

ARITH 18

Montpellier, France  
June 25-27, 2007



18th IEEE Symposium on Computer Arithmetic



# The Return of Silicon Efficiency

Simon Knowles  
Icera



# AGENDA

---

- Chip Market
- Design Objectives
- Modern Transistors
- Optimality of Designs
- The Design Hull
- Implications for Design

# “The Cheap will Outsell the Good”

---

## Customer Priorities

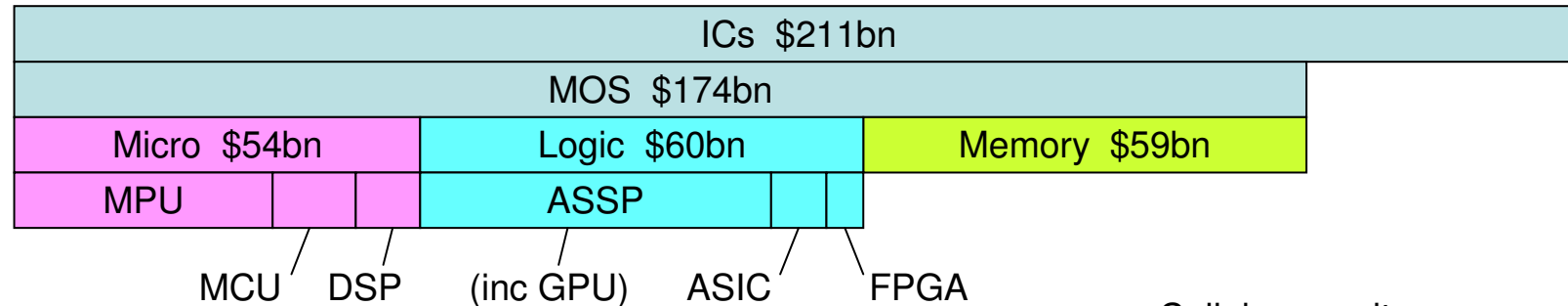
1. Availability
2. Price
3. { Features  
Power  
Size

## Design Choices

1. Hardware ↔ Software
2. Process Technology
3. MHz and Volts
4. Custom ↔ Synthesized

# The Chip Market

SIA 2006

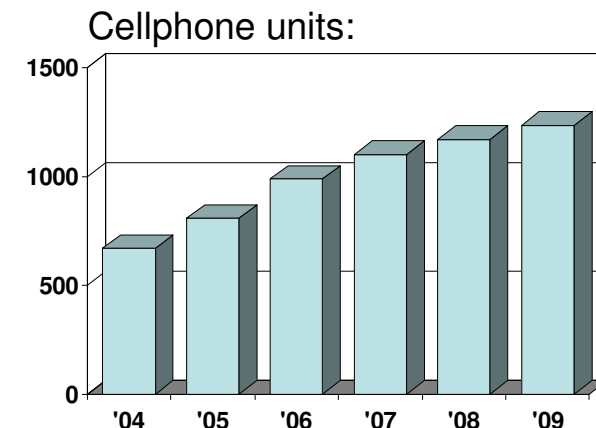


Most logic chips are processors or processor-centric.

Computers have driven MOS chip evolution – the first “product superclass”.

The second product superclass will be “cellphones”...

- Consumer driven – more cost sensitive than computers.
- High mobility – more power sensitive than computers.
- High (specialized) performance – eg. HSDPA cellular modem ~50GOPs.



# What's New(-ish) in Logic Chip Physics?

---

Small dopant populations ( $V_T$  variability)

Subthreshold leakage ( $V_T$  reduction is constrained, so  $V_{dd}$  likewise)

Gate leakage (end of  $\text{SiO}_2$   $T_{ox}$  scaling)

Channel Stress (engineered and STI-induced)

OPC rules

Well Proximity Effect

Line Edge Roughness (gate and wire)

NBTI (RAM  $V_T$  drift)

All important for designers, but not (directly) the  
main influencers of design evolution



# Top 2 Design Influencers

---

## 1. Most power-sensitive designs will set their own supply voltage.

- Design at typical, turn manufacturing speed spread (functionality risk) into power efficiency spread (battery life).
- Offset manufacturing leakage spread against speed spread.
- Dynamically trade performance vs power.
- Disconnect power in standby (efficient regulators are switch-mode).

## 2. Processors are evolving to displace fixed-function hardware

- Target multiple market sockets per design.
- Hardware-efficient (dynamic re-purposing of resources).
- Mitigation of specification uncertainty.
- Work around bugs by software change.
- Well-practiced idiom for speed.



# Modern Transistors



# Modern Transistors

---

ITRS distinguishes 3 roadmaps...

- High Performance (HP)
- Low Operating Power (LOP)
- Low Standby Power (LSTP)

TSMC offers 12 logic T's at 65nm...

- 2 x HP
- 6 x LOP
- 4 x LSTP

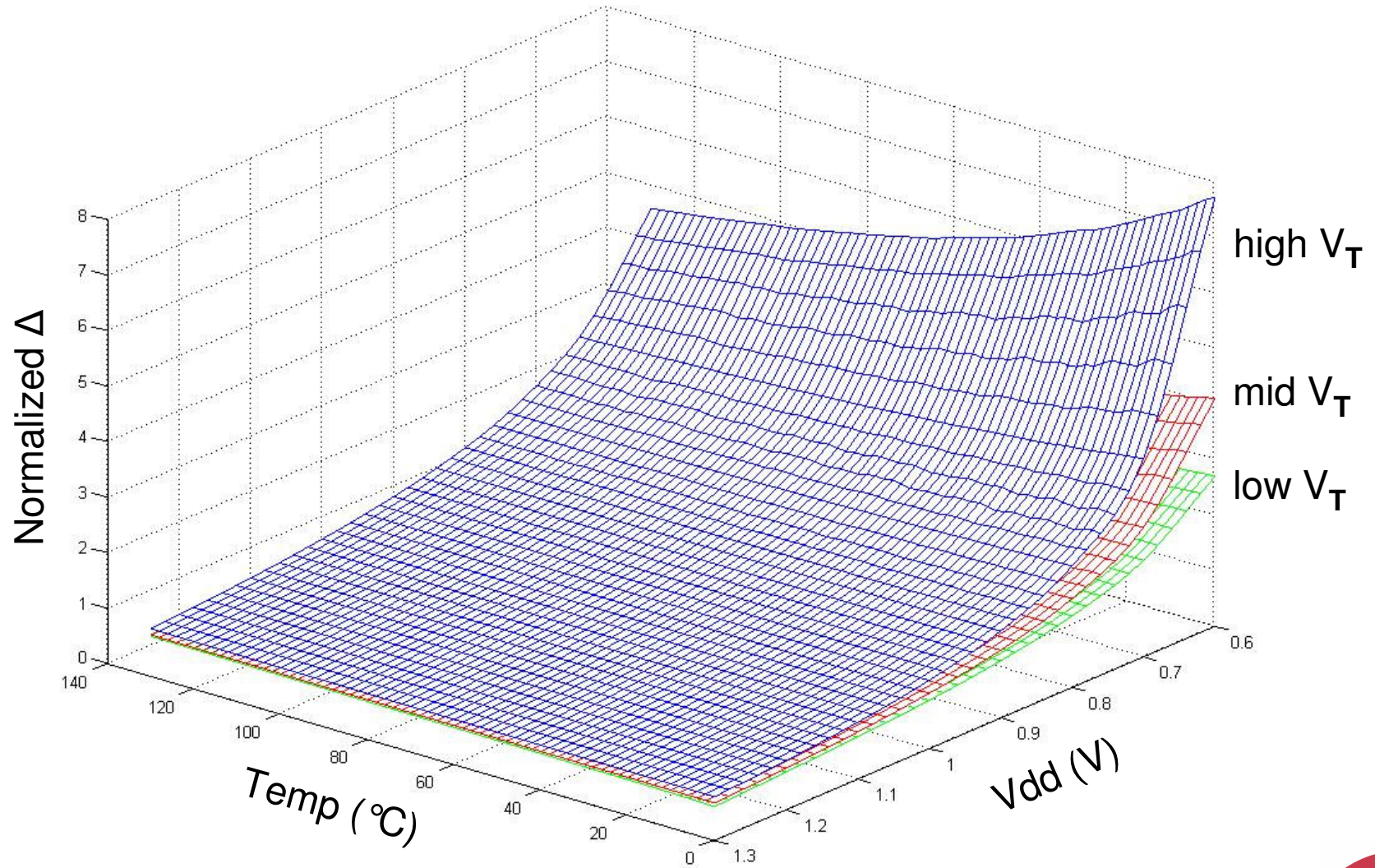
ITRS2005 for 65nm 2007

	$L_g$ (nm)	$T_{ox}$ (nm)	$V_T$ (mV)	
computers {	HP	25	1.1	165
"cellphones" {	LOP	32	1.2	285
	LSTP	45	1.9	524

