# Inphi ®
## *Think fast.*

# Inphi Moves Big Data Faster

# The Evolution, Pitfalls, and Cargo Cult Engineering of Advanced Digital Timing Sign-off

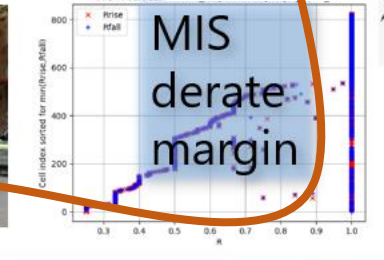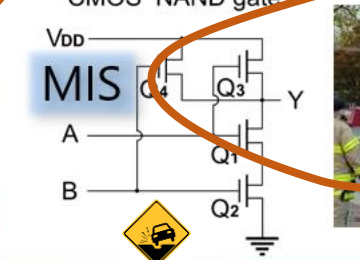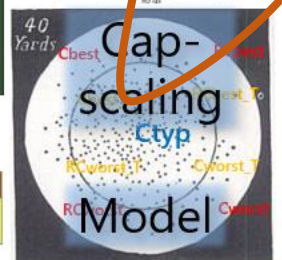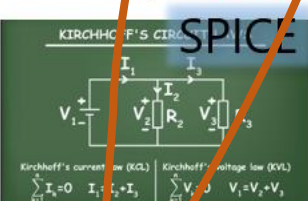Christian Lütkemeyer, April 2021

τ 2021

# Welcome on Board!

- This will be a fast-paced review of STA innovation over the last 15 years, from my vantage point in the trenches of SOC design.

- The audience is experts in the field, You!

- We will fly at 10,000 meters most of the time, cruising speed is 90 slides/h.

- Prepare for a few dive-downs for closer inspection at interesting vista points!

- Slides will be shared to walk along our route at your own pace again later.

Buckle up and enjoy the movie! ☺

# The Complex Maze of Robust Digital Design and Timing Sign-Off

# MOSFET Scaling and Disruptive MOSFET Device Innovation

▪ Moore's Law: Semiconductor manufacturing has gone through an amazing exponential scaling process for 50 years that has enabled the exponential growth of device complexity.

▪ Gate lengths reduced from 10um to 5nm (a 2000x reduction).

▪ Transistor complexity grew from Thousands (k) to Billions (G) (x4M)

▪ Disruptive innovation examples
- ■ Multi VT (65nm)
- ■ Strained silicon
- ■ Double patterning
- ■ High-k gates (leakage reduction)
- ■ FinFET

**MOSFET scaling**
(process nodes)
- 10 µm – 1971
- 6 µm – 1974
- 3 µm – 1977
- 1.5 µm – 1981
- 1 µm – 1984
- 800 nm – 1987
- 600 nm – 1990
- 350 nm – 1993
- 250 nm – 1996
- 180 nm – 1999
- 130 nm – 2001
- 90 nm – 2003
- 65 nm – 2005
- 45 nm – 2007
- 32 nm – 2009
- 22 nm – 2012
- 14 nm – 2014
- 10 nm – 2016
- 7 nm – 2018
- 5 nm – 2020
- Future 3 nm – ~2022



Moore's Law – The number of transistors on integrated circuit chips (1971-2018)
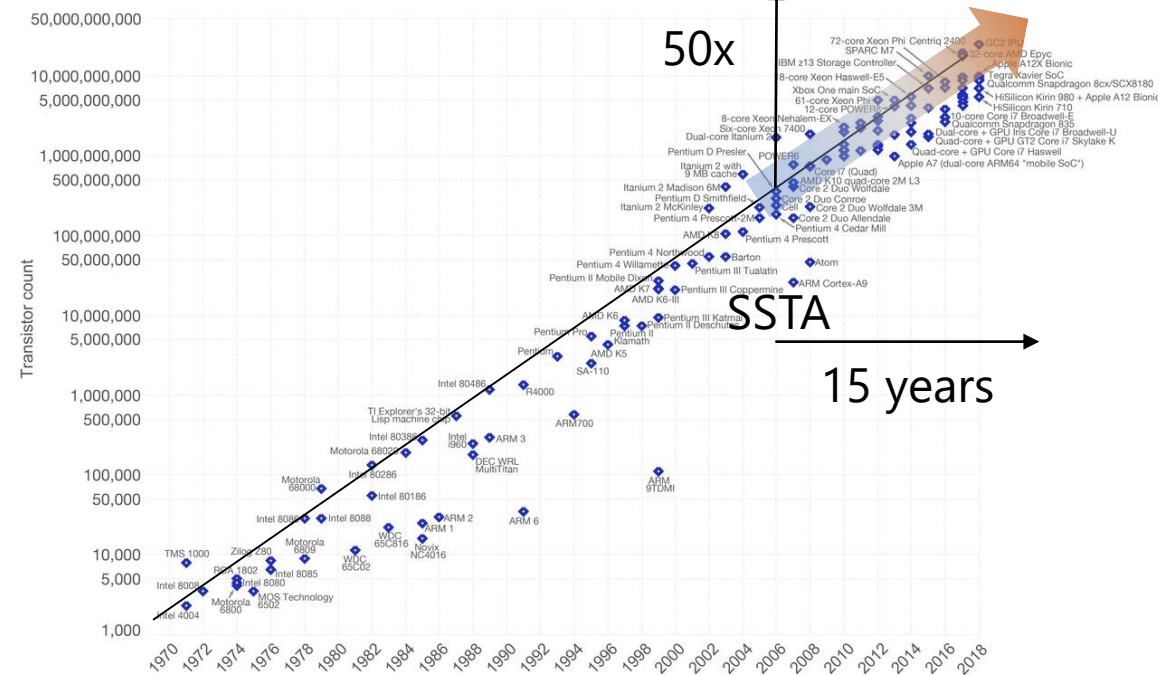Moore's law describes the empirical regularity that the number of transistors on integrated circuits doubles approximately every two years. This advancement is important as other aspects of technological progress – such as processing speed or the price of electronic products – are linked to Moore's law.

Data source: Wikipedia (https://en.wikipedia.org/wiki/Transistor_count)
The data visualization is available at OurWorldinData.org. There you find more visualizations and research on this topic.

Licensed under CC-BY-SA by the author Max Roser.

50x

SSTA

15 years

4

# The Innovation Domains of Digital Timing Verification

Primary drivers:

- MOSFET scaling and disruptive device innovation



- Interconnect scaling



Top mask cost: few $k



Bottom: few $100k [1]



