

Statistical Timing Methodology for Low-Power and Multi-Voltage Designs

Eric A. Foreman
Jeff Hemmett
Kerim Kalafala 
Steve Shuma
Natesan Venkateswaran
Vladimir Zolotov †

(eforeman@us.ibm.com)
(hemmett@us.ibm.com)
(kalafala@us.ibm.com)
(sshuma@us.ibm.com)
(natesan@us.ibm.com)
(zolotov@us.ibm.com)

IBM Systems
Essex Junction, VT and Fishkill, NY

IBM T.J. Watson Research †
Yorktown Heights, NY

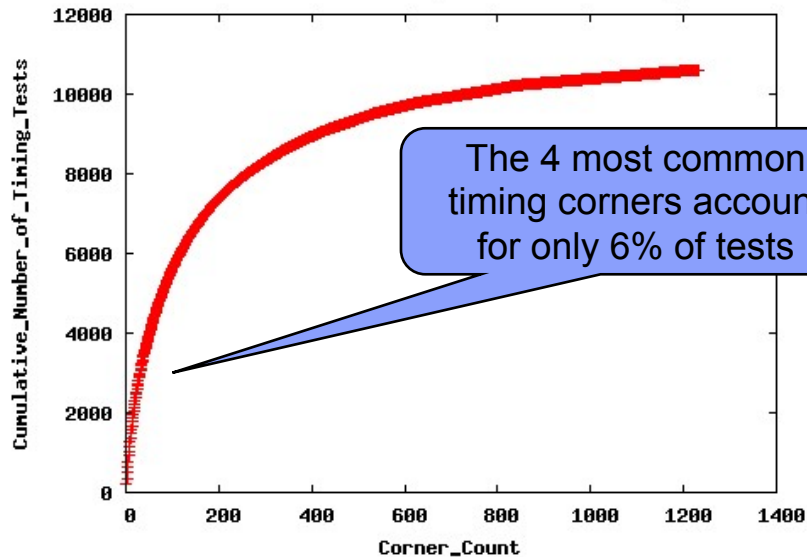


Member, IBM Academy of Technology

Introduction / Outline

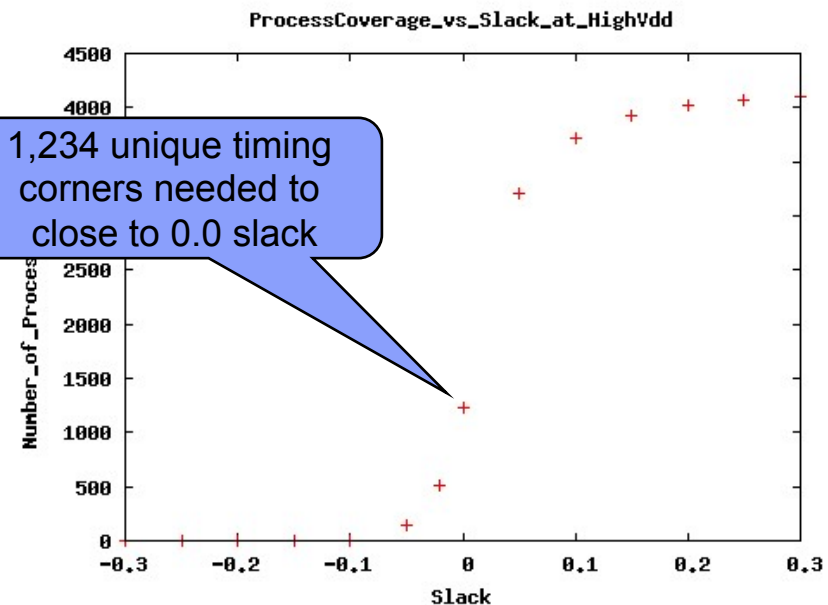
- Motivation
- Statistical Timing Overview
- Reliability Modeling
- Second Order Modeling
- Results
- Construction / Design Optimization
- Abstraction / Modeling Techniques
- Summary

Motivation for Statistical Timing - Avoiding the Exponential Trap of Corner-Based Timing



Note: 4 most common timing corners account
for 648 out of 10,616 negative slack tests

Number of unique corners at which
at least one end point had its worst slack



Statistical Timing - Canonical Model

- The canonical model is a bi-linear delay distribution consisting of multiple sources of variation and cross terms with respect to a single base variable
- Every delay quantity still has an Early (smallest) and Late (largest) value represented by Early and Late canonical models
- The canonical model form is defined as:

