation address many of the current low power requirements. Circuit techniques are improving to provide efficient means for on-die distributed regulations. There are substantial opportunities for research towards on-die regulations. Many interesting research topics are evolving under the umbrella of power distribution network. Many researchers will contribute to power integrity research in future from multiple fields such as theory and algorithm, software, architecture, circuit design, physical design, and electromagnetic.

Bibliography

- [1] <http://www.ansoft.com/products/hf/hfss/>.
- [2] http://www.ansoft.com/products/si/q3d_extractor.
- [3] <http://www.sigriety.com/products/powersi/powersi.htm>.
- [4] S. Alam, R. Jones, S. Rauf, and R. Chatterjee. Inter-strata connection characteristics and signal transmission in three-dimensional (3d) integration technology. In *8th International Symposium on Quality Electronic Design*, pages 580 –585, Mar 2007.
- [5] P. Arunasalam, F. Zhou, H. Ackler, and B. Sammakia. Thermo-mechanical analysis of thru-silicon-via based high density compliant interconnect. In *Electronic Components and Technology Conference*, pages 1179 –1185, Jun 2007.
- [6] S. I. Association. International technology roadmap for semiconductors, 2004, 2006, 2007.
- [7] M. M. Budnik and K. Roy. A Power Delivery and Decoupling Network Minimizing Ohmic Loss and Supply Voltage Variation in Silicon Nanoscale Technologies. *IEEE Transactions on Very Large Scale Integration (VLSI) Systems*, 14(12):1336–1346, Dec. 2006.
- [8] L. Chen, M. Marek-Sadowska, and F. Brewer. Buffer delay change in the presence of power and ground noise. *Very Large Scale Integration (VLSI) Systems, IEEE Transactions on*, 11:461 –473, Jun 2003.
- [9] D. Deshrijver and T. Dhaene. Broadband Macromodelling of Passive Components using Orthonormal Vector Fitting. *Electronic Letters*, 41:1160–1161, Oct. 2005.
- [10] V. Drabkin, C. Houghton, I. Kantorovich, and M. Tsuk. Aperiodic resonant excitation of microprocessor power distribution systems and the reverse pulse technique. In *Electrical Performance of Electronic Packaging*, pages 175 – 178, 2002.

- [11] P. Du, X. Hu, S.-H. Weng, A. Shayan, X. Chen, A. Ege Engin, and C.-K. Cheng. Worst-case noise prediction with non-zero current transition times for early power distribution system verification. In *11th International Symposium on Quality Electronic Design (ISQED)*, pages 624–631, Mar 2010.
- [12] W. Eisenstadt and Y. Eo. S-parameter-based ic interconnect transmission line characterization. *IEEE Transactions on Components, Hybrids, and Manufacturing Technology*, 15(4):483–490, Aug 1992.
- [13] B. Epstein and I. Weissman. *Mathematical Models for Systems Reliability*. 2008.
- [14] J. H. Friedman. *Multivariate Adaptive Regression Splines*. Annals of Statistics, 1991.
- [15] FWgrid project. <http://fwgrid.ucsd.edu>.
- [16] N. H. A. Ghani and F. N. Najm. Fast vectorless power grid verification using an approximate inverse technique. 2009.
- [17] M. Graziano, M. R. Casu, G. Masera, G. Piccinini, and M. Zamboni. Effects of temperature in deep-submicron global interconnect optimization in future technology nodes. *Microelectronics Journal*, (10):849–857, 2004.
- [18] M. Hashimoto, J. Yamaguchi, and H. Onodera. Timing analysis considering spatial power/ground level variation. In *IEEE/ACM International Conference on Computer Aided Design*, pages 814–820, 2004.
- [19] S. W. Ho, S. W. Yoon, Q. Zhou, K. Pasad, V. Kripesh, and J. Lau. High rf performance tsv silicon carrier for high frequency application. In *Electronic Components and Technology Conference*, pages 1946–1952, May 2008.
- [20] HSPICE. <http://www.synopsys.com/>.
- [21] X. Hu, W. Zhao, P. Du, Y. Zhang, A. Shayan, C. Pan, A. E. Egin, and C.-K. Cheng. On the bound of time-domain power supply noise based on frequency-domain target impedance. In *Proceedings of the 11th international workshop on System level interconnect prediction*, pages 69–76, 2009.
- [22] G. Huang, M. Bakir, A. Naeemi, H. Chen, and J. D. Meindl. Power delivery for 3d chip stacks: Physical modeling and design implication. In *IEEE Electrical Performance of Electronic Packaging*, pages 205–208, Oct. 2007.
- [23] P. K. K. Banerjee, S. J. Souri and K. C. Saraswat. 3D ICs: A Novel Chip Design for Improving Deep-submicrometer Interconnect Performance and Systems-on-chip Integration. *Proceedings of the IEEE*, 89(5):602–633, May 2001.