- SPICE model evolution and methodology changes

- Liberty library model innovation and sophistication

- Increasing level of detail in Timing sign-off margin modeling

# Sources of Timing Variation in Digital Circuits

## Process



SS

FF

## Operating voltage



-10%, -10%
$VDD_{nom}, VDDM_{nom}$
+10%, +10%

## Operating temperatures



$0^{o}C$

$125^{o}C$

## Local variation

$+"3\sigma"$

$-"3\sigma"$



## Spatial variation

$-x$
P, V, T
$+y$



## Interconnect

Cworst
RCbest
RCworst
Cbest

uncorrelated layers, correlated models ☹



## Coupling (SI)



$\frac{5}{3}C$

$\frac{1}{3}C$

## Simultaneous Switching

CMOS NAND gate



$1.1\ \tau$

$0.5\ \tau$

## Aging



$1.x\ \tau$

$1.0\ \tau$

# The Timing Variation Matrix

| Variation parameter | Range | Scenario multiplier |
|---|:---:|:---|
| Process | SS, FF | 2 |
| Operating voltage corners | $[VDD_{low}, VDD^{high}], [VDDM_{low}, VDDM^{high}]$ | 4 |
| Operating temperature corners | $0^oC, 125^oC$ | 2 |
| Interconnect corners | Cworst, RCbest, RCworst, Cbest | 4 |
| Coupling | $\frac{1}{3}C, \quad \frac{5}{3}C$ | |
| Local variation | $-3\sigma, +3\sigma$ | |
| Process gradients | +/- 10%/10mm | |
| Voltage gradients | +/-x mV | |
| Temperature gradients | +/-10K | |
| Interconnect variation (layer-to-layer) | 0.9, 1.1 | |
| Multi-Input Switching | 0.25, 1.1 | |
| Aging | 1.0, 1.x | |
| Modeling errors vs. SPICE | -3%, 3% | |
| Silicon to SPICE model gap | -3%, 3% | |

$$\prod = 64 \text{ scenarios}$$

# A Simple 65nm OCV Signoff Margin Model (circa AD 2005)

| Variation parameter | Range | Scenario multiplier |
|---|---|---|
| Process | SS, FF (total corner models) | 2 |
| Operating voltage corners | $VDD_{low}$, $VDD^{high}$ (no VDDM rail) | 2 |
| Operating temperature corners | 0°C, 125°C | 2 |
| Interconnect corners | Cworst, RCbest, RCworst, Cbest | 4 |
| Coupling | $\frac{1}{3}C$,  $\frac{5}{3}C$ | |
| Local variation | $-3\sigma, +3\sigma$ | |
| Process gradients | +/- 10%/10mm | |
| Voltage gradients | +/-x mV | |
| Temperature gradients | +/-10K | |
| Interconnect variation (layer-to-layer) | 0.9, 1.1 | |
| Multi-Input Switching | 0.25, 1.1 | |
| Aging | 1.0, 1.x | |
| Modeling errors vs. SPICE | -3%, 3% | |
| Silicon to SPICE (S2S) model gap | -3%, 3% | |

$\prod$ = 32 scenarios

- -/+10% flat OCV derate, independent of path length



- 100ps uncertainty margin

- This flow reduces pessimistic margin stacking by implicitly allowing for statistical cancelations between variation effects.

Inphi

# The Saving Grace of Fixed Voltage Design in the Good Old OCV Days



- The foundries asked customers to close timing with +/-3 sigma total corner transistor models, combined with +/-3 sigma interconnect models. The actual manufacturing precision was tighter, close to ~ +/-1.5 sigma.
  => The foundries kept significant safety margin for themselves and for their customers!

- Only a small number of chips could be marginal (close to SS corner silicon), if hold timing was robust!
  - =>Aging is low risk, only noticeable on actual SS silicon.
  - =>Supply noise related setup deficiencies mostly appear on SS silicon.

Foundries leveraged their manufacturing margins to improve yield by tweaking process targets, or to create second source silicon for customers who owned their GDS.

# The Complex Maze of Robust Digital Design and Timing Sign-Off

# CMOS SPICE Model Evolution and Methodology Changes for Digital Timing Verification

- Very early MOSFETs had very few parameters:
  - Width (w) and length (l)

- Later, area and perimeter of source and drain junction region were added to the models
  - AS, AD, PS, PD

- When strain was introduced, more parameters were needed to describe the environment of the MOSFET.
  - OD extension from gate (strain)
  - Distance to edge of well (proximity effect)

- A 16nm FinFET instance has >50 parameters
  - Think about the complexity of making an accurate model with 50 parameters when you can only measure individual MOSFETs in silicon which include significant local manufacturing variations.

Planar MOSFET



FinFET



14-nm FinFET by UMC

14-nm FinFET by Intel

# Local Variation Emerges as a Major Concern in Digital Timing at 65nm



- MOSFETs have several parameters (physical, and electrical) that determine performance.

- Manufacturing imperfections create global variations (chip to chip, wafer to wafer), as well as uncorrelated local variations (between neighboring transistors).

- As MOSFETs were scaled down, the total process variation increased due to uncontrollable and uncorrelated random dopant fluctuations and line edge roughness of the gate.

# Total Corner SPICE Model Pessimism

- When local variation was small, total corner models only included a small amount of pessimism.



- With growing local variation, the pessimism became significant.



- Total corner model[1]:
  "A 3-sigma transistor of all manufactured transistors."

$$V_{T,SS} = V_{T,nom} + 3 \cdot \sqrt{\Delta V^2_{T,\sigma_{global}} + \Delta V^2_{T,\sigma_{local}}}$$

$$V_{T,FF} = V_{T,nom} - 3 \cdot \sqrt{\Delta V^2_{T,\sigma_{global}} + \Delta V^2_{T,\sigma_{local}}}$$

- Global Corner model:
  "The average transistor on a 3-sigma wafer."

$$V_{T,SSG} = V_{T,nom} + 3\ \Delta V^2_{T,\sigma_{global}}$$

$$V_{T,FFG} = V_{T,nom} - 3\ \Delta V^2_{T,\sigma_{global}}$$

1) The statistical distributions for total corner models can be measured directly.

2) Global corner parameters require the statistical subtraction of local variation.

# Modeling of Local Variation on Digital Timing Paths

$\tau_1, \sigma_1$  $\tau_2, \sigma_2$  $\tau_3, \sigma_3$

$\tau_4, \sigma_4$

$\tau_5, \sigma_5$

$\tau_6, \sigma_6$

- Local variation adds uncorrelated random variation to the nominal gate delay at the respective modeling corner.

pdf

$3\sigma_i$

delay

$\tau_i$

- Nominal path delay: $\tau_{path} = \sum_{i,path} \tau_i$

- Because local variation is uncorrelated, it accumulates statistically over a timing path:

$$\boldsymbol{\sigma_{path}} = \sqrt{\sum_{i,path} \sigma_i^2}$$

- Effectively, delay grows ~linear with path length $\boldsymbol{n_{stages}}$, whereas variability grows ~ $\sqrt{\boldsymbol{n_{stages}}}$

  - Shorter critical paths require more relative margin!

Inphi

# Evolution of Total Corner SPICE Model Optimism and Pessimism

sigma_local=0.4*sigma_global

sigma_local=sigma_global

# The Complex Maze of Robust Digital Design and Timing Sign-Off



## Evolution in Local Variation Models

# The Dawn of Statistical STA at 65nm

- Designers and EDA houses realized that the total corner pessimism became excessive.
  - ARM cores with flat OCV timing closed at 1GHz using total corner SPICE models and 3-sigma interconnect models ran at 1.3GHz in the lab!

- The first generation of SSTA tools appeared
  - EinsTimer Statistical (IBM ~2004)
  - Goldtime (Extreme DA, founded in 2003 to create a SSTA tool, acquired by Synopsys in 2011).
  - PrimeTime VX (Synopsys, 2006), Tempus (Cadence, 2013)

- Local variation was removed from the SPICE corner models for the timing libraries
  - => Global corner models (SSG and FFG).

- Instead, margin against local variation was created inside the SSTA margin model in the STA tool.

- This resulted in less overdesign for long paths, and less risk of underdesign for short paths.
  - This methodology change was especially important for hold margin on short paths that cannot be embedded in the SPICE model!