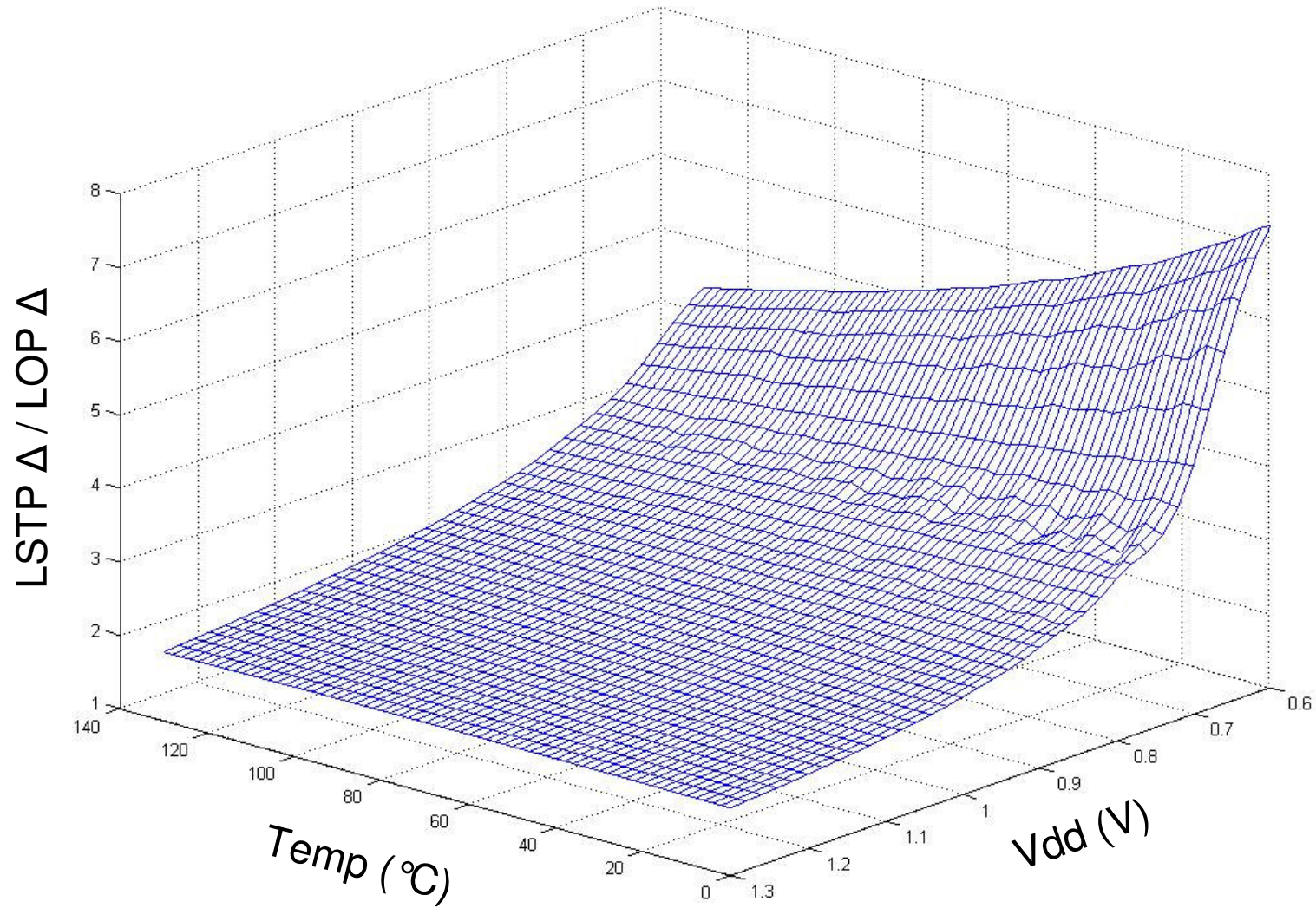




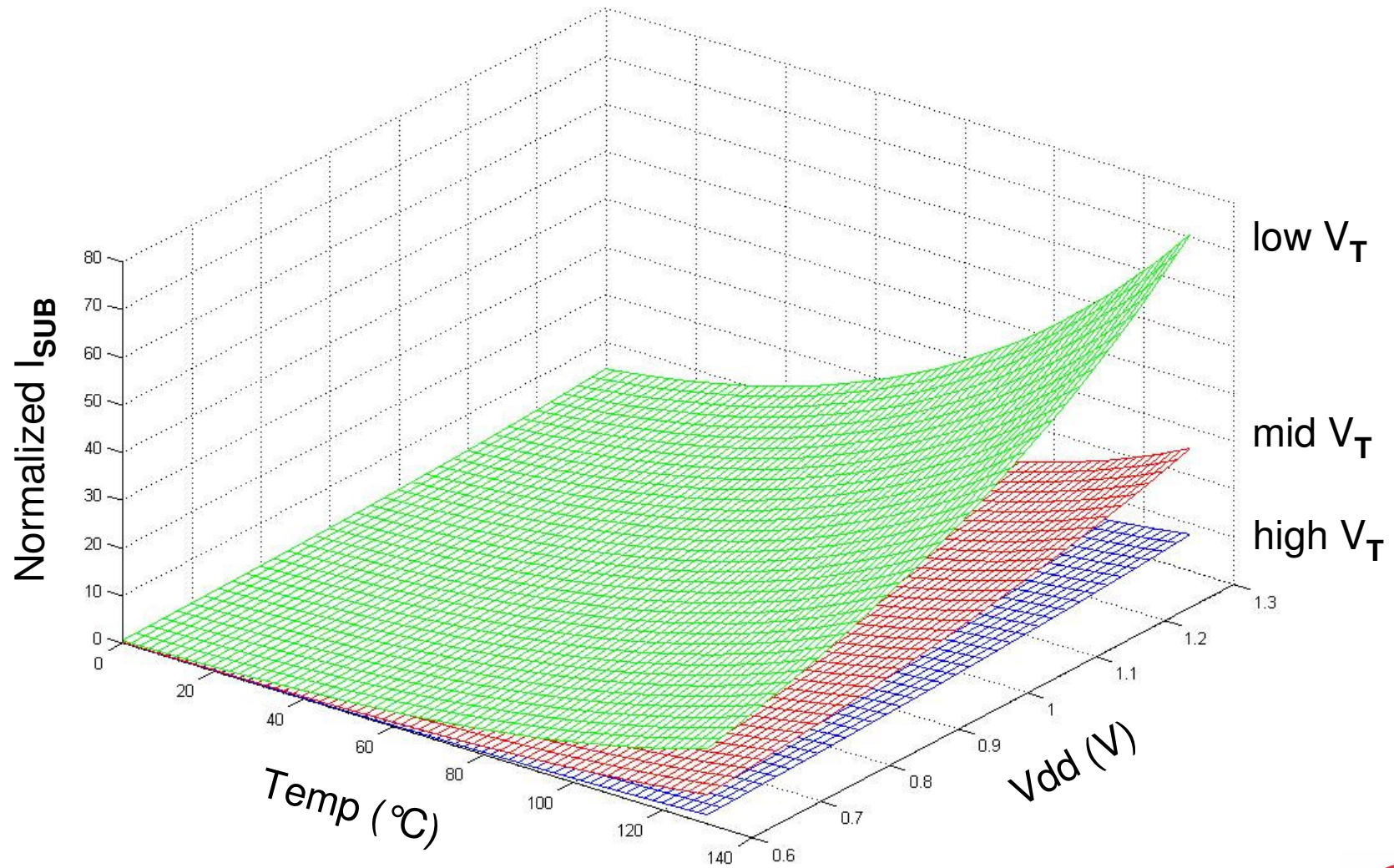
# 65nm LOP Delay



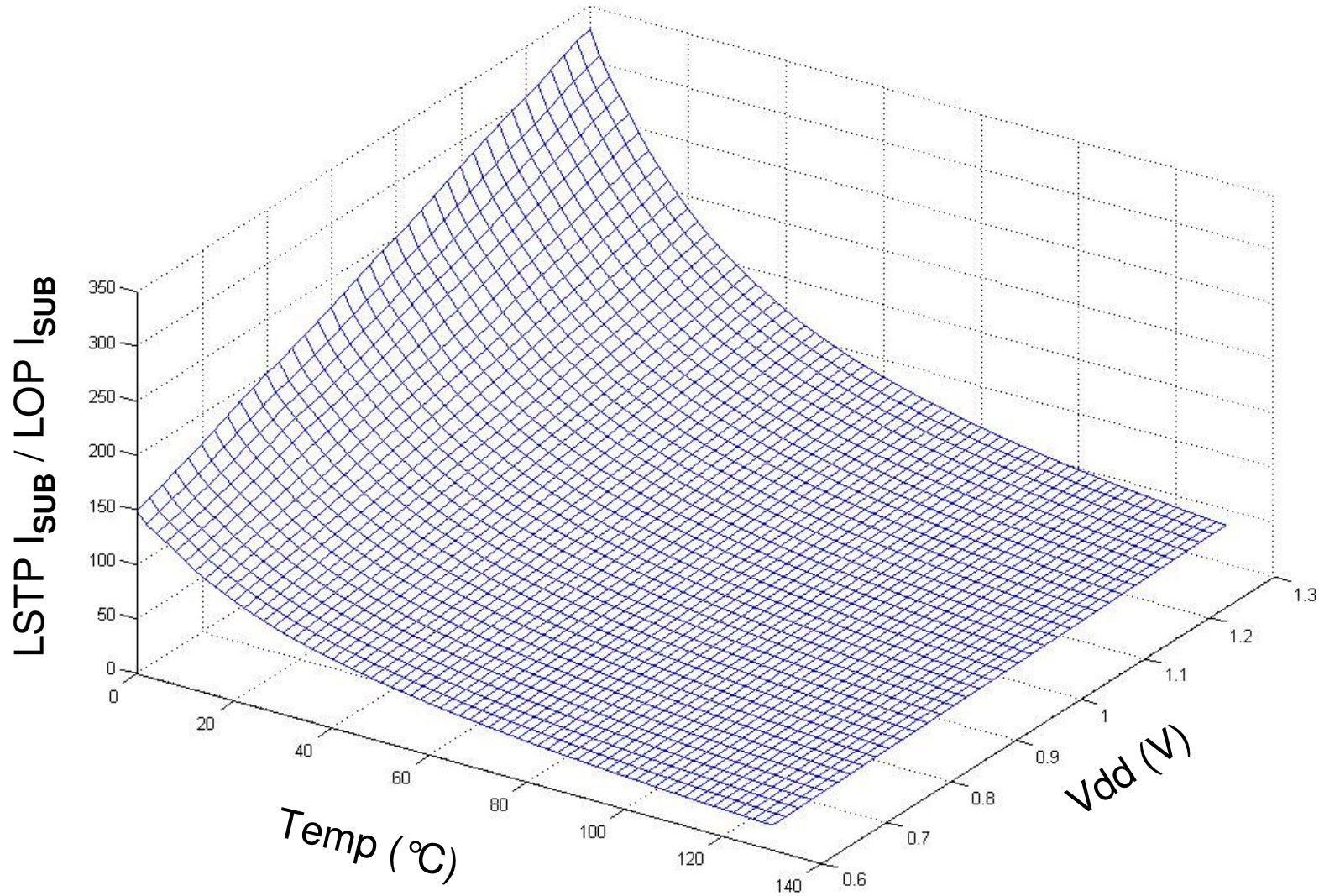
# 65nm LSTP vs LOP Delay



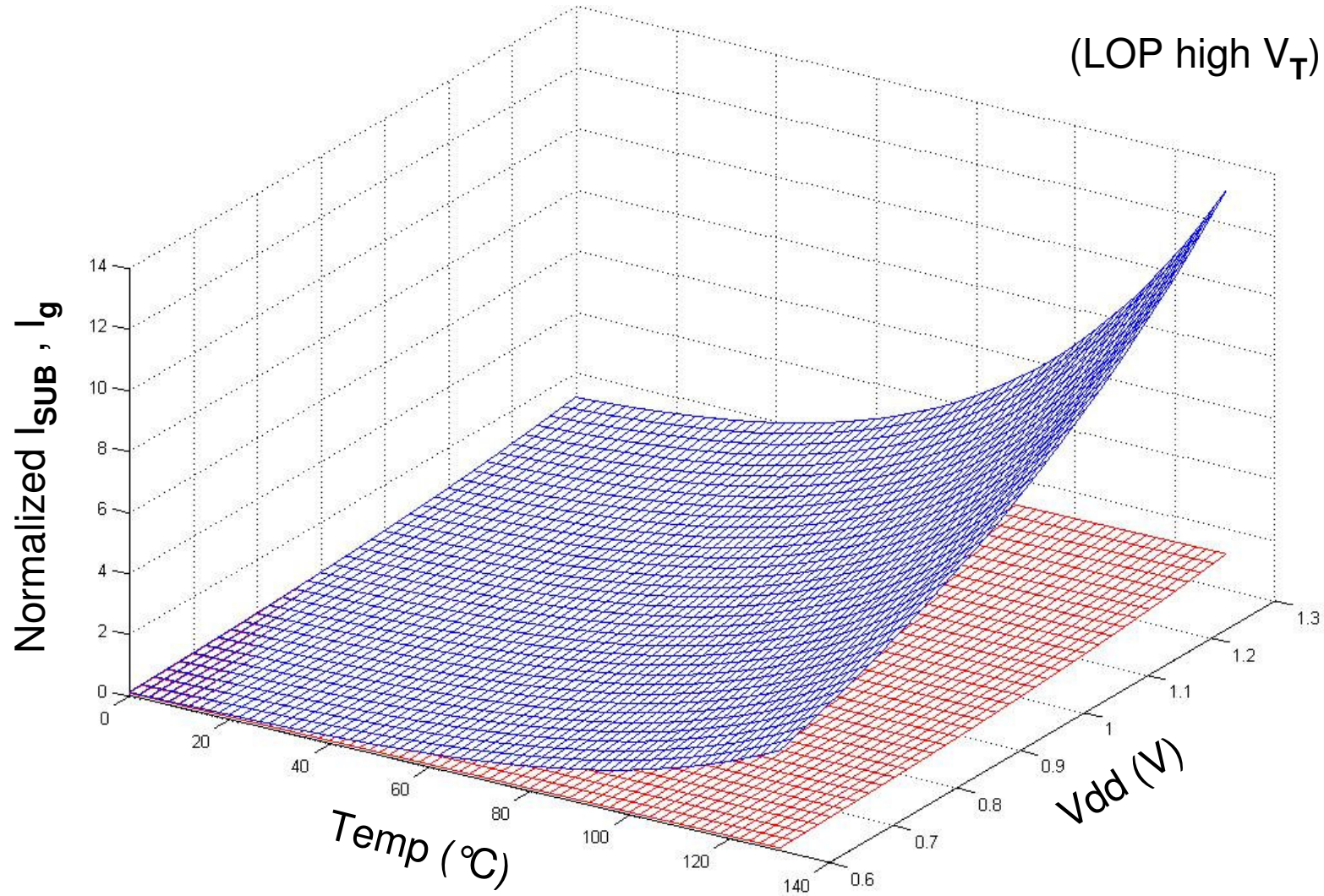
# 65nm LOP Sub-Threshold Leakage



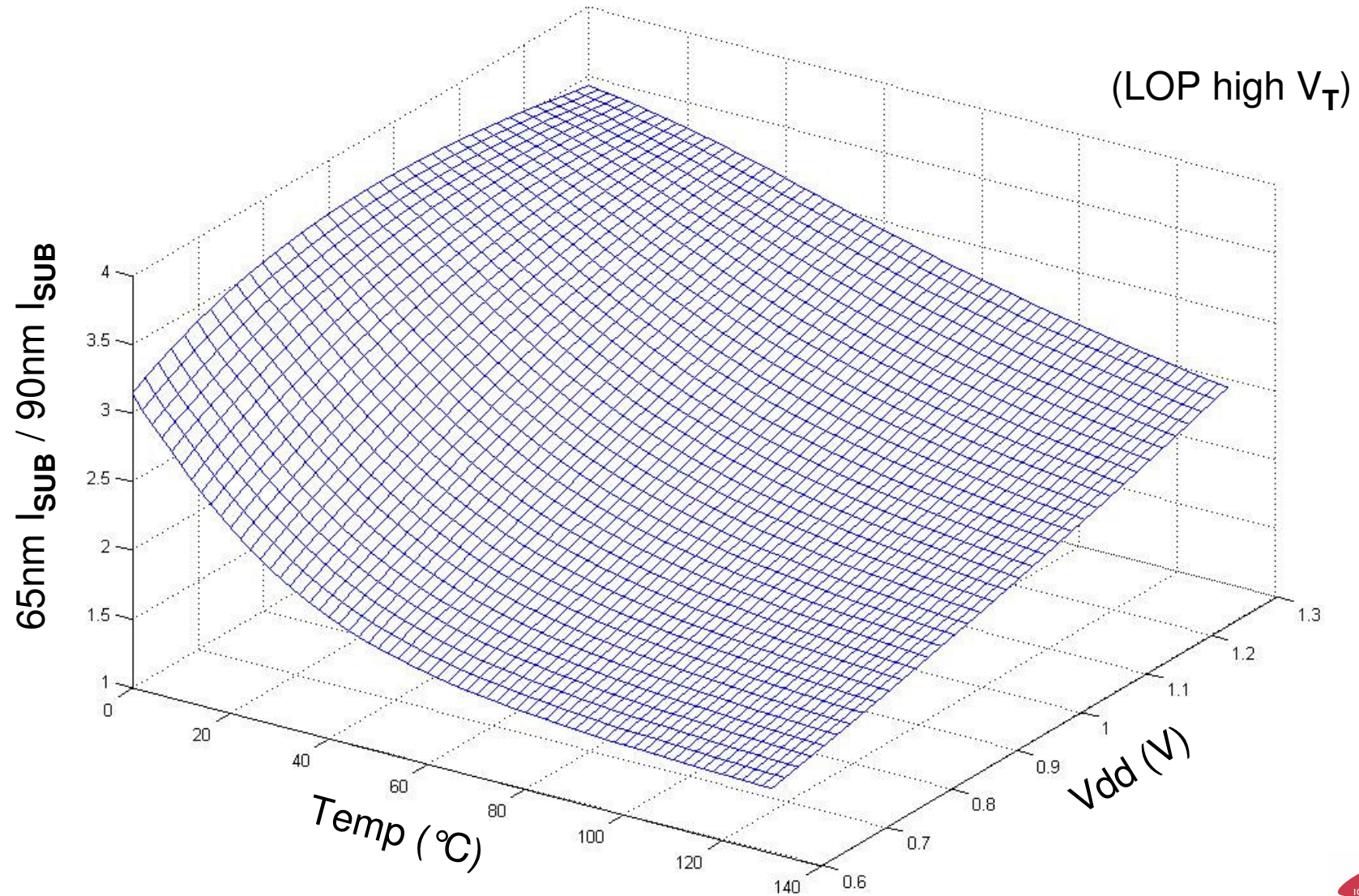
# 65nm LSTP vs LOP Sub-Threshold Leakage