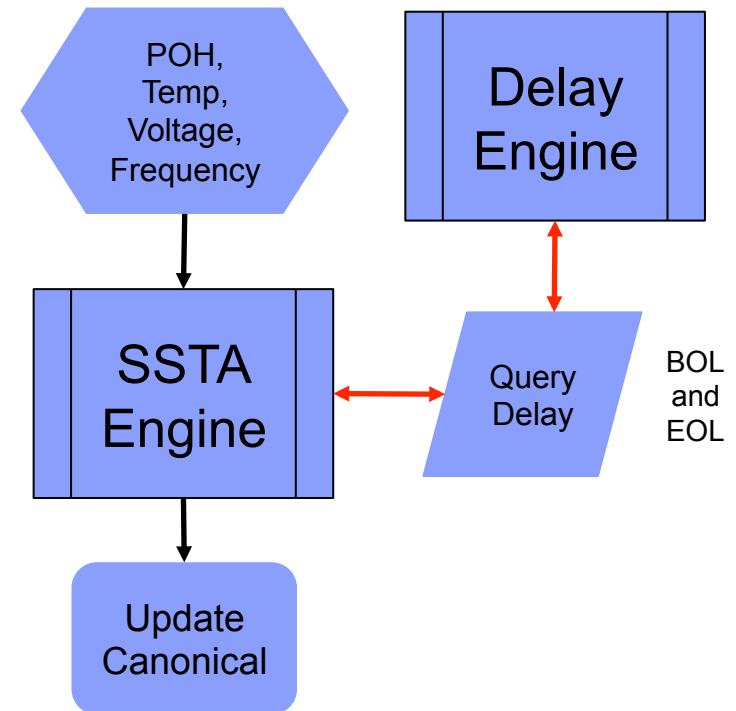
$$a_0 + a_1 \Delta X_1 + a_2 \Delta X_2 + \boxed{a_{1,2} \Delta X_1 \Delta X_2} \cdots + a_{n+1} \Delta R_a$$

The diagram illustrates the components of the canonical model equation. Arrows point from descriptive boxes to specific terms in the equation:

- A yellow box labeled "Constant (nominal) value in the absence of variations" points to a_0 .
- A green box labeled "Sensitivities" points to the coefficients a_1 and a_2 .
- A blue box labeled "Global random variables; these are probability distributions" points to the variables ΔX_1 and ΔX_2 .
- An orange box labeled "2nd order cross term" points to the boxed term $a_{1,2} \Delta X_1 \Delta X_2$.
- A purple box labeled "Independently random uncertainty" points to ΔR_a .