- Creating high quality statistical libraries was a significant hurdle.
  - The EDA stakeholders climbed statistical margining on an incremental path:
    AOCV | POCV => LVF

# Innovation in Sign-Off Margin Modeling for Local Variation



Established technology for over 40/65nm → Mainstay technology for 40/28/20nm → Emerging technology for sub 20nm

**OCV** — Flat global margin — Extra pessimism

**AOCV** — Depth based margin — Reduced pessimism

**POCV** — Statistical delay distribution — Least pessimism

Flat derate

Derate =f(path length, spatial spread)
IMHO "poor man's" statistical model

=> LVF

# AOCV, POCV, and Liberty Variation Format (LVF)

- Advanced OCV (Or stage-based OCV): Scales margin down for longer paths by counting gates on the path. This was poor-man's statistical analysis as it did not account for the absolute magnitude of the variation.
  - AOCV analysis opens a significant gap between graph-based analysis and path-based analysis.
- Parametric OCV (POCV): Extreme-DA's simplified statistical margin model in their STA tool "Goldtime". It reduced the significant long/short path merging pessimism in AOCV.
  - POCV accumulates $\sigma^2$ variation over the timing graph, leading to a smaller gap between path-based and graph-based analysis. The variation amount for each gate is assumed proportional to the gate delay.
- Liberty Variation Format: (Expensive library characterization)
  - Absolute timing variation is captured in Slew x Load tables similarly to the NLDM data. Different data for early and late variability is supported to account for skewness of the statistical distribution.
    - $\sigma_{late} > \sigma_{early}$

# Why was it important to innovate to absolute variation data in LVF? Derating as a margin mechanism has a hole!



- Thought experiment:
  A small gate with a small output load receives an input signal with a large RC delay. The input transition time is very large. The trip point (Vout=Vin) of the gate may be off-center so that the nominal delay is close to zero. We could even see negative delays for the small gate.

- When local threshold voltage variation shifts the trip point up or down, the point in time where the gate output switches moves significantly to early or late. The variation grows roughly proportional to the input transition time around the trip point. => Slow transitions create timing variation hot spots.

- Margin mechanisms that are based on derating of nominal delays (like OCV, AOCV, or POCV) create no meaningful margin for such a hot spot if the delay is close to zero![1] LVF fixes this hole. ☺

1) This is the reason why derating of constraints is not a reliable margin mechanism. Constraints are differences between longer paths. Margin needs to cover uncorrelated variation on the sum of the paths.

# Slew Dependence of Constraint Variation, Tau 2014

**variation hotspot**

## Hold time met?

$t_{rf,D}$

$t_{r,CK}$

D

CK

State of the art:

$$\tau_{hold} = f(t_{r,CK}, t_{rf,D})$$

Not modeled at that time:

$$\sigma_{hold} = f(t_{r,CK}, t_{rf,D})$$

- Hold time variation also shows a significant increase for slow data or clock transition times.[1]
  - Slew-dependent constraint variation was later added to the Liberty library data (~2016).
- By modeling slew-dependent hold time constraint variation explicitly, this concern can be removed from the blanket uncertainty budget.
  - => Constant margin can be reduced.
  - However, there are still good reasons to maintain a backstop of constant margin!
    - Silicon to SPICE gap and other modeling errors.

- Modeling this variation also provides an incentive for the design tools and designers to maintain good signal transition times at flip-flop inputs that are essential in avoiding variation hot spots.

1) Christian Lütkemeyer and Praveen Ghanta (Broadcom): "Modeling slew dependent constraint arc variation in Static Timing Analysis". Tau Workshop 2014.

Inphi

# The Abstraction Layer Error Pile in Digital Circuit Timing Verification

**STA timing results**

All slacks ≥0.0? ☺. Path with slack ≤0.0? ☹. It is not that simple.
**We can pass all timing checks and have non-robust silicon.**

**STA Tool**

**SPEF Interconnect**

Different STA tools will almost never produce the same slack values. (Assumed input waveform shapes, interpolation in load/slew tables and CCS waveforms, SI modeling, Advanced Waveform Processing). It is important to be able to distinguish significant differences from analysis noise!

**.lib model**

**Sign-off recipe**

The Foundry provides most components for a sign-off recipe, but there are significant holes (MIS, layer-to-layer interconnect variation).

SPICE to .lib model errors (simplified model of non-linear input capacitance, discretization of loads, slews and CCS output current waveforms).

**SPICE model**

Silicon to SPICE model error. The FinFET SPICE Models have ~50 parameters. Making an accurate model with this complexity from silicon measurements is a challenge!
**And the process may shift vs. the early models if you are designing at the bleeding edge!**

**Early adopter risk**

Inphi

# Resolution, Accuracy, and Pessimism in STA



- **Resolution:**
Slack numbers[1] in timing checks are typically reported with a resolution of 0.1ps. The EDA tools fix timing until slack is zero or positive.
  - Slack<0.0: ☹
  - Slack >=0.0: ☺



pdf     A     B     A<B? 90%     C     A<C? 100%

- **Accuracy:**
Due to the abstraction from SPICE to the .lib library models an error of +/- 2% is quite good.
  - Designers need to consider the margin in the analysis, and the error due to accuracy limitations, when they react to timing "violations". A "small" 10ps violation could be an indication of significant yield fallout, or just noise in the analysis. The timing and margin context matters.

- **Pessimism:**
STA tools are pessimistic in their treatment of signal coupling. They only track the early and late arrival time of aggressor signals. Depending on the topology of the circuit, an aggressor signal might not even switch close to the transition of a critical timing signal.



early     aggressor  switches here?     late

victim receives SI penalty when its transition falls into the early-late aggressor window

time

23

1) difference between latest arrival, and latest arrival that meets timing for a timing check.

Inphi

# STA Modeling Error Characteristics

- The modeling error vs. SPICE is a function of
  - Cell types: (Single stage gates show more error than multi-stage gates).
  - Supply voltage: (timing data at lower supply voltages is less accurate).
  - Signal transition times: Slow transition times create large sensitivity to waveform tail discrepancies.
  - Input signal transition shapes: The signal transition shape in a circuit deviates from the characterization waveform due to differences in the driving gate (VTs, stack differences), RC delays, coupling, kickback from side-loads (back-Miller coupling).
  - The specific circuit.



- When waiving negative slack violations, designers need to take the accuracy of STA and the overall margins in the sign-off recipe into account.
  - Setup problems go away for faster silicon, or increased voltage. Hold problems can strike anywhere.

# Additional Important Verification Aspects

- **Signal integrity**
  - No glitches on asynchronous sets/resets?
    - Does the glitch analysis cover for local variation tails?
- **Completeness of timing constraints**
  - Only what gets checked is verified.
  - If critical paths get masked with a set_false_path, trouble is waiting.
- **Simplicity increases confidence**
  - The more complicated a model is, the harder it is to validate all aspects of it.
  - Tool bugs, human errors, …

# The Complex Maze of Robust Digital Design and Timing Sign-Off



Moore's Law · Process · STA .lib Models · Aging · Supply noise · EDA Models · Interconnect · Measured Variation · LVF · DCC Model · Measurements · Cargo cult engineering · Cargo cult innovation · Variation · SPICE · Coupling · Cap-scaling Model · MIS · MIS derate margin · Waste · $$ Profits · CMOS NAND gate · DLY250 · CKDLY250

Inphi

# A Story About Strained Silicon and Cargo Cult Engineering

calteches.library.caltech.edu/id/eprint/51/2/CargoCult.pdf

## Cargo Cult Science

by RICHARD P. FEYNMAN

Some remarks on science, pseudoscience, and learning how to not fool yourself. Caltech's 1974 commencement address.

During the Middle Ages there were all kinds of crazy ideas, such as that a piece of rhinoceros horn would increase... But even today I meet lots of people who sooner or later get me into a conversation about UFO's, or astrology, or some form of mysticism, expanded consciousness, new types of awareness, ESP, and so forth. And I've concluded that it's *not* a scientific world.