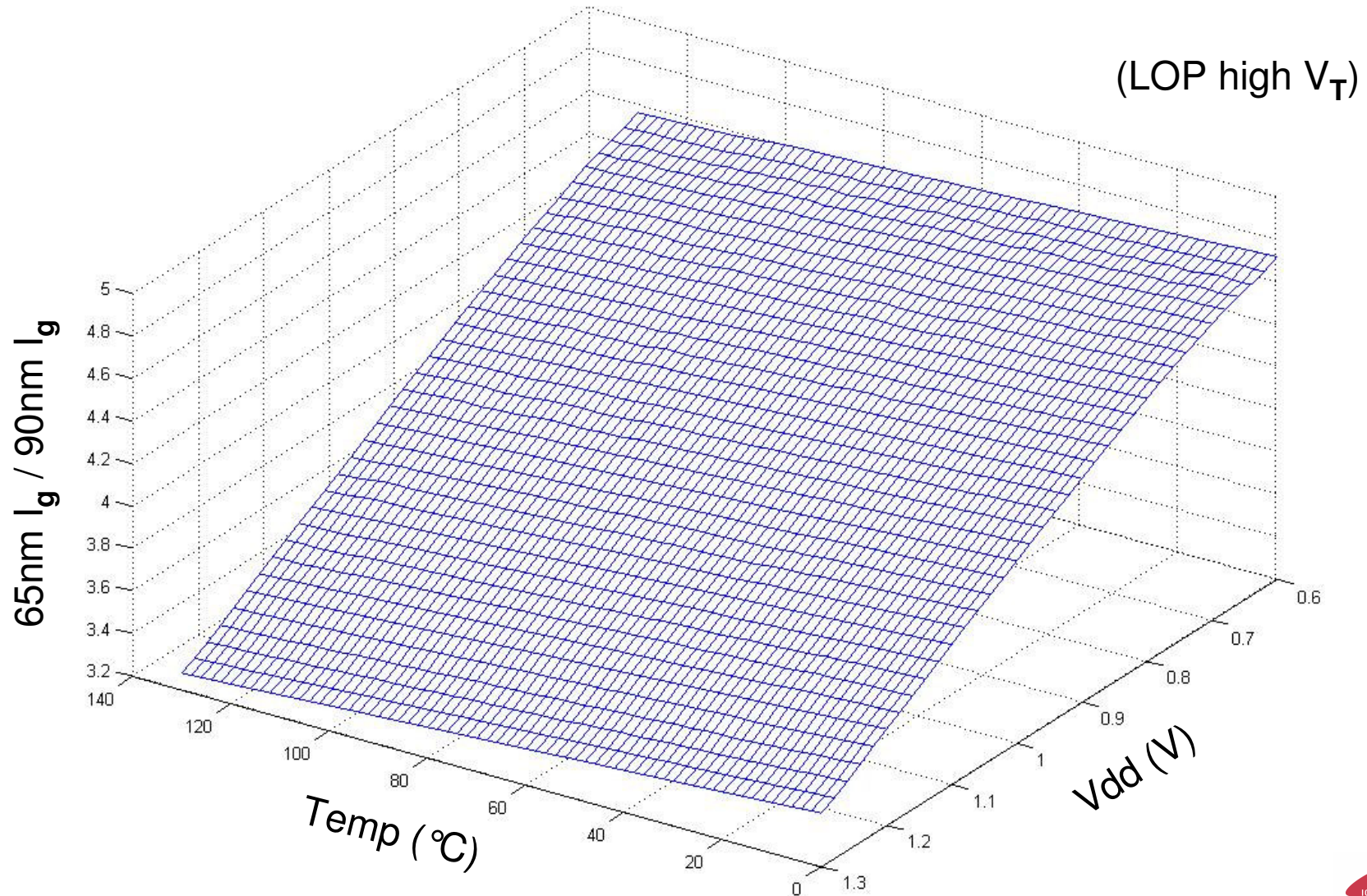
# High Temp $I_{SUB}$ Still Dominates $I_g$ at 65nm (LOP and LSTP)



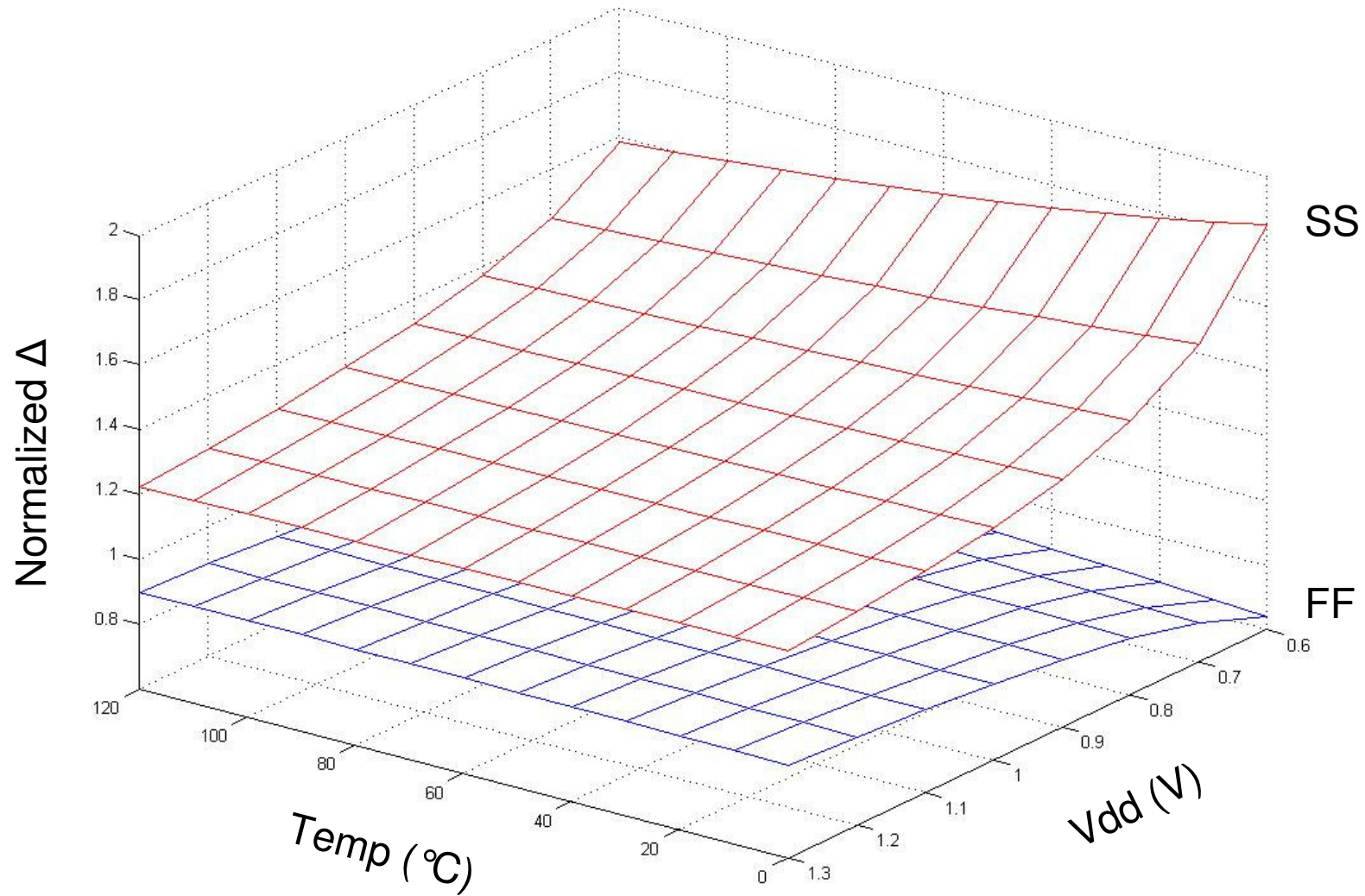
# LOP 65nm $I_{SUB}$ vs 90nm $I_{SUB}$



# LOP 65nm $I_g$ vs 90nm $I_g$

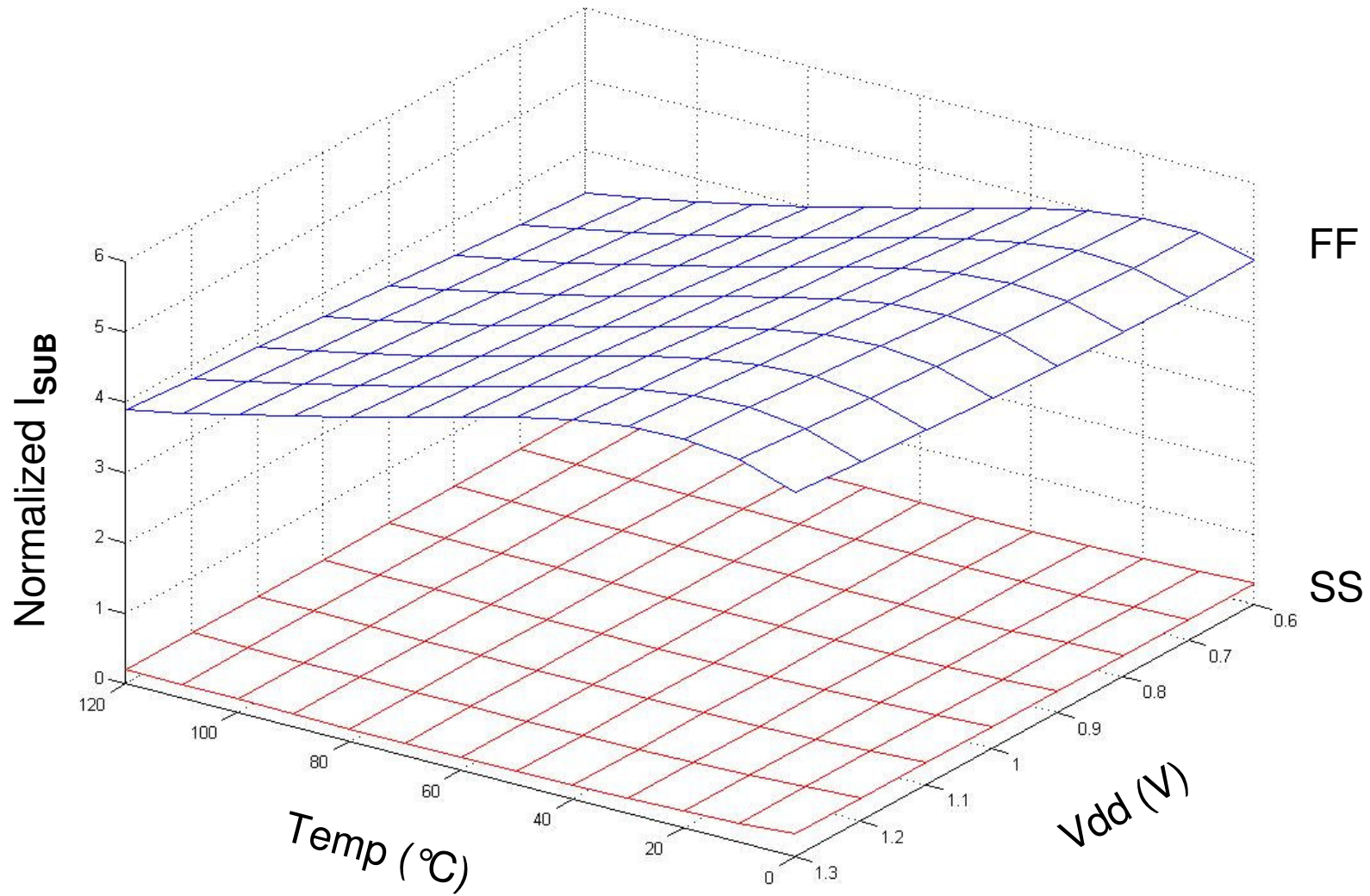


# Manufacturing Spread: 65nm mid $V_T$ LOP Delay





# Manufacturing Spread: 65nm mid $V_T$ LOP $I_{SUB}$



# Optimality of Design



# Power, Money and Performance

---

Excellent recent work on Power vs Performance (among many)...

- Srinivasan, Brooks, Gschwind, Bose, Zyuban, Strenski and Emma, “Optimum Pipelines for Power and Performance”, Micro '02.
- Zyuban and Strenski, “Unified Methodology for Resolving Power-Performance Tradeoffs at the Microarchitectural and Circuit Levels”, ISLPED '02.
- Harstein, Puzak, “Optimum Power/Performance Pipeline Depth”, Micro '03.
- Dao, Zeydel, Oklobdzija, “Energy Optimization of Digital Pipelined Systems Using Circuit Sizing and Supply Scaling”, Trans. VLSI Systems '05.

But outside the computer market, few chips sell on MHz.

Most designs set out to minimize cost (and power) at fixed performance...

- Decode an MPEG video without glitching.
- Execute an ADSL modem at 8Mbps.
- etc...