- [24] I. Kantorovich and C. Houghton. Maximum tolerable power supply noise for data-clock synchronization. In *IEEE Electrical Performance of Electronic Packaging*, pages 167 –170, Oct 2006.
- [25] R. Labie, W. Ruythooren, K. Baert, E. Beyne, and B. Swinnen. Resistance to electromigration of purely intermetallic micro-bump interconnections for 3d-device stacking. In *International IEEE Conference on Interconnect Technology Conference*, pages 19 –21, Jun 2008.
- [26] S. K. Lau, K. N. Leung, and P. Mok. Analysis of low-dropout regulator topologies for low-voltage regulation. In *IEEE Conference on Electron Devices and Solid-State Circuits*, pages 379 – 382, Dec 2003.
- [27] S. Lim. Physical design for 3d system on package. *Design Test of Computers, IEEE*, 22(6):532 – 539, Nov-Dec 2005.
- [28] S. Lim and A. Huang. Low-dropout (ldo) regulator output impedance analysis and transient performance enhancement circuit. In *Twenty-Fifth Annual IEEE Applied Power Electronics Conference and Exposition (APEC)*, pages 1875 –1878, 2010.
- [29] E. Long, W. Daasch, R. Madge, and B. Benware. Detection of temperature sensitive defects using ztc. In *VLSI Test Symposium, 2004. Proceedings. 22nd IEEE*, pages 185 – 190, Apr 2004.
- [30] MARS. Mars user guide. <http://www.salfordsystems.com/>.
- [31] F. Martorell, M. Pons, A. Rubio, and F. Moll. Error probability in synchronous digital circuits due to power supply noise. In *Design Technology of Integrated Systems in Nanoscale Era, 2007. DTIS. International Conference on*, pages 170 –175, Sep 2007.
- [32] R. Milliken, J. Silva-Martinez, and E. Sanchez-Sinencio. Full on-chip cmos low-dropout voltage regulator. *IEEE Transactions on Circuits and Systems I: Regular Papers*, 54(9):1879 –1890, Sep 2007.
- [33] T. Okumura, F. Minami, K. Shimazaki, K. Kuwada, and M. Hashimoto. Gate delay estimation in sta under dynamic power supply noise. In *Proceedings of the Asia and South Pacific Design Automation Conference*, pages 775 –780, 2010.
- [34] M. Popovich, A. V. Mezhiba, and E. G. Friedman. *Power Distribution Networks with On-Chip Decoupling Capacitors*. Springer, 2008.
- [35] H. Qian, S. R. Nassif, and S. S. Sapatnekar. Power Grid Analysis Using Random Walks. *IEEE Transaction on Computer-Aided Design of Integrated Circuits and Systems*, 24(8), Aug. 2005.

- [36] Y. Qui. High-Frequency Modeling and Analysis Buck and Multiphase Buck Convertors. *Virginia Tech Ph.D. Dissertation*, 2001.
- [37] Y. S. M. W. R. Yang, C.-Y. Hung and H.-W. Wu. Loss characteristics of silicon substrate with different resistivity. In *Microwave and Optical Technology Letters*, pages 1773–1776, 2006.
- [38] T. Rahal-Arabi, G. Taylor, M. Ma, and C. Webb. Design and validation of the pentium reg; iii and pentium reg; 4 processors power delivery. In *Symposium on VLSI Circuits Digest of Technical Papers*, pages 220 – 223, 2002.
- [39] A. Rahman, J. Trezza, B. New, and S. Trimberger. Die stacking technology for terabit chip-to-chip communications. In *Custom Integrated Circuits Conference*, pages 587 –590, Sep 2006.
- [40] READHAWK. Power noise analysis for next generation ics. *Apache Design Solutions*.
- [41] T. Rosing, K. Mihic, and G. De Micheli. Power and reliability management of socs. *IEEE Transactions on Very Large Scale Integration (VLSI) Systems*, 15(4):391 –403, Apr 2007.
- [42] C. Ryu, D. Chung, J. Lee, K. Lee, T. Oh, and J. Kim. High Frequency Electrical Circuit Model of Chip-to-chip Vertical Via Interconnection for 3D Chip Stacking Package. In *IEEE 14th Topical Meeting on Electrical Performance of Electronic Packaging*, pages 151–154, Oct. 2005.
- [43] M. Saint-Laurent and M. Swaminathan. Impact of power-supply noise on timing in high-frequency microprocessors. *Advanced Packaging, IEEE Transactions on*, 27(1):135 –144, Feb 2004.
- [44] S. S. Sapatnekar and H. Su. Analysis and Optimization of Power Grids. *IEEE Design & Test of Computers*, 20:7–15, May 2003.
- [45] C. Selvanayagam, J.H. Lau, X. Zhang, S. K. W. Seah, V. Kripesh, and T. C. Chai. Nonlinear Thermal Stress/Strain Analyses of Copper Filled TSV (Through Silicon Via) and Their Flip-Chip Microbumps. In *Electronics Components and Technology Conference*, pages 1073–1081, May 2008.
- [46] A. Shayan, K. Bowles, S. Dobre, M. Popovich, X. Chen, and C. Pan. Resonance-aware methodology for system level power distribution network co-design. In *IEEE 18th Conference on Electrical Performance of Electronic Packaging and Systems*, pages 29 –32, 2009.
- [47] A. Shayan, X. Hu, H. Peng, M. Popovich, W. Zhang, C. Cheng, L. Chua-Eoan, and X. Chen. 3d power distribution network co-design for nanoscale stacked silicon ics. *Electrical Performance of Electronic Packaging*, 2008.