I think the educational and psychological studies I mentioned are examples of what I would like to call Cargo Cult Science. In the South Seas there is a Cargo Cult of people. During the war they saw airplanes land with lots of good materials, and they want the same thing to happen now. So they've arranged to make things like runways, to put fires along the sides of the runways, to make a wooden hut for a man to sit in, with two wooden pieces on his head like headphones and bars of bamboo sticking out like antennas —he's the controller—and they wait for the airplanes to land. They're doing everything right. The form is perfect. It looks exactly the way it looked before. But it doesn't work. No airplanes land. So I call these things Cargo Cult Science, because they follow all the apparent precepts and forms of scientific investigation, but they're missing something essential, because the planes don't land.

# The Story About Cargo Cult Science from WWII

The US were taking back the islands in the Pacific from Japanese occupation. Soldiers arrived on amphibian ships, created a landing strip, built a tower, put up antennas, placed a man with flags on the runway, and planes with cargo started to arrive. A while later they left for the next island...

The local people on the islands tried to recreate everything perfectly later, but cargo did not arrive...

# Strain Engineering and Proximity Error (40nm Node)

NMOS

PMOS


Stress layer

INV  

CKINV  


8F
20F


Strain defining endcap
endcap

- Strain increases $I_{D,sat}$ in MOSFETs
- Strain depends on the layout adjacent to the MOSFET, i.e. varying neighbor cells in a standard cell design.


perimeter | DUT

- This "proximity effect" variability is not captured in standard cell .lib models. Cells are only characterized in a "typical" neighborhood (NWELL, OD).
- To reduce the strain variability for clock inverters or buffers, regular inverters/buffers received "endcap" cells that provided a defined OD neighbor pattern to the active devices.

# Clock Insertion Delay Balancing

20ps / CKINV

2ns

- Balancing clock insertion delays in the nanosecond range with CKINV cells required long chains of tens of clock inverters.

- For data signals the library provided special delay cells ("DLY*") to fix hold violations. These implemented much larger delays in a significantly smaller footprint.

DLY250

Cargo cult innovation

Strain defining endcap | CKDLY250 | Strain defining endcap

Suddenly, "CKDLY" cells appeared in an in-house 40nm library update.
I stumbled on them by accident when I reviewed a path report for sign-off waiver questions.

Inphi

# Schematic View: INV vs. DLY



- INVD4
  - Parallel devices to minimize output resistance.
  - Optimized to switch as fast as possible.

- DLY
  - Stacked long channel devices maximize R.
  - Inputs tied together to maximize C.
  - => Maximizes RC delay per unit area with a maximum stack height of 3.



Red + 4 wheels!
But do they match?

# The Cargo Cult Story Behind the "CKDLY" Cells



- A creative manager in a product team had the idea to create the CKDLY cells by adding the strain-controlling endcaps. He knew about the strain impact from the neighboring cells…

- The helpful library team supported their customer by providing the layouts and timing views, without consulting the timing sign-off experts. Their goal: "Keep the customer happy!"

- The cargo cult problem:
<span style="background-color:red;color:yellow">Matching between the CKINV cells and the long channel devices in the DLY cells is poor. There may also be a large silicon to SPICE model gap for these devices.</span>

- What the library team gives, the sign-off team had to take away!
We must follow robust design principles!

  "Only the paranoid survive!"[1]

1) Guiding principle of Andy Grove, 3rd employee (COO) at Intel, later CEO.

# The Complex Maze of Robust Digital Design and Timing Sign-Off



## The Multi-Input Switching Modeling Hole

# Multi-Input Switching

A1
A2
A3
A4 — ZN

A1
A2 — C
A3
A4

A1
A2, A3, A4
ZN

$\tau_{MIS}$

significant speed-up => hold risk

4  3  2  1

For N falling inputs:

$$\tau \sim \frac{C \cdot VDD/2}{I_{sat}} \sim \frac{1}{N}$$

| # inputs falling: N | MIS derate: $R_{rise}$ |
|---|---|
| 1 | 1 |
| 2 | 1/2 |
| 3 | 1/3 |
| 4 | 1/4 |

1)

The faster Multi-Input Switching events (#2 to #4) are not captured in the Liberty model of the library. **UNMODELED HOLD RISK!**

1) The "Mystery Bag & Blindfold Set is a good learning tool for the Montessori classroom. When we make $10M mask sets, mysteries should be avoided.

34

Inphi

# Example of a Potential MIS Failure Scenario



"slow NAND4" delay (SPICE):
SIS: 90.5ps
MIS: 29.6ps

Multi-bit FF
STA: Hold time met!

Silicon: Violation! (-61ps)

**How do we know this circuit does not exist in our 100M+ gate design?**

clk

4-bit counter wrap-around (1'b1111 to 1'b0000). Zero detect fires early.

- All common clock tree => no margin from derating of divergent clocks.
- Capture clock delayed internally in the MB FF (higher internal clock load).
- If the hold check at the multi-bit Flip-Flop meets timing with zero slack in STA when MIS is ignored, then MIS speed-up of the NAND4 may cause significant hold time violations in silicon.
- This example with -61ps MIS slack shift shows that fixing the MIS hole with constant uncertainty would be very expensive!
  It may have been ok when we had 100ps uncertainty in the simple 65nm OCV sign-off.

# History of Multi-Input Switching Margining (My View)

- Broadcom started to use aggressive derating to cover MIS risk at the 40nm node when advanced timing sign-off was established. In-house, SPICE simulation-based, characterization was developed to provide the worst-case derating data. No EDA characterization solutions were available.

- In 2014, I developed a more refined derating model that could reduce pessimism by leveraging the information of the input timing arrival windows on the derating values.[1]
  This model was adopted by Cadence in their Tempus STA tool. The model also prompted fixes in the PrimeTime MIS margin model that had been developed by Extreme DA per Broadcom's request.