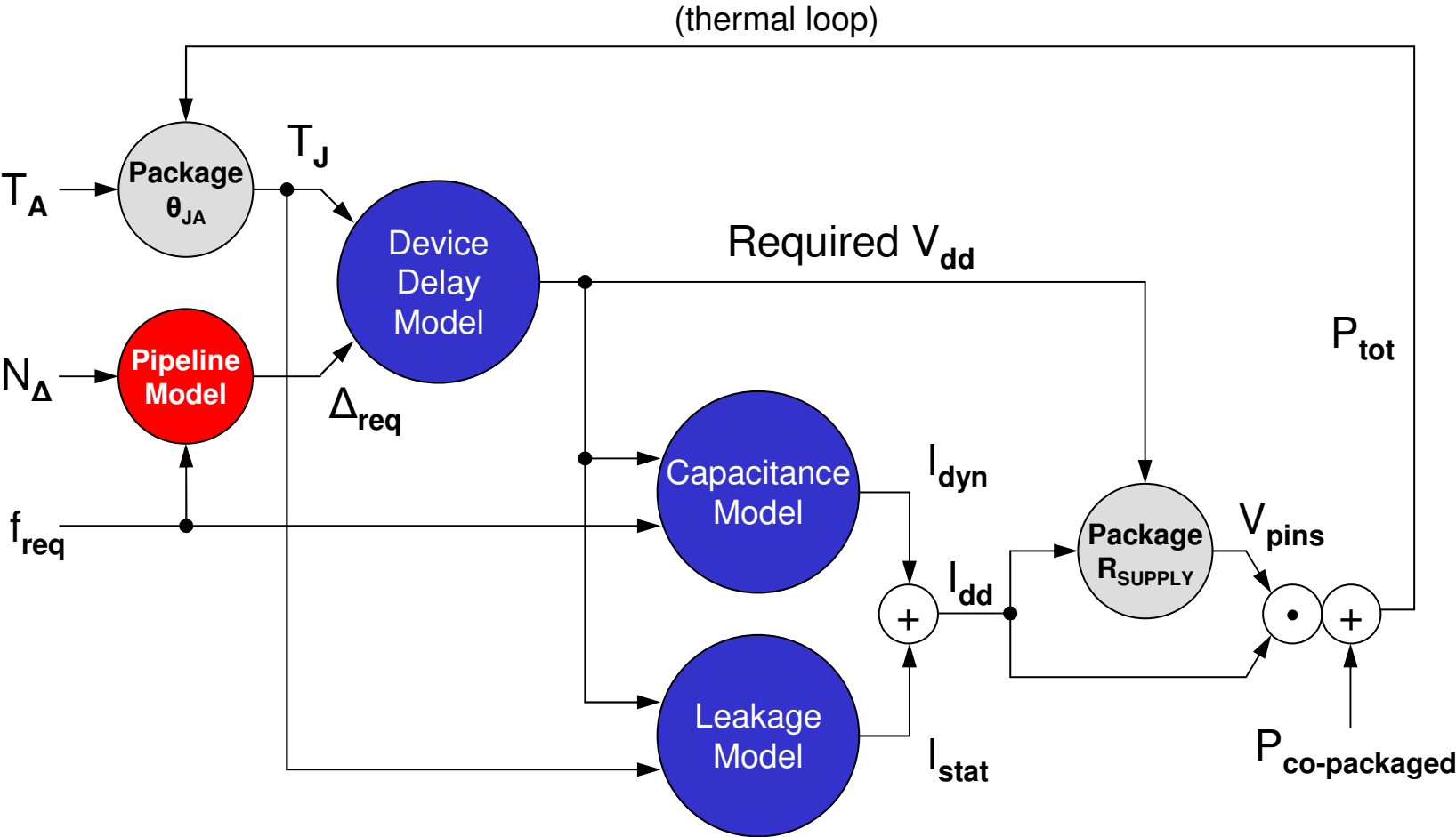


# Design Levers

- |                     |   |  |
|---------------------|---|--|
| (micro)architecture | { | <ol style="list-style-type: none"><li>1. Pipelining</li><li>2. Replication (parallelism)</li><li>3. Speculation</li><li>4. Logical redundancy (eg. carry-skip adder)</li><li>5. Arithmetic redundancy (eg. carry-save)</li></ol> |
|                     |   | 6. Vdd   |
| transistors         | { | <ol style="list-style-type: none"><li>7. <math>V_T</math></li><li>8. <math>t_{ox}</math></li><li>9. <math>L_g</math></li></ol>   |
| geometry            | { | <ol style="list-style-type: none"><li>10. Manhattan cell height (transistor sizing)</li><li>11. Wire geometry</li><li>12. Wire screening</li></ol>   |
|                     |   | 13. Structured layout (wire control)   |
| circuits            | { | <ol style="list-style-type: none"><li>14. Logic circuit engineering</li><li>15. Clock circuit engineering</li></ol>  |
|                     |   | 16. Package engineering  |

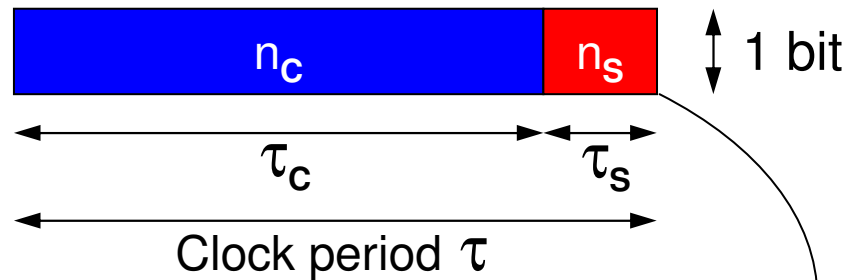


# Electro-Thermal Model



# Pipeline Model

K “pipettes”...



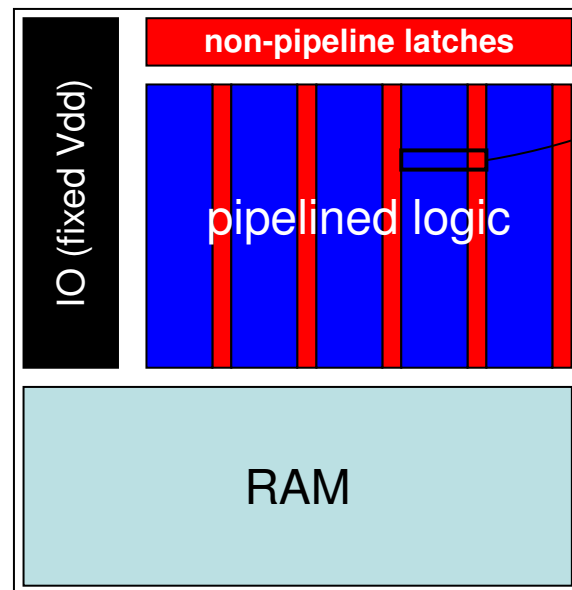
Average  $n_c$  transistors of combinatorial logic per pipe stage per pipeline bit.

$n_s$  transistors per pipeline latch.

Units of  $\tau$  are “i”  $\equiv$  “FO4”; 1i =  $\Delta$  seconds.

Arbitrary width/depth mix of pipettes make up the pipelined logic.

Chip...



Define...

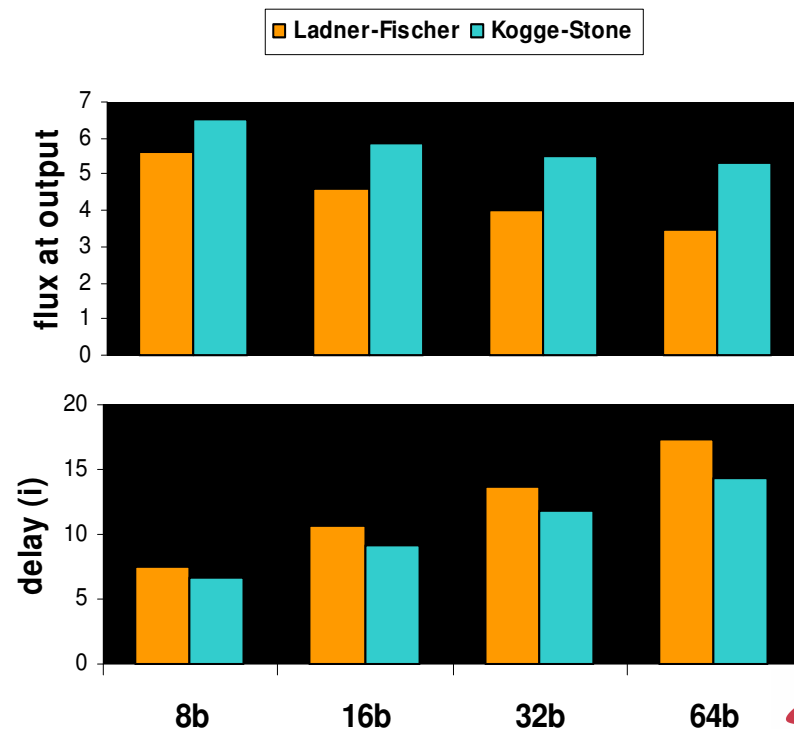
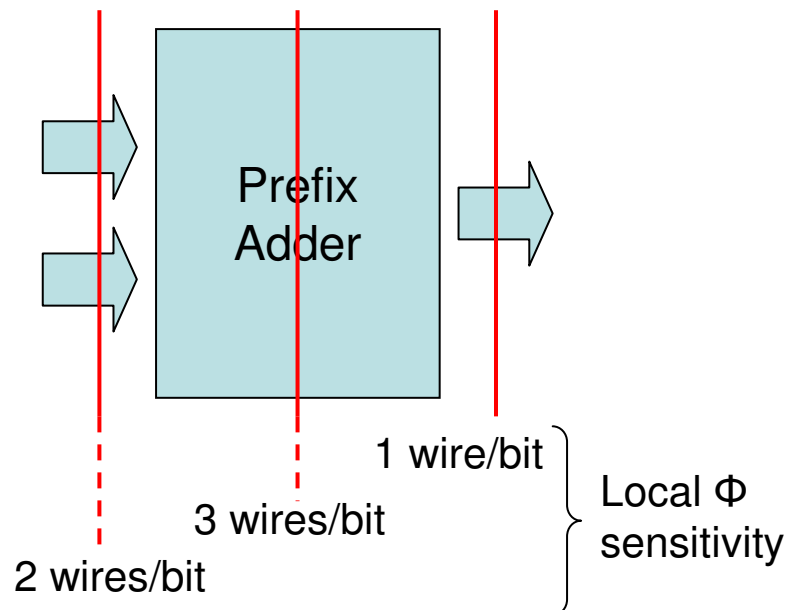
Flux	$\Phi = n_c / \tau_c$	(T/i/b)
Performance	$\Psi = K n_c f$	(T/s)
Cost	$N = K(n_c + n_s)$	(T)
Frequency	$f = 1 / \tau \Delta$	(Hz)



# Reality not Included...

- Performance  $\Psi$  ( $T_{\text{COMB}}/s$ ) tends to over-state effective performance at low  $\tau$ , because logical redundancy is introduced to speed up logic paths and mitigate latency by speculation.
- Flux  $\Phi$  (T/i/b) varies a lot locally; tends to rise at low  $\tau$  for the same redundancy reason, and because pipe cuts have less freedom to avoid high-flux logic.

Eg. adders...

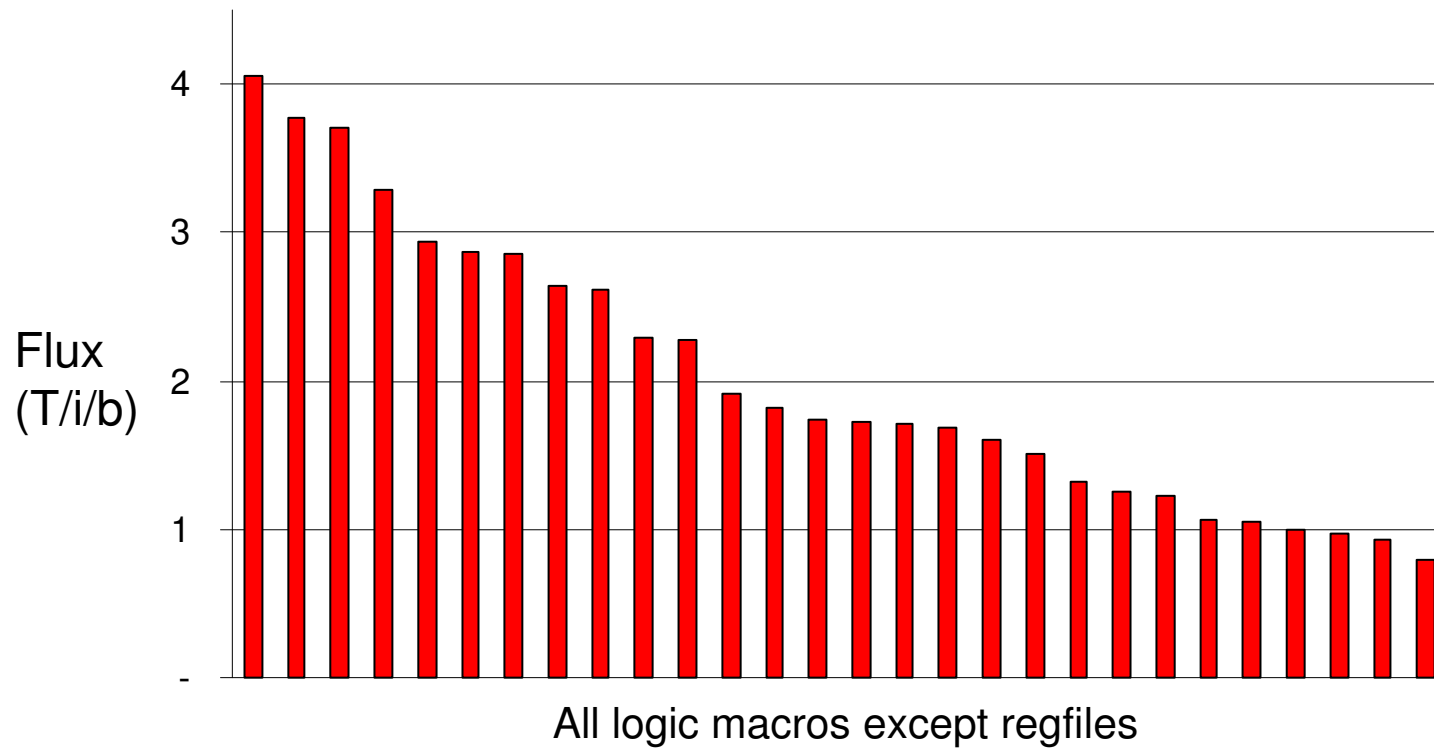


# Pipeline Flux in Icera's DXP®

---

$\Phi \sim 3-4$  for arithmetic and random translation logic.