- [48] A. Shayan, X. Hu, H. Peng, W. Yu, W. Zhang, C.-K. Cheng, M. Popovich, X. Chen, L. Chua-Eaon, and X. Kong. Parallel flow to analyze the impact of the voltage regulator model in nanoscale power distribution network. In *Quality of Electronic Design*, pages 576–581, Mar 2009.
- [49] J. Shi, Y. Cai, S. X. Tan, and X. Hong. High Accurate Pattern Based Precondition Method for Extremely Large Power/Ground Grid Analysis. *International Symposium on Physical Design*, pages 108–113, 2006.
- [50] L. D. Smith, R. E. Anderson, D. W. Forehand, T. J. Pelc, and T. Roy. Power Distribution System Design Methodology and Capacitor Selection for Modern CMOS Technology. *IEEE Transactions on Advanced Packaging*, 22(3):284–291, Aug. 1999.
- [51] J. Sun, J.-Q. Lu, D. Giuliano, T. P. Chow, and R. J. Gutmann. 3D Power Delivery for Microprocessors and High-Performance ASICs. In *IEEE 22nd Applied Power Electronics Conference*, pages 127–133, Mar. 2007.
- [52] M. van Heijningen, M. Badaroglu, S. Donnay, G. Gielen, and H. De Man. Substrate noise generation in complex digital systems: Efficient modeling and simulation methodology and experimental verification. *IEEE Journal of Solid-State Circuits*, 37(8):1065 – 1072, Aug 2002.
- [53] A. Waizman. Cpu power supply impedance profile measurement using fft and clock gating. In *Electrical Performance of Electronic Packaging*, pages 29 – 32, Oct 2003.
- [54] S. Wane and A.-Y. Kuo. Chip-package co-design methodology for global co-simulation of redistribution layers (rdl). In *Electrical Performance of Electronic Packaging*, pages 59 –62, Oct 2008.
- [55] L. M. Wedepohl and L. Jackson. Modified Nodal Analysis: an Essential Addition to Electrical Circuit Theory and Analysis. *Engineering Science and Education Journal*, (3):84–92, June 2002.
- [56] S.-H. Weng, Y.-M. Kuo, S.-C. Chang, and M. Marek-Sadowska. Timing analysis considering ir drop waveforms in power gating designs. In *IEEE International Conference on Computer Design*, pages 532 –537, Oct 2008.
- [57] J. Xiong and L. He. Full-chip multilevel routing for power and signal integrity. *Integr. VLSI J.*, 40:226–234, Apr 2007.
- [58] K. Yao, M. Xu, Y. Meng, and F. C. Lee. Design Considerations for VRM Transient Response Based on the Output Impedance. *IEEE Transaction on Power Electronics*, 18(6), Nov. 2003.

- [59] W. Zhang, L. Zhang, A. Shayan, W. Yu, X. Hu, Z. Zhu, E. Engin, and C.-K. Cheng. On-chip power network optimization with decoupling capacitors and controlled-esrs. In *15th Asia and South Pacific Design Automation Conference*, Jan 2010.
- [60] W. Zhang, L. Zhang, R. Shi, H. Peng, Z. Zhu, L. Chua-Eoan, R. Murgai, T. Shibuya, N. Ito, and C. K. Cheng. Fast Power Network Analysis with Multiple Clock Domains. In *25th International Conference on Computer Design*, pages 456–463, Oct. 2007.
- [61] W. Zhang, Y. Zhu, W. Yu, A. Shayan, R. Wang, Z. Zhu, and C.-K. Cheng. Noise minimization during power-up stage for a multi-domain power network. In *Proceedings of the Asia and South Pacific Design Automation Conference*, 2009.
- [62] M. Zhao, R. V. Panda, S. S. Sapatnekar, and D. Blaauw. Hierarchical Analysis of Power Distribution Networks. *IEEE Transaction on Computer-Aided Design of Integrated Circuits and Systems*, 21(2):159–168, Feb. 2002.
- [63] Y. Zhou and H. Leung. Predicting object-oriented software maintainability using multivariate adaptive regression splines. *J. Syst. Softw.*, 80:1349–1361, Aug 2007.
- [64] C. Zhuo, J. Hu, M. Zhao, and K. Chen. Power Grid Analysis and Optimization Using Algebraic Multigrid. *IEEE Transaction on Computer-Aided Design of Integrated Circuits and Systems*, 27(4):738–751, Apr. 2008.