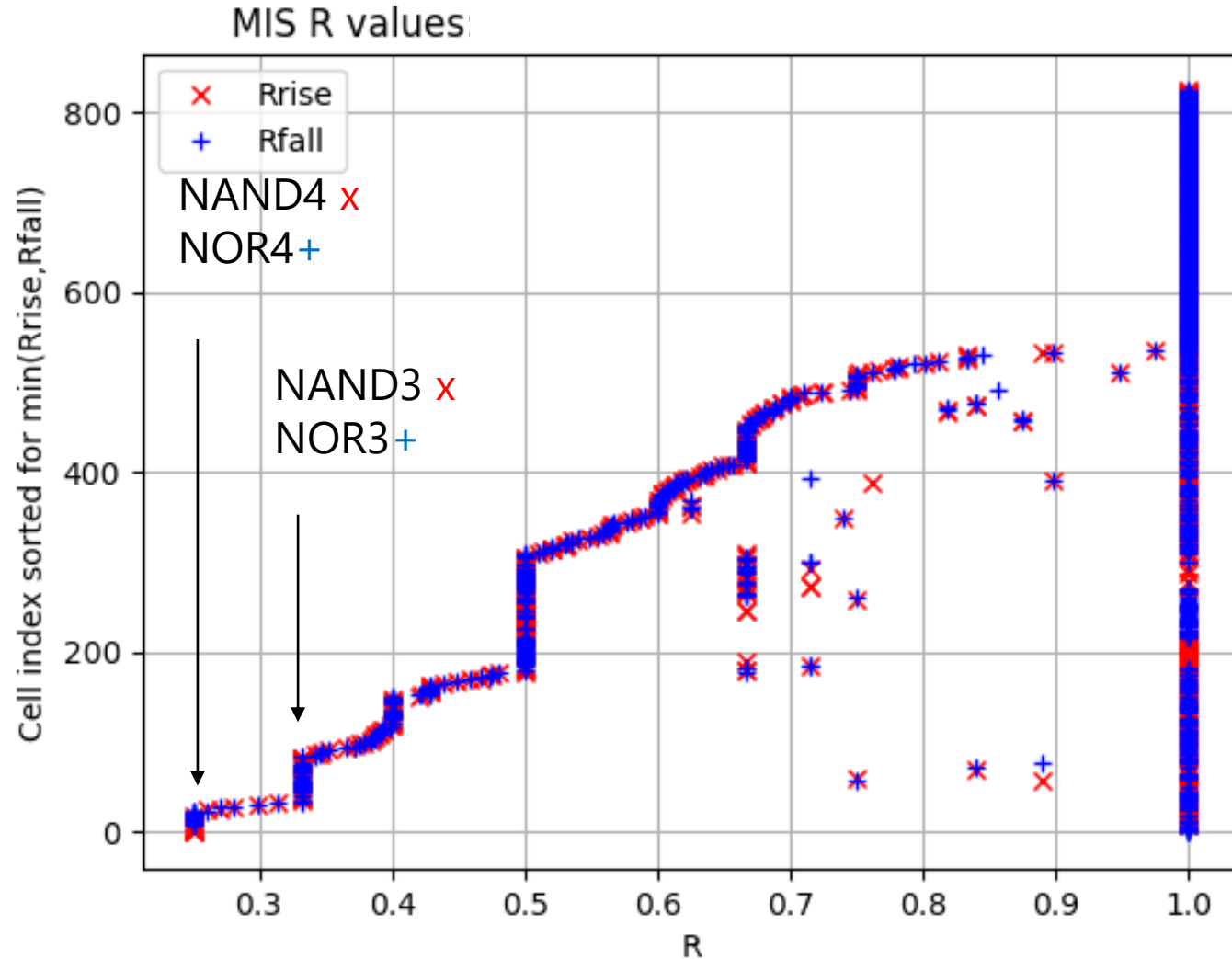- Synopsys is working on a solution that does not require characterization. (See Tau 2019 MIS Panel) However, currently this solution does not support gates where some inputs have an added input inversion. Such "Swiss cheese solutions"[2] still leave holes, just fewer.

- Until a no-hole solution is available commercially I plan to rely on the derating mechanisms in PT and Tempus, and in-house characterization of conservative derating values.

1) Christian Lütkemeyer (Broadcom): "A practical model to reduce margin pessimism for Multi-Input Switching in Static Timing Analysis". Tau Workshop 2015. http://www.tauworkshop.com/2015/slides/Lutkemeyer_TAU15_PPT.pdf

2) The holes in Swiss cheese are vanishing as they were created by tiny hay particles in the milk collection buckets.

# Example of Multi-Input Switching Derating Data



MIS R values

- More than 500 of the 800+ standard cells in a library show significant MIS speed-up.

- The MIS derating data is created from an analysis of the MOSFET sizes and topology in the cdl netlists of the standard cells. It is much faster than a simulation-based characterization.
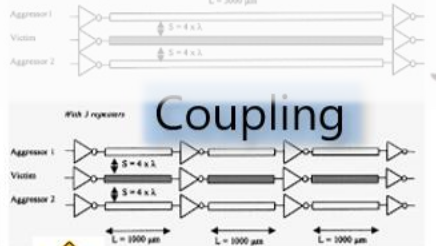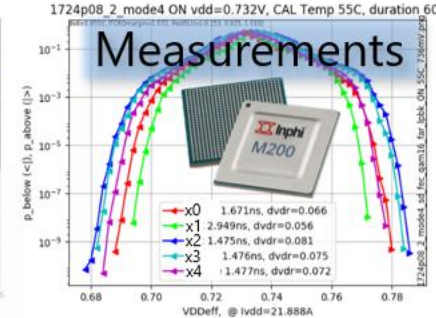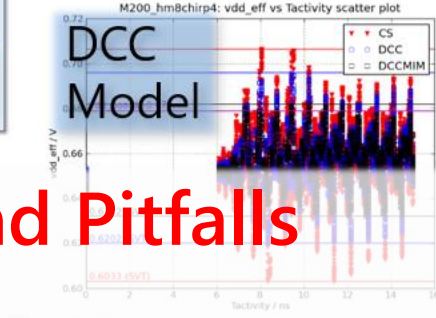
# MIS Fixing Cost Example from M200


M200
16 nm
LightSpeed III™

- M200 closed timing with moderate fixed uncertainty w/o MIS modeling until close to tape-out.

- When MIS modeling became available, we allowed a significant reduction of the constant margin since the MIS risk was now explicitly modeled.

- MIS derating identified about 300 risk locations in total. They were fixed in simple ECOs.

- <mark>In the context of a complex chip the fixing cost was negligible.</mark>

- We have heard from a customer that they had to fix a chip after they had run into a MIS problem.

  - Think about the fun of identifying functional failures in your complex chip when your sign-off analysis indicates your chip is supposed to be robust! Debug, ECOs, tape-out... Plus the delay to revenue...

- <mark>It is mind-boggling that the STA MIS hole, known for more than twenty years, is still not comprehensively covered in the EDA and Foundry modeling world.</mark>

# The Complex Maze of Robust Digital Design and Timing Sign-Off



**Interconnect Modeling and Pitfalls**

# Fully Correlated Interconnect Modeling Corners

**Inter layer dielectric thickness bias**

$\sigma_C$

$t_{i,nom}$

$t_{nom}$

$\sigma_M$

$w_{nom}$

$t_{i,nom}$

$\sigma_C$

Cbest

RCbest

$3\sigma_C$

Cbest_T ✕

RCbest_T

✕ Ctyp

**Wire cross section bias**

$3\sigma_M$

RCworst_T ✕

✕ Cworst_T

RCworst

Cworst

For setup modeling, *_T corners are 1.5 sigma models that reduce interconnect pessimism. In older technologies, designers would combine SS total corner models with 3 sigma interconnect.

The interconnect corner models suffer from the unrealistic assumption that the variations on all layers point in the same direction.[1]

1) CMP is a force vector that will create a thinner layer if the previous layer was thicker => reverse correlation is more likely.

⬭ **Inphi**

# Complexity of Interconnect Modeling

- Each layer has at least four physical variation parameters
  - Width, thickness, inter metal dielectric thickness, and via resistance.
    => 16 possible corners per layer; 4 dominant corners.
- For 15 metal layers this would result in $4^{15}$ = 1,073,741,824 dominant corner combinations.
  - Clearly it is not possible to provide corner models comparable to the established Cworst, Cbest, RCworst, RCbest scheme for so many combinations.
  - It is also impossible to manufacture meaningful metal split lots for so many layers.
- As a result, chip designs that are created with insufficient hold margin against interconnect variation are likely to see occasional yield collapse as the random manufacturing variations expose paths with marginal hold robustness.
  - Having a few split lot wafers with good yield early on gives only limited confidence that a design is robust to future metal manufacturing surprises.

$\sigma_c$

$t_{i,nom}$

$R_{via}$

$\sigma_{Mt}$

$\sigma_{Mw}$

$t_{nom}$

$w_{nom}$

Inphi

# A Clock Balancing Example Over Different Metal Layers



**Launch path Cycle n (early)**

Metal9/10 -||- ↓

clk

**Capture path Cycle n (late)**

Metal11/12 -||- ↑

**Unexposed hold risk**

- When clock branches are balanced over different metal layers the layer-to-layer variation can cause hold failures that are not revealed by analysis with fully correlated models.[1]

**Hold time met?**

☹ Data arrives before clock. ☹

1) Fully correlated interconnect models are pessimistic when metal loads are matched against gate loads. If the model is accurate!