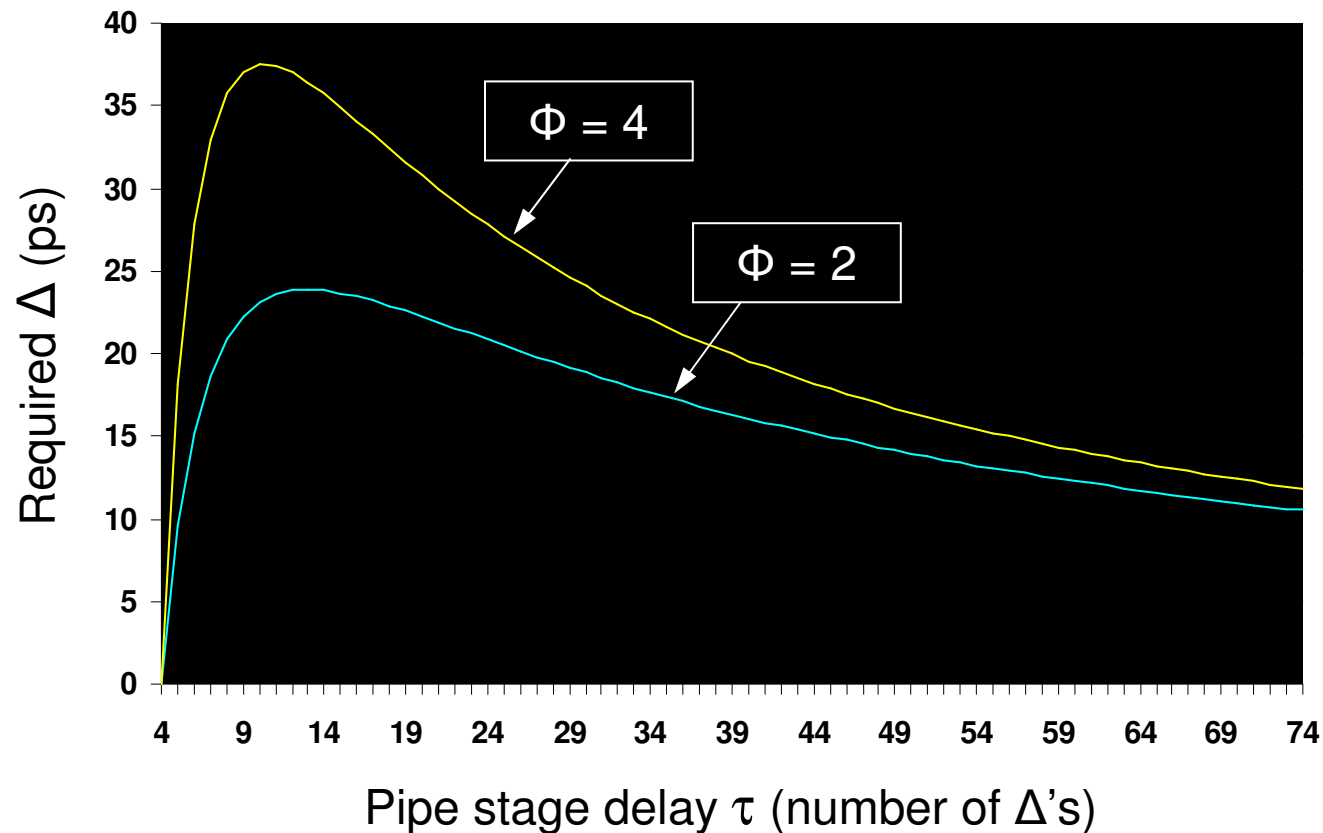
$\Phi \sim 1-2$  for steering logic.





# Required $\Delta$ to meet specified Performance $\Psi$ and Cost $N$

$$\Delta_{\text{required}} = \left( \frac{N}{\Psi\tau} \right) / \left( 1 + \frac{n_s}{\Phi(\tau - \tau_s)} \right)$$



## Specification

Latch delay 4i

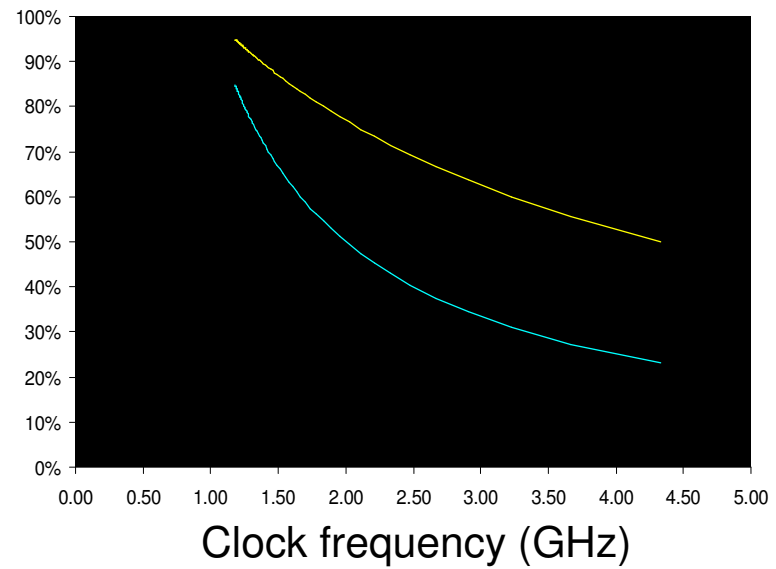
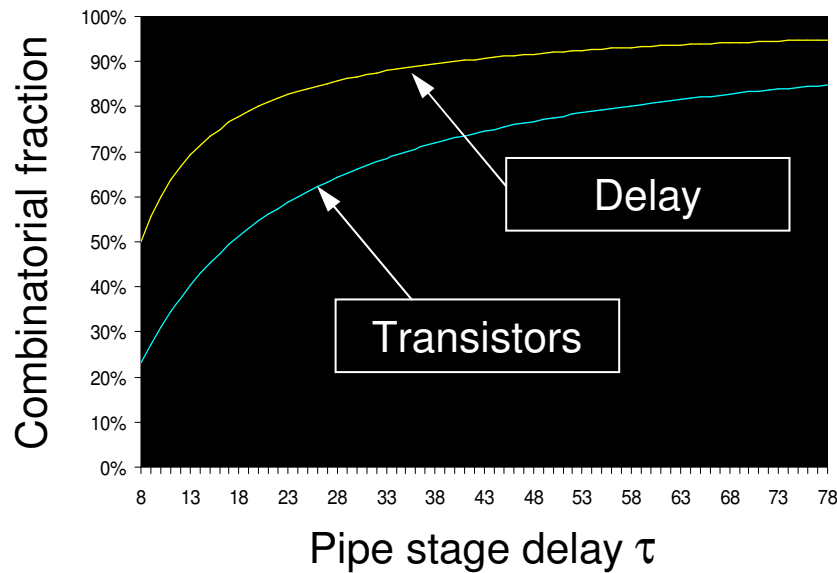
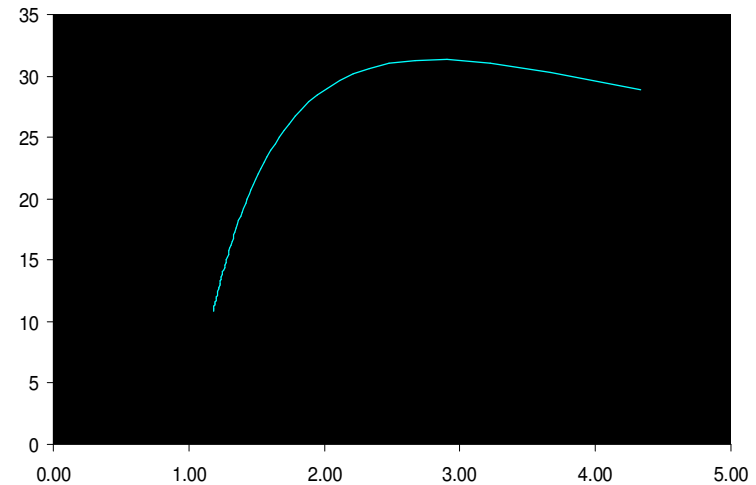
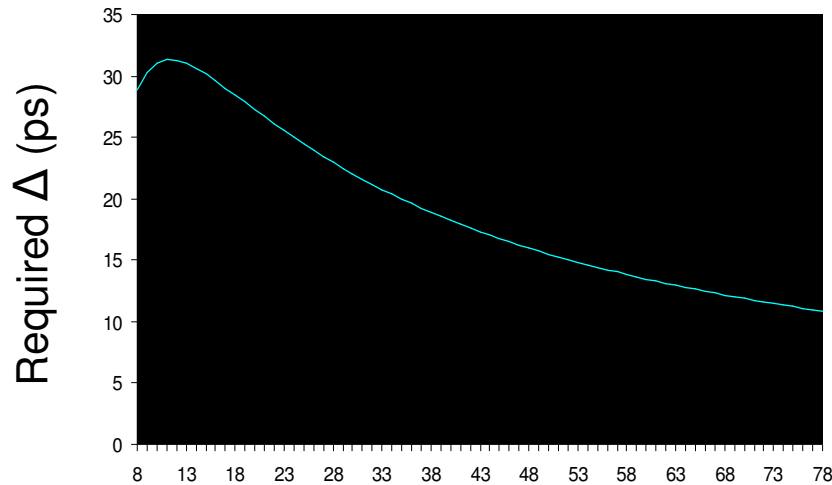
Latch cost 40T

Performance 10PT/s

Design Cost 10MT



# Constant Flux Implications ( $\Phi = 3, \tau_S = 4i$ )

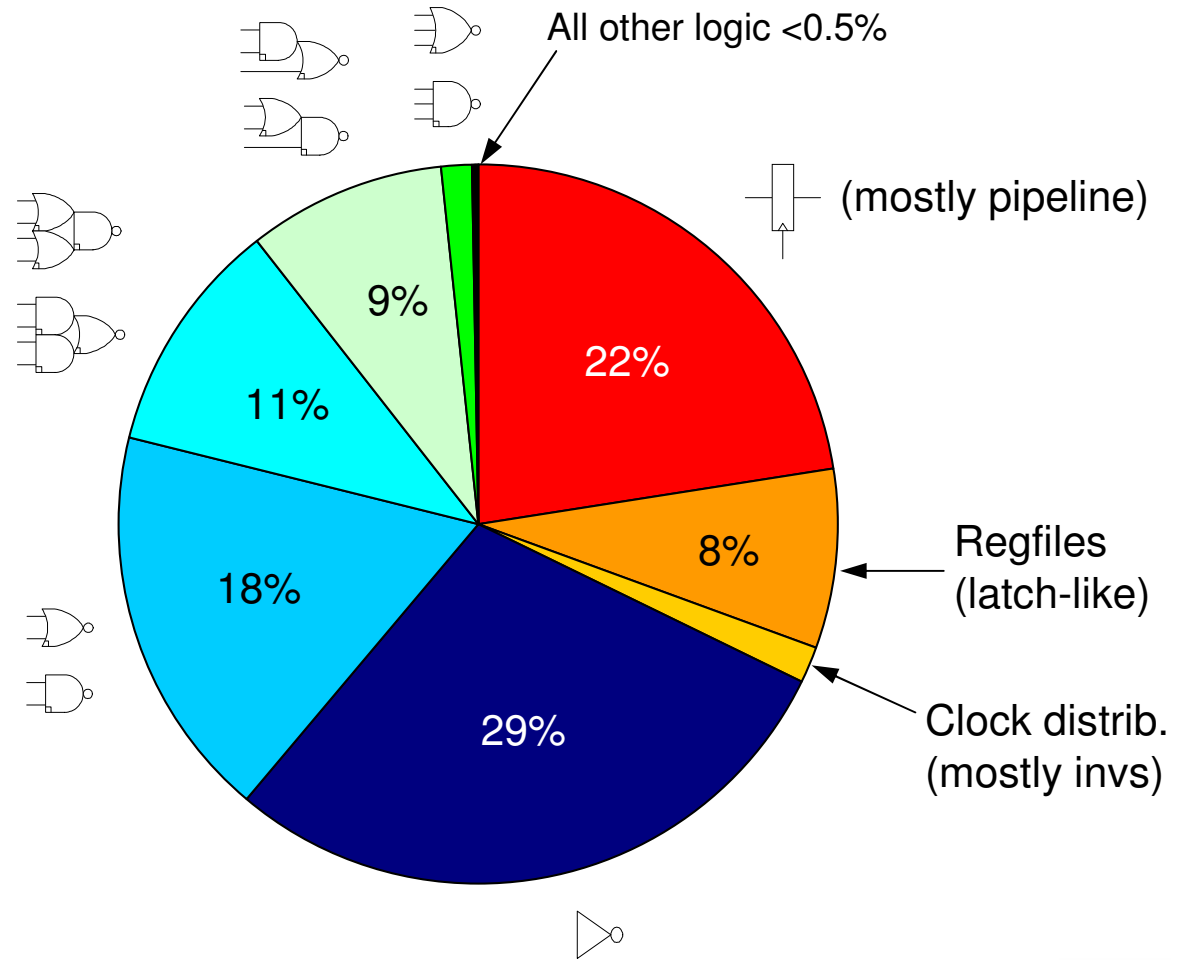


# DXP® Logic Stage Population, by #Transistors

Manhattan cells are ubiquitous in all digital design methodologies.

Complex cells can be important for performance but account for negligible % of total transistors.