# Interconnect Ring Oscillator (ITCRO)

Tau Workshop 2016: Layer-to-layer interconnect variation is a significant but unmodeled source of hold time optimism in conventional BEOL corner
Christian Lutkemeyer, Ali Anvar, Broadcom; http://www.tauworkshop.com/2016/slides/1_TAU_2016_Christian_LTLIV_SBUS-ETP101-Rshort.pdf

- **Capacitance change => change in period.**
- **From period changes vs. "unloaded" cases we can calculate capacitance ratios with ~1% accuracy.**



sl[1]

not(sl[1])  test load 1

sl[0]

not(sl[0])  test load 0

sl[1]

not(sl[33])

sl[0]

not(sl[32])

run  n0  n3  n6  n9  n13  o

weak inverter

# 16FF+ Data: Interconnect Capacitance Ratio-of-Ratios

Ratio of Ratios vs. layer M7



- Test structures:
  - 1w (single width, single spaced)
  - 2w (double width, double spaced)
  - p (plate wider than 2w)

- We see a significant increase for the M2_2w / M7_2w, M2_p / M7_p, M3_2w / M7_2w, M3_p / M7_p ratios.
  - We suspect that the extracted capacitance values are missing capacitances to the bottom (inside the standard cells). The "gray" box extraction may not work properly!
    This is consistent with more impact on M2 vs. M3 as M2 is closer to the cells.

# 7nm Data: Interconnect Capacitance Ratio-of-Ratios

Ratio of ratios normalized to M3



Set A   Set B

- 2 Similar experiments, Set A and Set B close to each other

- M2_1w/M3_1w trends high

- M2_2w/M3_2w trends high

- M3_2w/M3_1w trends high

- Mx with x>3 ratios trend low, consistent with a modeling error that increases M3 capacitance on the chip vs. the model.

- Suspicion that the "gray" box extraction may miss coupling to the inside of the cells.

- Layer-to-layer ratio band shows about 20% variability.

**M2 and M3 show significantly increased "2w" capacitances. The "gray" box extraction may still be off.**

# The Sensitivity SPEF (Standard Parasitic Exchange Format) in IEEE Standard 1481-2009

- This standard included mechanisms to parameterize the parasitic elements in the SPEF file with sensitivity parameters for individual metal layers.

- SSTA tools were supposed to find the worst-case variation parameter set for each timing check.

- TSMC developed support once for the 40nm node.

- SSTA tools were not ready, the technology did not get used, and the industry fell back to the fully correlated interconnect corner models.

- I attempted to resuscitate support for the SSPEF model at the 2016 Tau Workshop.[1]

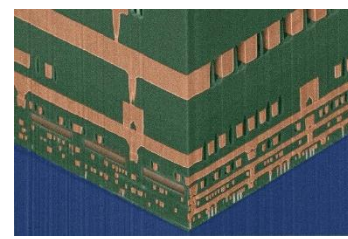  - Sharing data with >10% capacitance ratio variation in a mature 28nm CMOS process was not scary enough to motivate change. ☹

SSPEF RIP

- The SSPEF was carried forward into the 1481-2019 Standard, but still there is no actual support in the EDA and Foundry World that I live in.

1) Christian Lütkemeyer and Ali Anvar (Broadcom): "Layer-to-layer interconnect variation is a significant but unmodeled source of hold time optimism in conventional BEOL corner models", http://www.tauworkshop.com/2016/slides/1_TAU_2016_Christian_LTLIV_SBUS-ETP101-Rshort.pdf

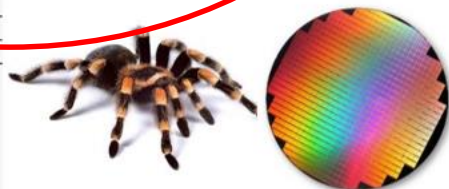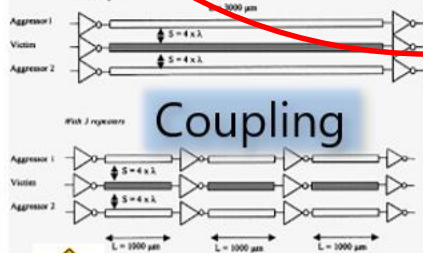# Capacitance Scaling to Cover Interconnect Layer-to-Layer Variation

Cap-scaling

| Corner | Scale early | Scale late |
|--------|-------------|------------|
| Cworst (C↑) | 0.8 | 1.0 |
| Rcbest (C↑) | 0.8 | 1.0 |
| RCworst (C↓) | 1.0 | 1.2 |
| Cbest (C↓) | 1.0 | 1.2 |

- Some EDA tools enable independent scaling of the capacitance from the SPEF file for early and late analysis.

- Scaling only the capacitance without adjusting the reverse correlated resistance is a non-physical, but pessimistic margin mechanism that creates margin against layer-to-layer variation.

- Gate-load dominated paths receive less pessimism.

- By cap-scaling, designs are hardened against variations on delays that are matched over different metal layers.

# The Complex Maze of Robust Digital Design and Timing Sign-Off



**Chip Power Integrity Modeling**

# The Complexities of the Chip Power Delivery System

# The Power Integrity Modeling Complexity Challenge



- How does SPICE work?
  - The SPICE engine analyzes circuits based on the Kirchhoff's current or voltage law. It assigns nodes to a circuit and attempts to solve the current and voltage values at the respective nodes. The SPICE simulator first generates nodal equations in the matrix format before solving them to obtain the values.[1]

- Today's SOCs combine Millions to even Billions of gates in a single chip.

- The switching current of each gate is a function of its local voltage.

- Since all the gates are connected to the same supply rail, the supply voltage and current of each cell depends on all the other cells in the circuit, and their activity.

- SPICE would have to solve an enormously complex matrix, ultimately connecting all gates with all others, and iteratively solve voltage and current equations to find an accurate solution for the time varying voltages.
  => This problem cannot be solved with SPICE in the computers we have access to!

1) https://resources.pcb.cadence.com/blog/2019-how-does-spice-simulation-work

# A Cargo Cult Engineering Example from the "Flat Earth" Power Integrity World

- I complained about the modeling problem on the following page to (then) Apache (RedHawk) many years ago.

- The problem is related to how the capacitive cell load is connected when the demand current of a standard cell is characterized.

  - For Electromigration (EM) modeling the load must be lumped to either VDD or VSS to show the full current on the local rail.

  - For power integrity modeling where the chip demand current is simulated, the loads should be split 50/50 between VDD and VSS.

**Electromigration**

**Chip Level Power Integrity**

# CMOS Demand Current Characterization with Lumped-VSS Load



- The demand current $i_{VDD}$ is characterized for all standard cells with a <u>lumped capacitive load from "o" to "VSS"</u>.

- With this model, we observe a large demand current when the output is rising, and a small demand current when the output is falling.

- For the clock tree, this results in a <mark>large imbalance of the demand current between rising and falling transitions</mark> at the end points.
  => <u>Significant supply noise at the clock frequency.</u>

- Some engineers asked: "How can we reduce this imbalance to reduce the supply noise?"

# Power Integrity Cargo Cult Engineering Example: Minimizing peak current via opposite-phase clock tree.



Figure 4. Two-level binary clock tree.



Figure 5. Our approach (using opposite phases).