Designs are characterized (capacitance, leakage, delay, area) by a few simple Manhattan stages.



# The Design Hull



# Example Design Spec

---

Performance  $\Psi = 15\text{PT/s}$

Flux  $\Phi = 4\text{T/i/b}$

Latch Cost  $n_s = 40\text{T}$

Latch Delay  $\tau_s = 4i$

Non-Pipe Logic =  $5\text{MT}$

RAM =  $1\text{MB}$

Latch Activity =  $3\text{x Combinatorial Activity}$

RAM Activity =  $0.2\text{x Combinatorial Activity}$

65nm high  $V_T$  LOP Process

Ambient Temp =  $85^\circ\text{C}$

Max Die Temp =  $125^\circ\text{C}$

Package  $\theta_{JC} = 15^\circ\text{C/W}$

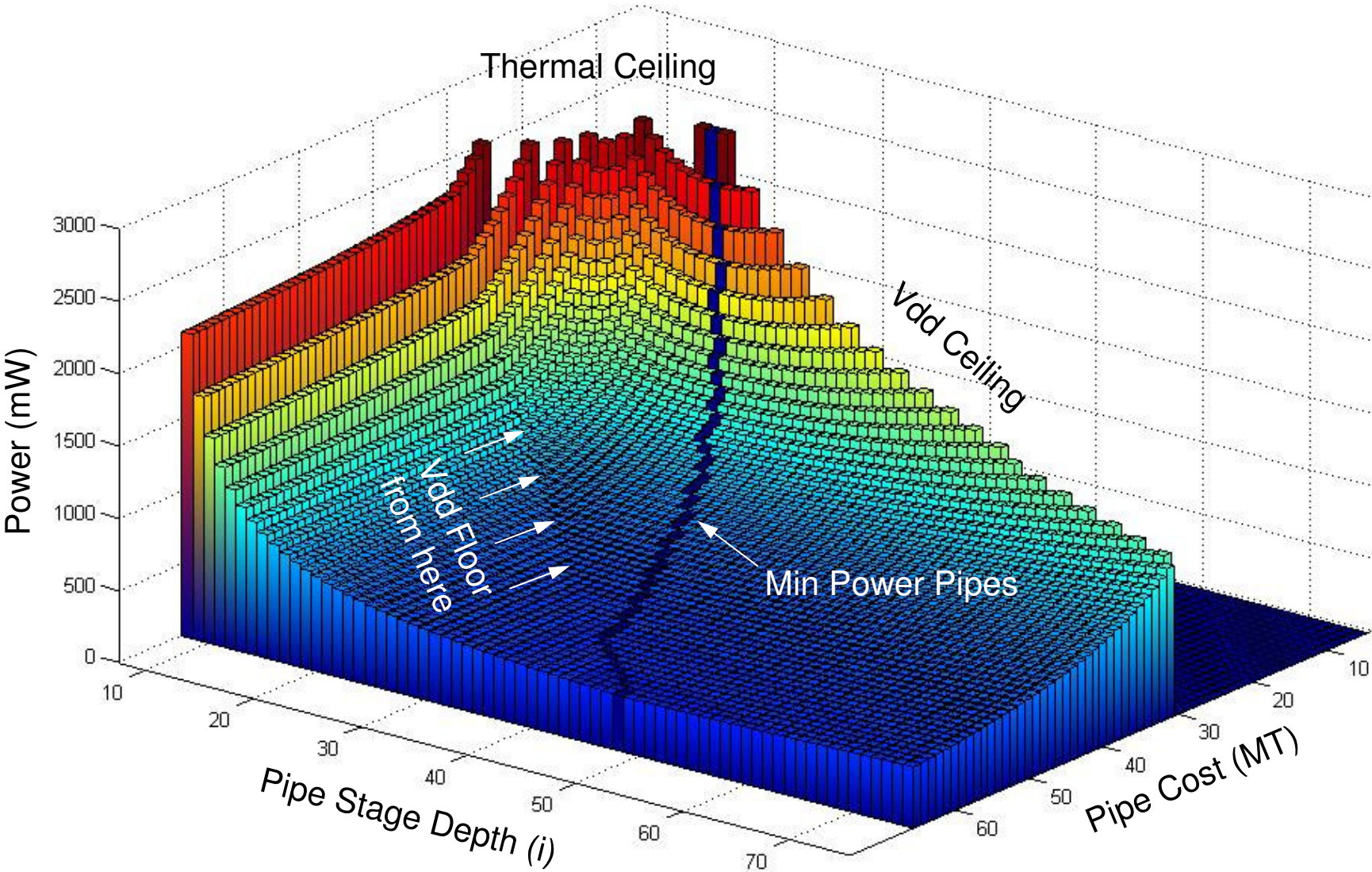
Supply R =  $50\text{m}\Omega$

Fixed Parasitic Power =  $100\text{mW}$

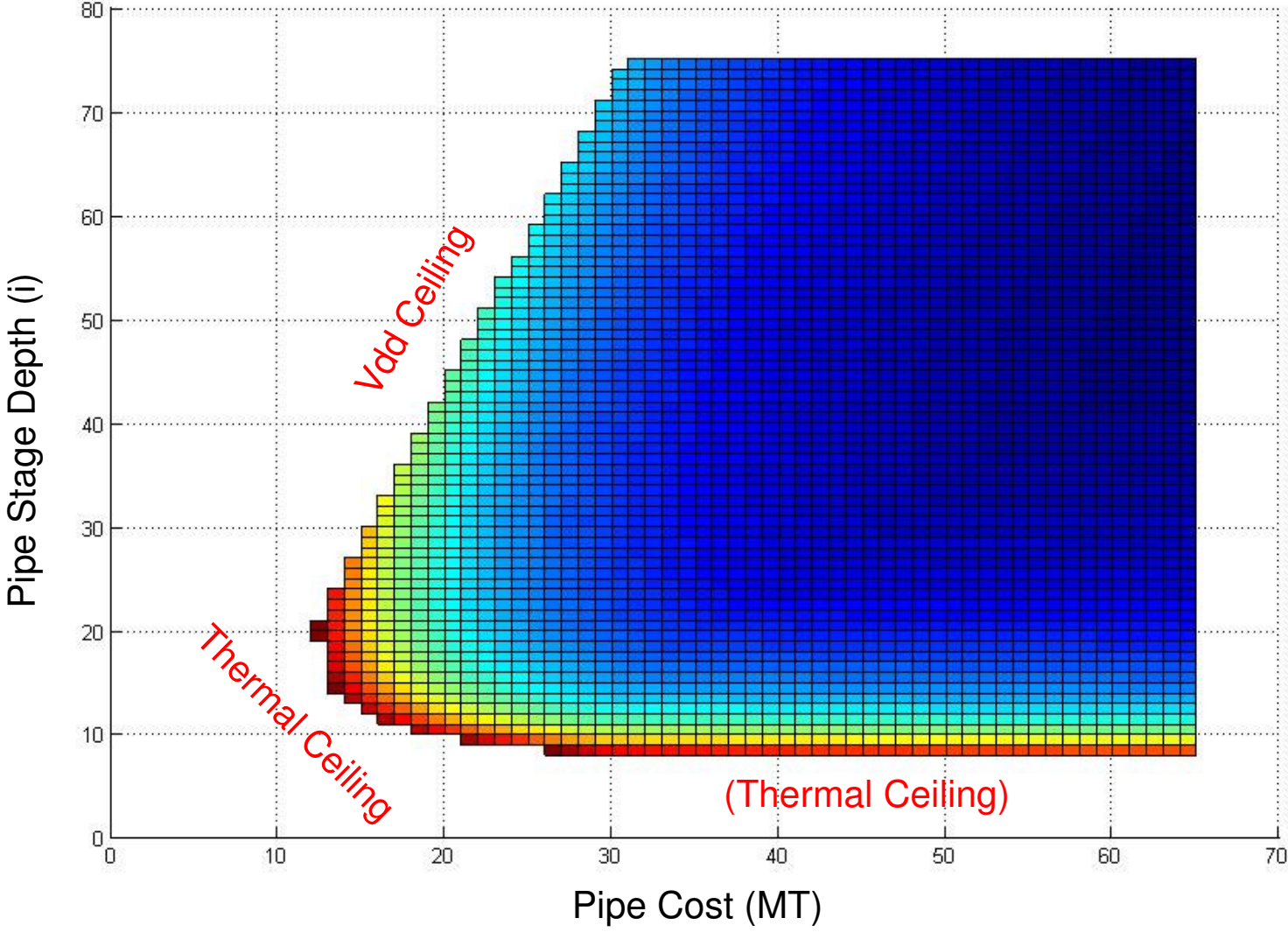
Vdd Limit Range =  $0.7\text{V} - 1.2\text{V}$



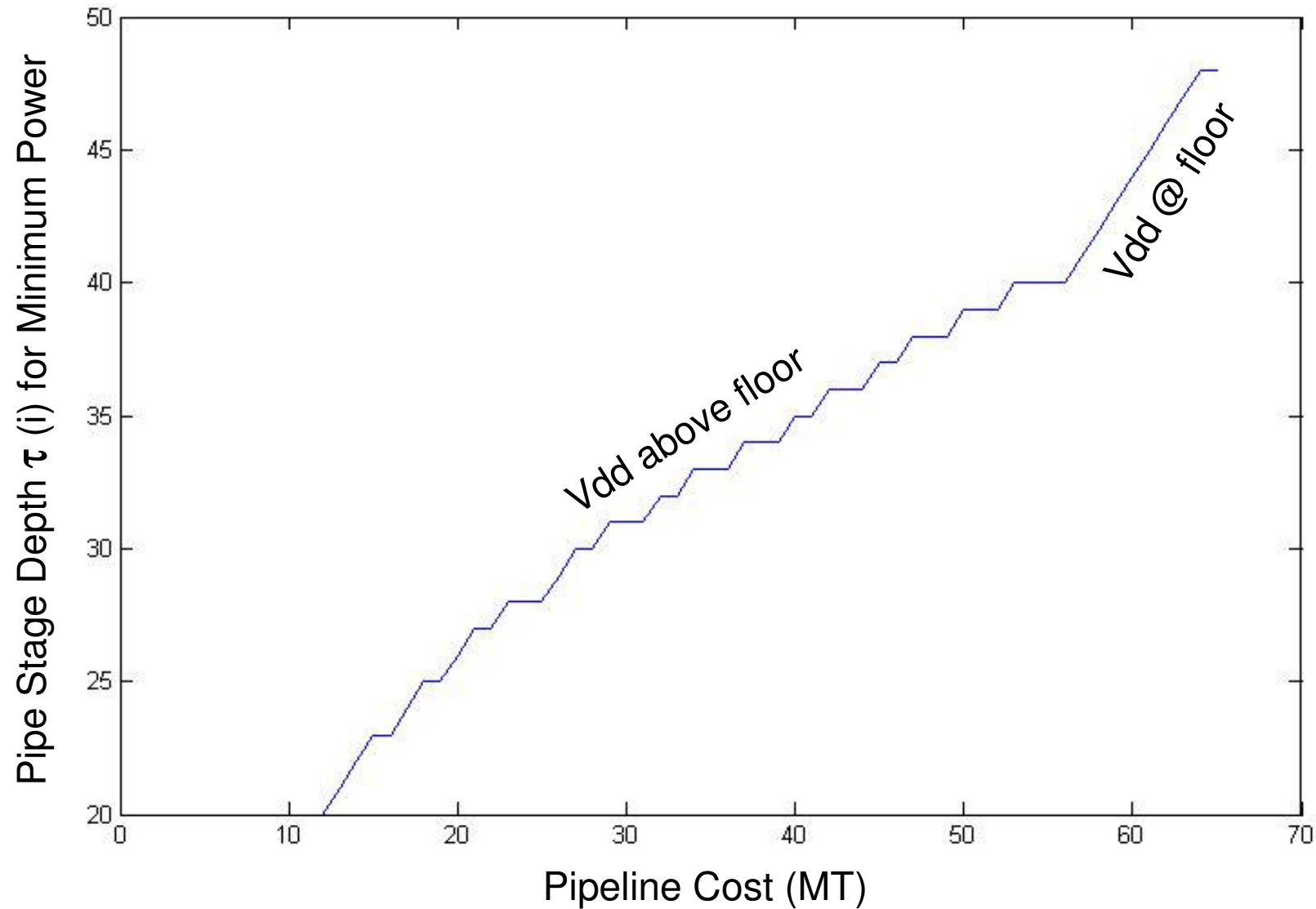
# Design Hull



# Feasible Designs

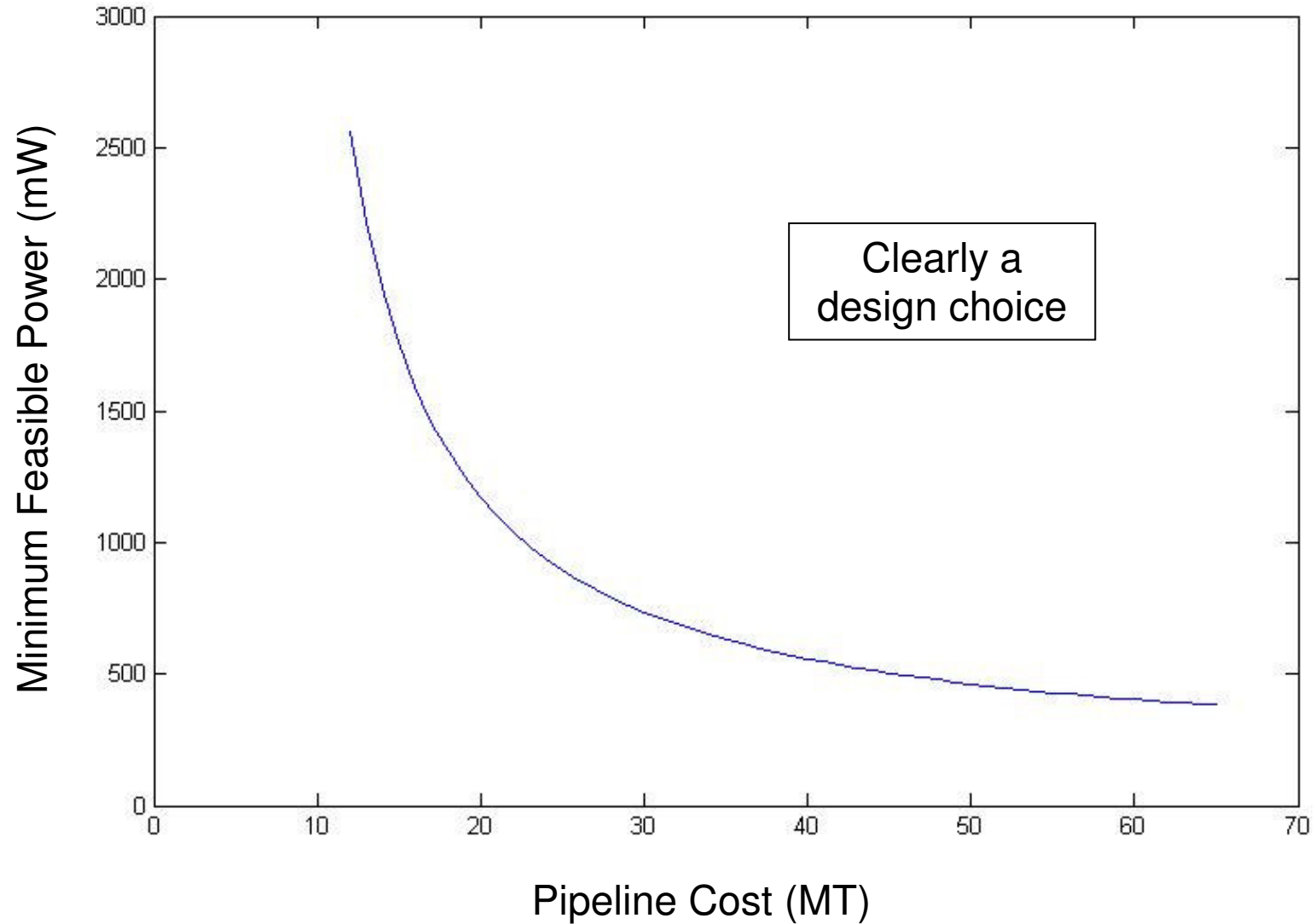


# Minimum Power Pipelines

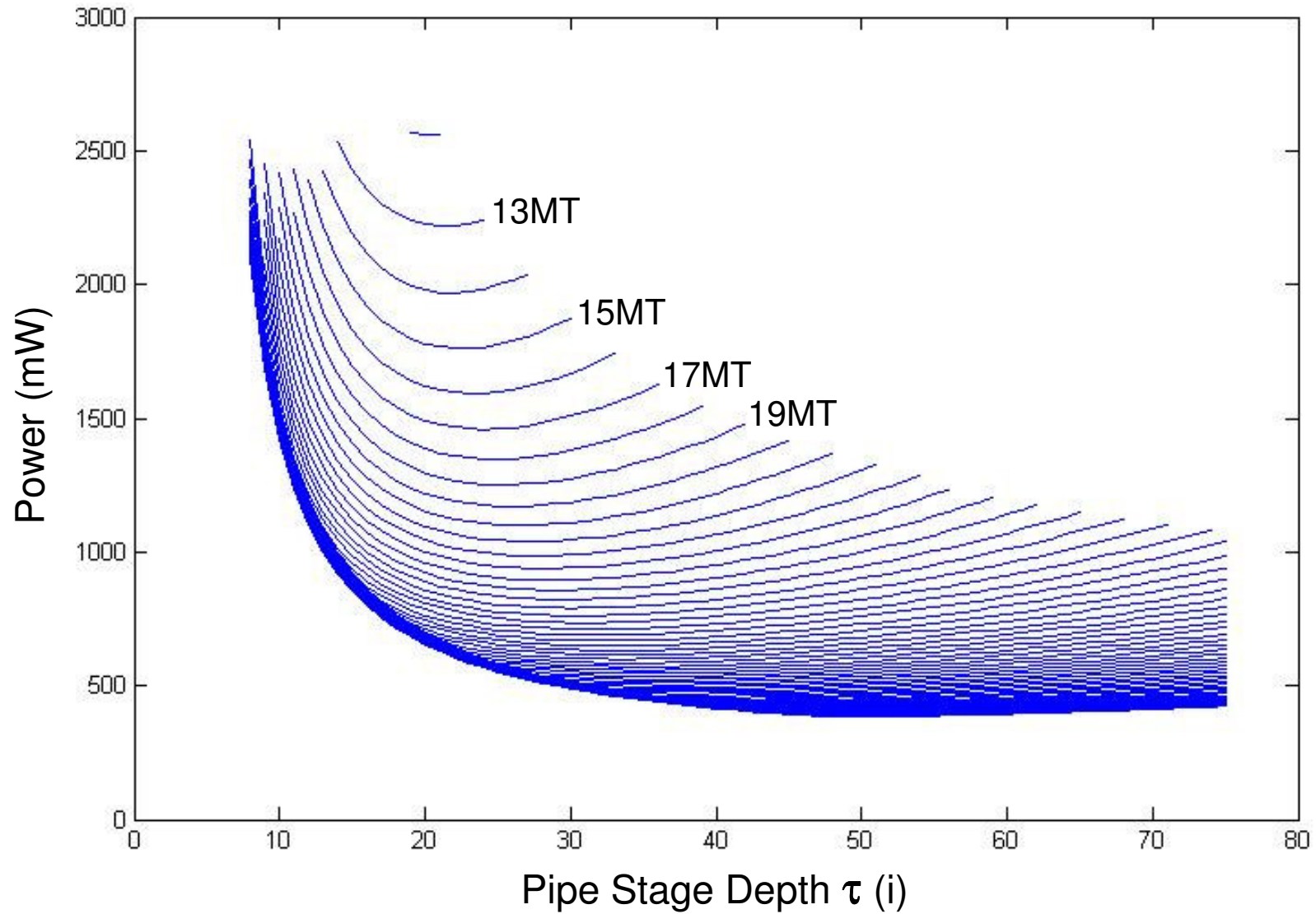




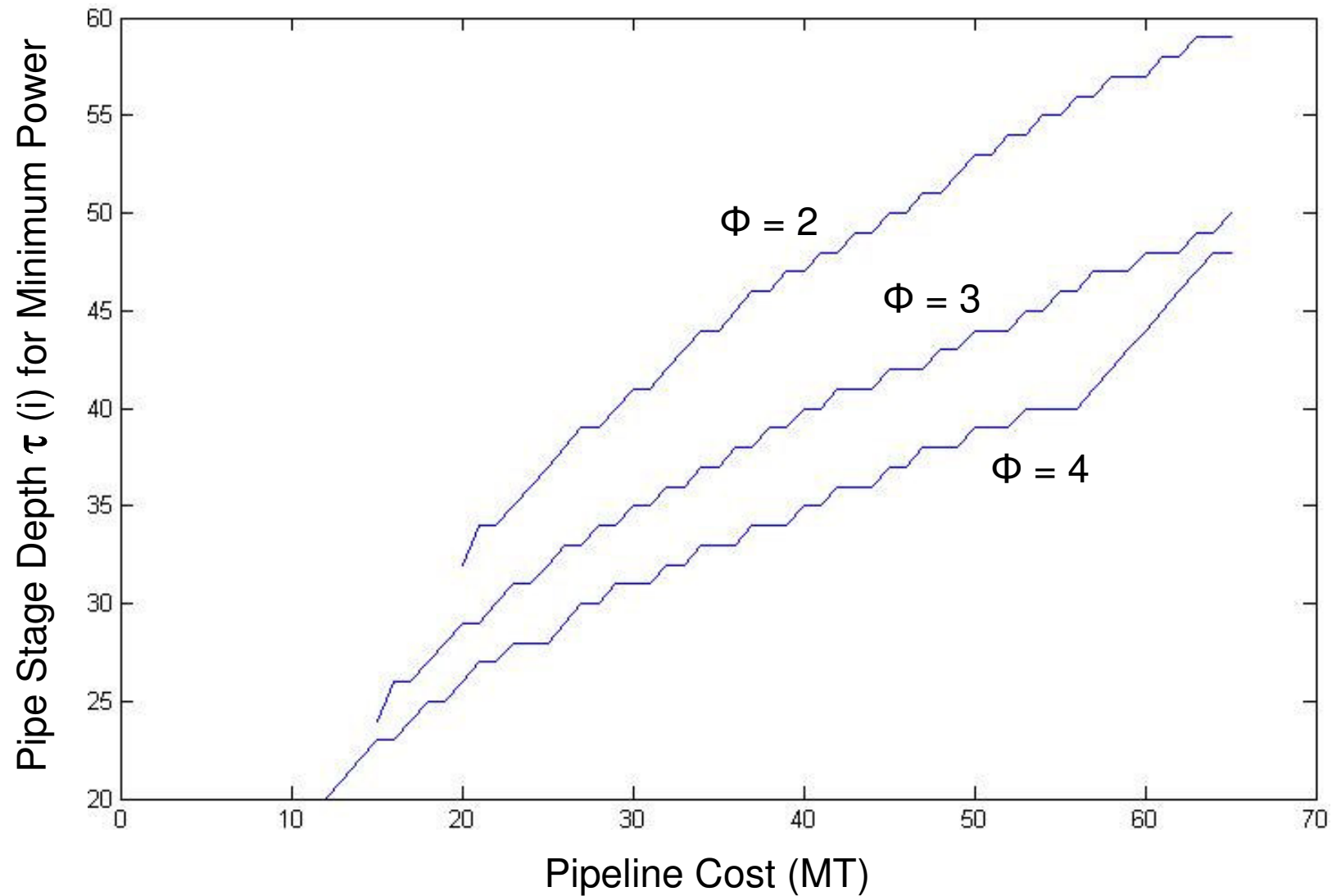
# The Cost of Minimizing Power



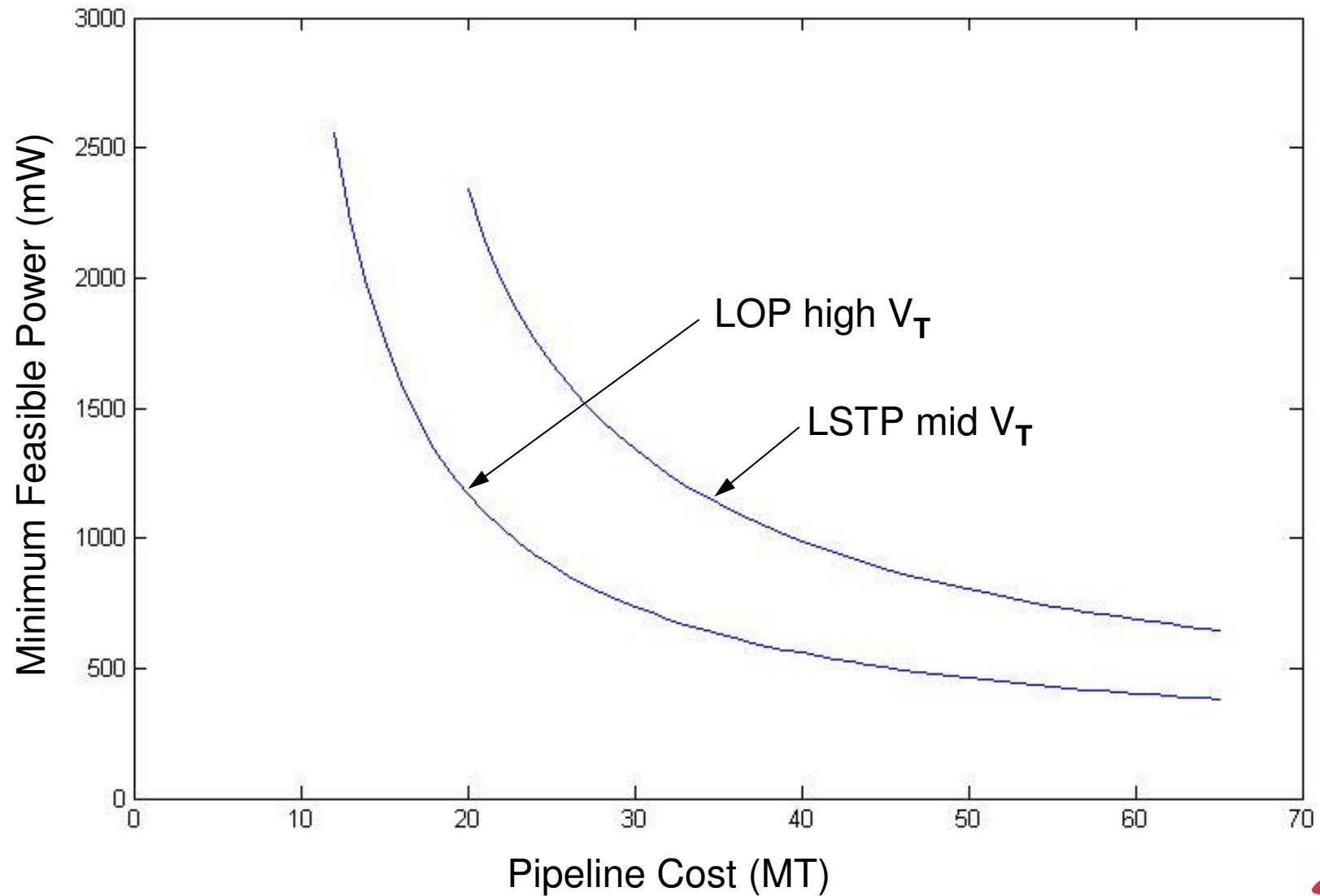
# Minimum Power Pipelines



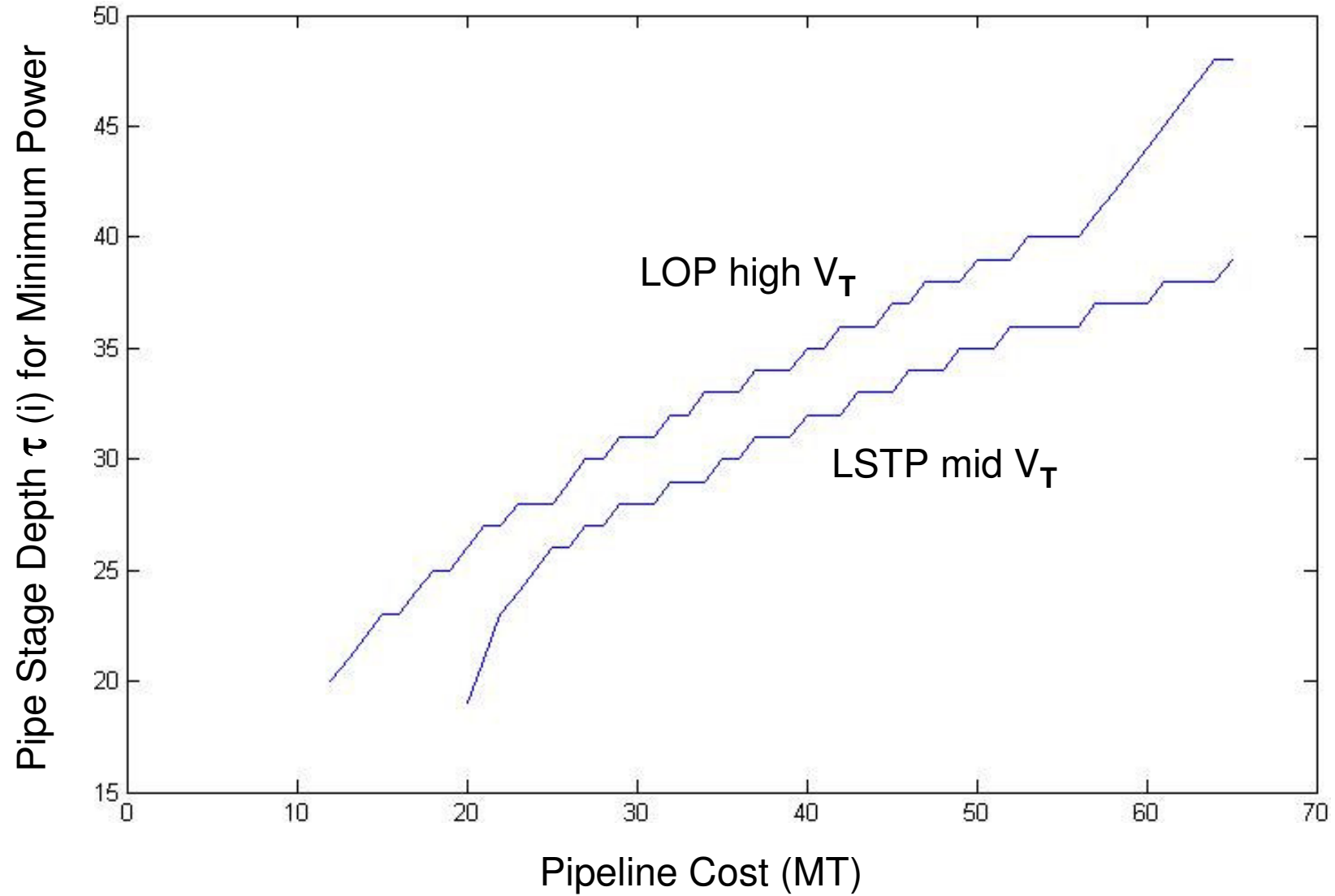
# Sensitivity to Pipeline Flux



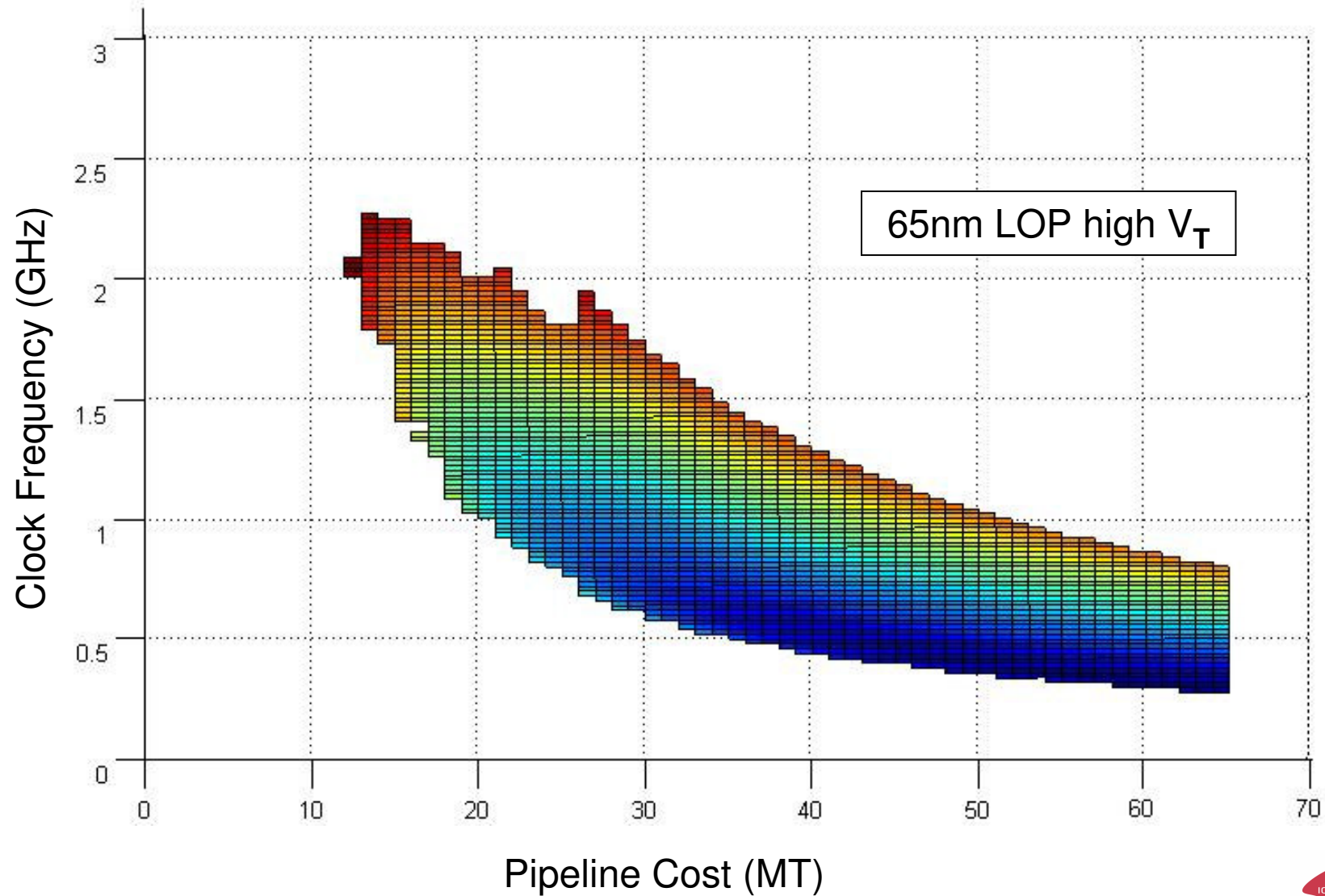
# LOP vs LSTP



# LOP vs LSTP



# Feasible Solutions in Cost-Frequency Space



# Implications for Design



# Implications for Chip Design

---

**Every logic chip will have control over its supply voltage.**

- Integrated regulators?