▪ There are several papers which propose to use negative edge triggered flip-flops for half of the registers, and clock inversion for those registers, to reduce the demand current imbalance between the rising and falling edge of the clock tree.

=> <u>Goal: Reduction of the supply noise amplitude, and shift of the energy to twice the clock frequency.</u>

• Registers remain clocked at the same time.

• Nieh, Yow-Tyng, Shih-Hsu Huang, and Sheng-Yu Hsu. "Minimizing peak current via opposite-phase clock tree." Proceedings of the 42nd annual Design Automation Conference. ACM, 2005.

• Y. Ryu and T. Kim, "Clock buffer polarity assignment combined with clock tree generation for power/ground noise minimization," in Proc. IEEE/ACM Int. Conf. Comput.-Aided Design, 2008, pp. 416–419.

# A Lot of Work to Fix a Non-Existent Problem



VSScratio=1.0

CMOS layout symmetry

VSScratio=0.5

$i_{VDD}$

$C_{load}$

VSS

"o" can couple to:
Well(VDD) | Substrate(VSS),
Power nets VDD | VSS,
Victim nets @ VDD | VSS

$i_{VDD}$

$C_{load}/2$

$C_{load}/2$

VSS

v(o)

$i_{VDD}$

**The imbalance is just a modeling artifact![1)]**

v(o)

$i_{VDD}$

t

t

54

1) Shielding clock wires with VSS-only can create this imbalance. It is better to use VDD | VSS shields!

Inphi

# Supply Resistance of a Switching CMOS Inverter



- The average supply current increases ~linear with VDD:

$$i_{VDD,avg} = f \cdot C \cdot VDD = \frac{VDD}{R_{inv,avg}}$$

- A switching CMOS inverter appears as an average resistance of $R_{inv,avg} = \frac{1}{fC}$ on the power grid.

  If the inverter does not switch, the inverter and its output load appear as a capacitor with a small parasitic series resistor.

  - The RC time constant is approximately the switching delay of the gate.
  - This is in the range of ps to 10s of ps, i.e. much shorter than the resonant period of the power delivery system (~100MHz for a FCBGA).

- Switching CMOS circuits are non-linear and time varying. Their average current over a clock cycle behaves like a resistor with finite resistance.

# EDA Solution to Handle Power Integrity for 100M+ Instances: Clobber Resistive CMOS Logic Into Current Sources

"Resistive" CMOS Circuit

Current Source



- $R_{CMOS} = \dfrac{1}{f \cdot C}$

- $I_{CMOS} = \dfrac{VDD}{R_{CMOS}}$

- $R = \infty$

- $I_{CS} \neq f(VDD)$

**Current source modeling eliminates the need to iterate to find** $vdd(t) = \sum i(t) \cdot Z$

# Experiment to Show the Beneficial Damping of Switching Logic on a Single Supply Plane in a CMOS SOC With Load Steps

VDDboard=0.6V, Lpkg=0.1nH, Rpkg=0.25mΩ, C=60nF [1]



Resistive Model

Current Source Model

1) Parameters chosen to highlight the CS-problem.

# Load Step Response

Resistive Model



Current Source Model



- The load step response amplitude is smaller. Oscillations are dampened as the system has a lower Q-factor.

- The load step amplitude is increased, and the high Q-factor of the RLC power delivery system leads to a very long settling time.

# State of the Art in Commercial SOC Power Integrity Modeling Tools

- High Spatial Resolution (100s of M instances).

- Pessimistic current source modeling for the cells clobbers the resistive essence of switched capacitor CMOS circuits.

- Circuit activity can be modeled with
  - vectors (VCD) or
  - vector-less (switching probabilities are propagated).
  - Only a very short real-time window can be simulated, i.e. just a few clock cycles.



- Since modeling is so unreliable and severely limited, how can we know how much digital supply noise may appear in reality?

# Example of SOC Power Integrity

VDDeffective measurement data (M200)



**Each data point represents the probability that the voltage will be below (◄) or above (►) VDDeff, measured over 60s of operation.**

- Dynamic Voltage Drop Ratio (dvdr) is the percentage drop below the average supply at the lowest valley of the supply ripple, averaged over full clock cycles.

- To observe worst-case supply voltage events in a complex SOC, measurements need to be done over Billions of clock cycles.

- For power integrity modeling work, this means that it is essentially impossible to find the worst-case time window of a few nanoseconds that we can simulate in vector data, even if we could generate and process activity data that represents seconds of real-world operation!

- Is modeling data good enough for back-annotation in STA? Use at your own risk!

# A 7nm LVF Signoff Flow

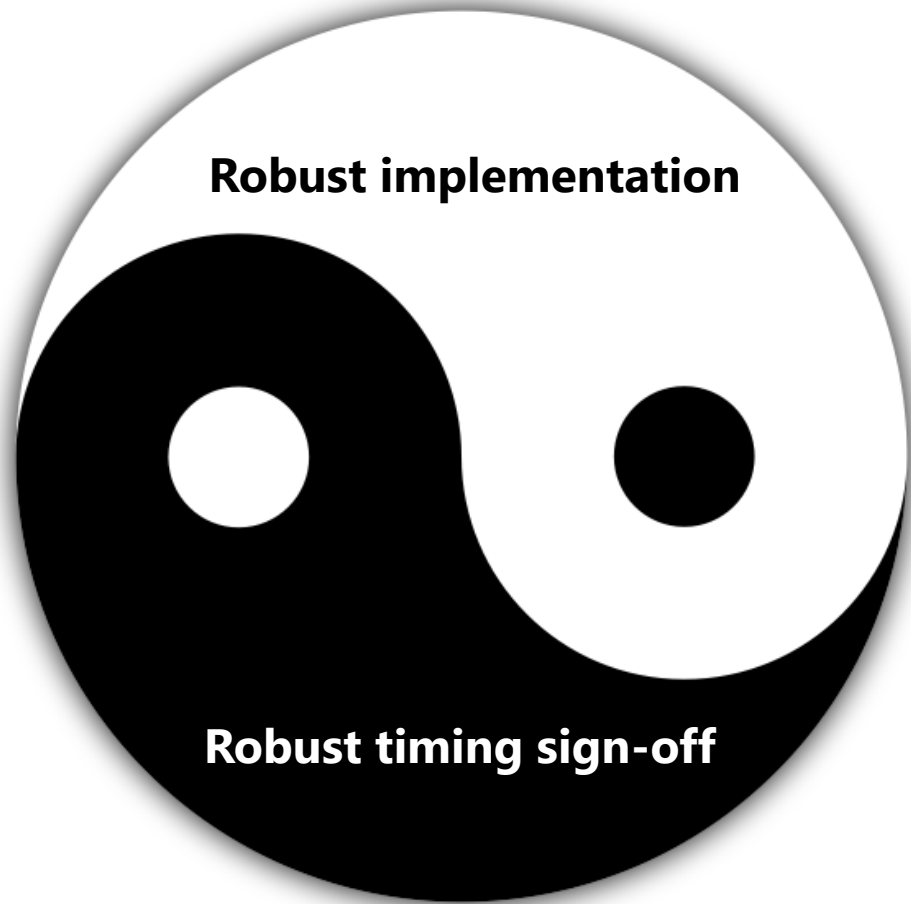| Variation parameter | Range | | Scenario multiplier |
|---|---|---|---|
| Process | SSG, FFG (global corner models) | | 2 |
| Operating voltage corners | $[VDD_{low}, VDD^{high}], [VDDM_{low}, VDDM^{high}]$ | | 4 |
| Operating temperature corners | 0°C, 125°C | | 2 |
| Interconnect corners | Cworst, RCbest, RCworst, Cbest (1.5 $\sigma$ for setup) | | 4 |
| Coupling | $\frac{1}{3}C,\qquad \frac{5}{3}C$ | | |
| Local variation | $-3\sigma, +3\sigma$ | LVF variation model | |
| Process gradients | +/- 10%/10mm | Diagonal of bBox derate | |
| Voltage gradients | +/-x mV | Derating or voltage scaling with differential early/late voltage | |
| Temperature gradients | +/-10K | | |
| Interconnect variation (layer-to-layer) | 0.9, 1.1 | **Capacitance scaling** | |
| Multi-Input Switching | 0.25, 1.1 | **MIS derating** | |
| Aging | 1.0, 1.x | Waiting for accurate statistical aging support. | |
| Modeling errors vs. SPICE | -3%, 3% | Additional derating and uncertainty | |
| Silicon to SPICE model gap | -3%, 3% | | |