- Switch off in standby – limited call for LSTP processes?

**We can't minimize cost and power simultaneously.**

**Good cost-power points require fairly high speed.**

- Beyond synthesis and auto-P&R today → more Structured Custom?

**Short pipelines expose the basic TA assumption (single path sensitization).**

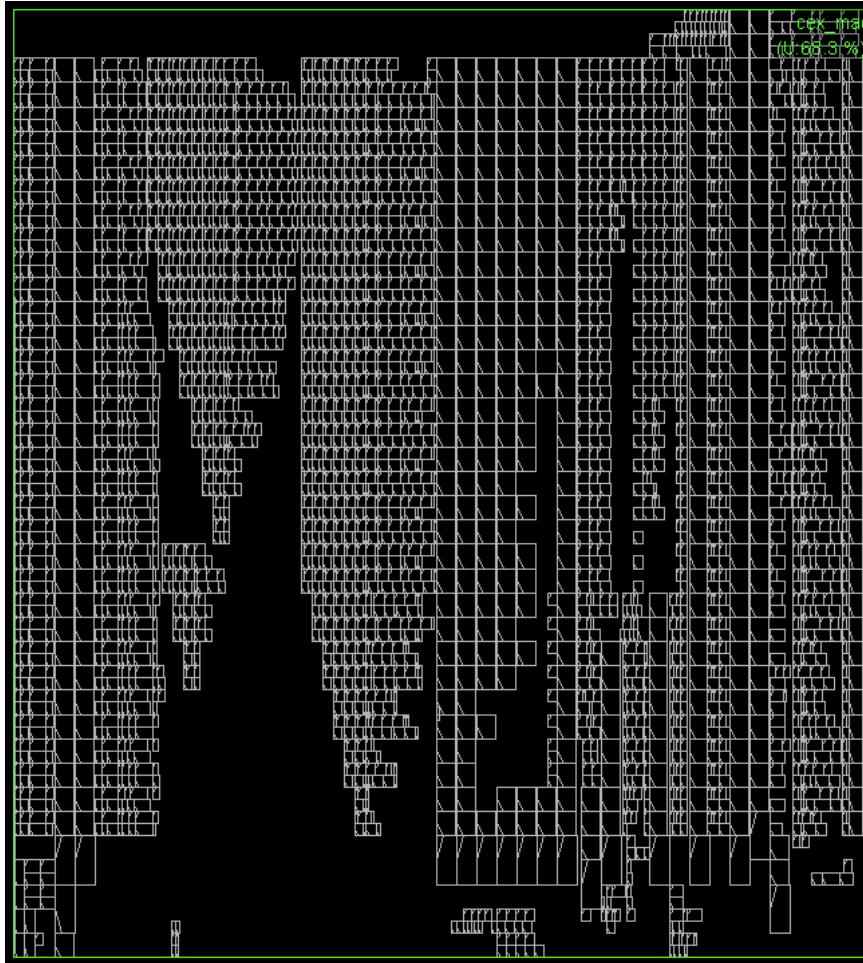
- Wavefront timing analysis tools?
- Voltage agility makes TA sign-off harder.



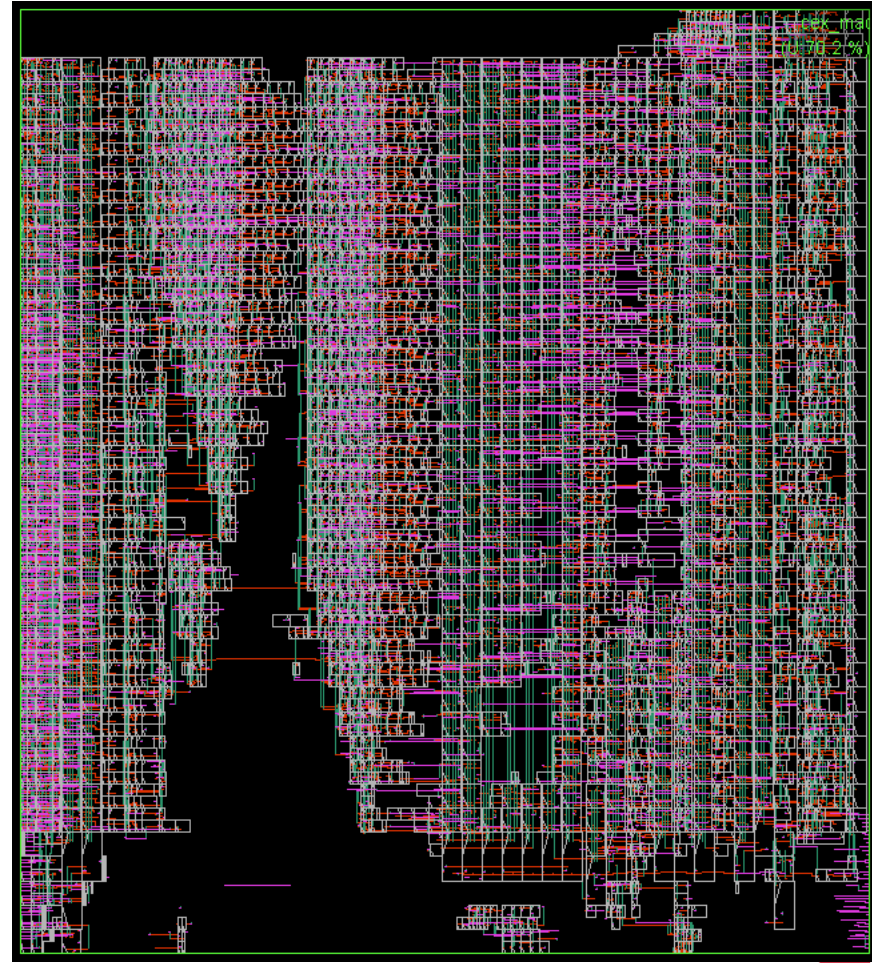


# Structured Custom

Natural cell placement by engineers...



...is enough for predictable auto-routing



Thank You

---

Enjoy the Symposium 😊