$$\prod = 64 \text{ scenarios}$$
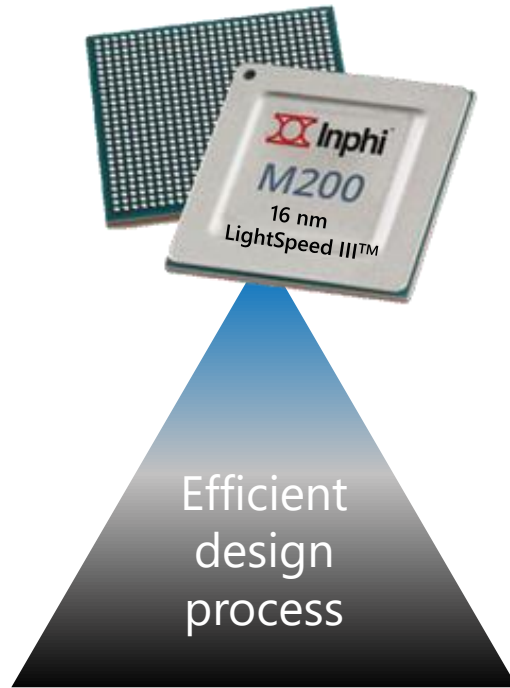
Up to 500 scenarios?

# Holistic Optimization of Implementation and Sign-Off



- Clock trees
  - Control routing layers to ensure matching.
  - Do you allow useful skew?
    - Matching between skewed clock tree and data path increases hold risk.
  - Clock meshes for high clock frequency designs.
    - Minimize non-common clock delays.
  - ...

- Low voltage designs
  - Vt choices for logic and clock tree.
  - Transition time limits for clock and data.
  - Avoidance of cells with high variability.

- If the margin cost is too high, the reduction of margins is not always the right answer!

- Maybe the design style is wrong for the problem at hand.

# Closing the Loop Over Sign-Off, Implementation, and Silicon
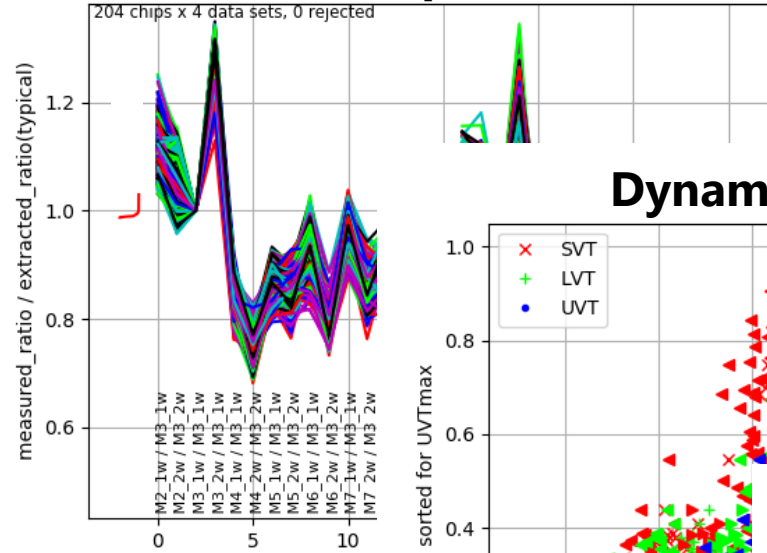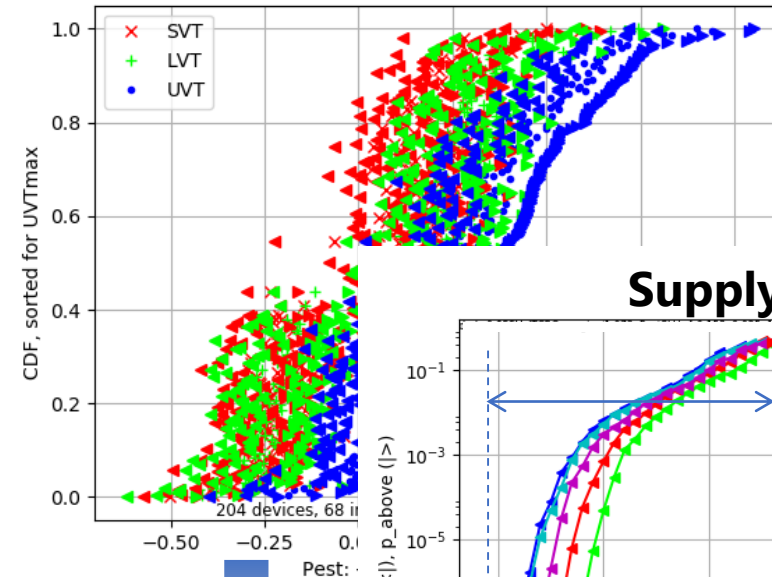
**Silicon Measurements in Products**

**16 nm LightSpeed III™**

**M200**

Efficient design process

**Robust timing sign-off**

**Robust design style**

**learn**

**Interconnect Capacitance Ratios**



204 chips x 4 data sets, 0 rejected

**Dynamic Performance**



SVT
LVT
UVT

**Supply Noise**



| | | |
|---|---|---|
| x1 | 1.671ns, | dvdr=0.066 |
| x2 | 2.949ns, | dvdr=0.056 |
| x3 | 1.475ns, | dvdr=0.081 |
| x4 | 1.476ns, | dvdr=0.075 |
| x5 | 1.477ns, | dvdr=0.072 |

VDDeff, @ Ivdd=21.888A

# Summary

- Improved modeling of local variation and other timing variation effects has enabled a significant reduction of pessimism for long paths and fixed uncertainty.

- Beware of pitfalls! Foundry and EDA support leave exposure.
  - Multi-Input Switching speed-up => Derating mechanisms exist. Data may have to be home-grown.
  - Interconnect layer-to-layer variation => Capacitance scaling can provide robust cover.

- State-of-the-art chip power integrity models clobber the beneficial damping of switching CMOS circuits.

- There are modeling artifacts when cells are characterized with lumped-VSS loads.
  - Don't waste your time to fix "problems" that do not exist!

- Back-annotating low confidence IR drop data into STA seems risky IMHO.

- Be realistic about overall modeling inaccuracies and the Silicon-to-SPICE gap.

- Know what you don't know!
  - Known unknowns!
  - Unknown unknowns!

- Avoid Cargo Cult Engineering!

Reasonable margins can safe your good-nights sleep!

**And your job!**

**Inphi**

# The Complex Maze of Robust Digital Design and Timing Sign-Off

**Thank You for flying with us today!**

Ops uncertainty

As we chip away the margins
in our chips
it is good to remember:

"Only when the tide goes out
do you discover who's
swimming naked."

Warren Buffet

Inphi