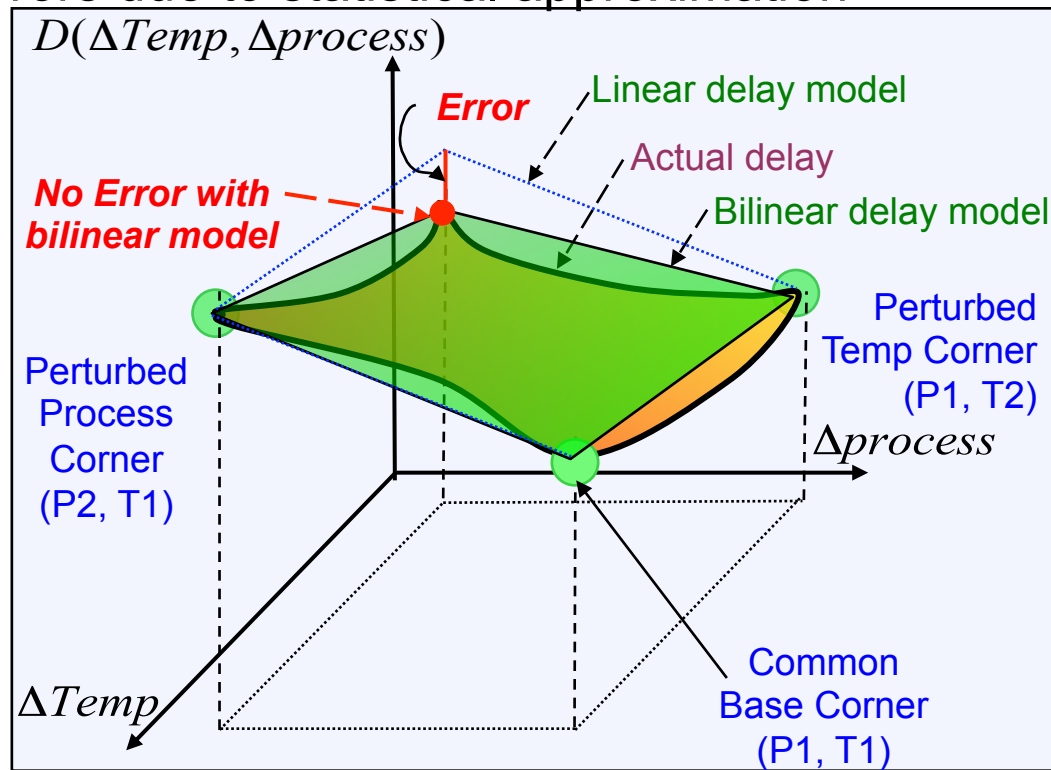
Product Reliability Modeling in SSTA

- **NBTI and Hot-e impacts transistor delay over time**
- **Delay Model Characterizes this**
 - Typically 4-10% at 3σ
- **SSTA incorporates source of variation using delay model**
- **SSTA computes variability to represent Beginning of Life and End of Life Conditions**
- **SSTA can tell you which condition is worse**



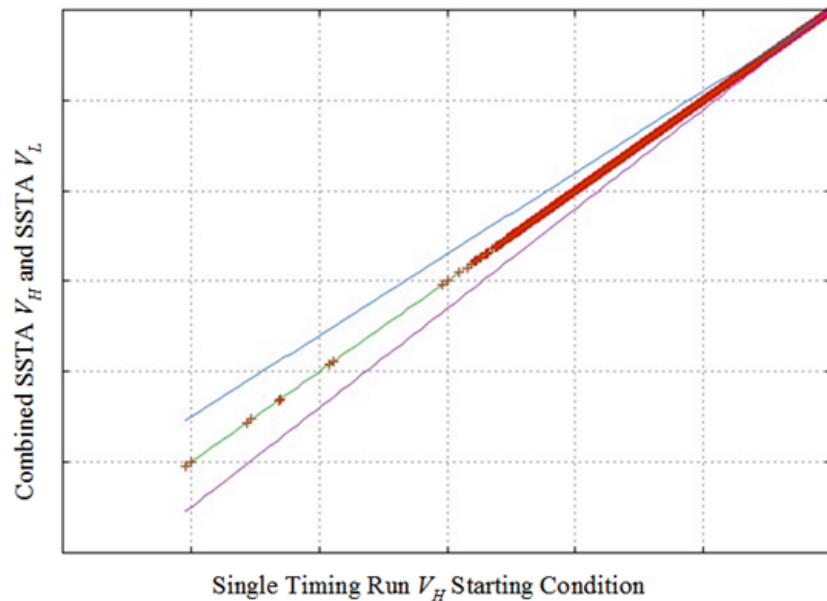
Bi-linear Modeling

- Use of 2nd order cross terms to create a Single Timing Run (STR)
 - Bounds errors due to statistical approximation

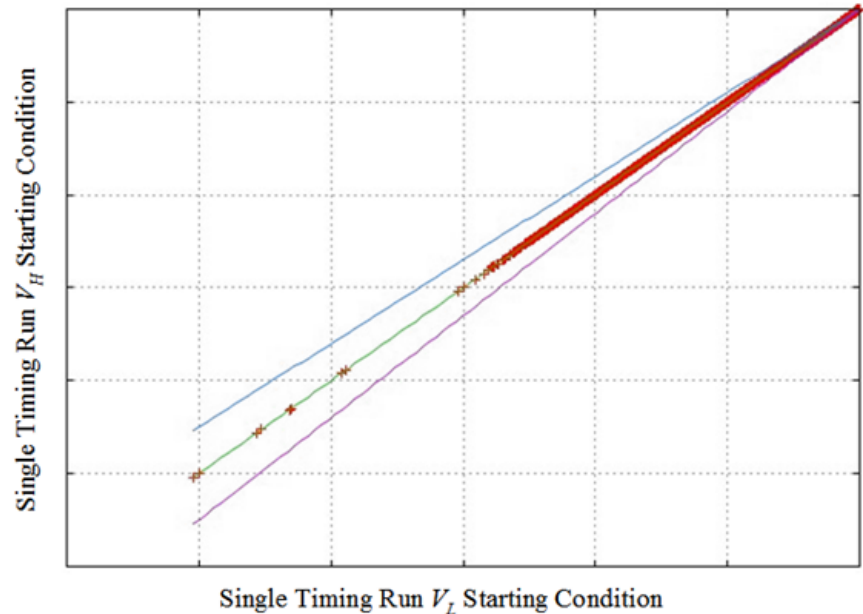


Actual delay shows variable ΔX is different at different Vdd values
The bi-linear model accounts for different ΔX across vdd space

Results



- X-axis: STR timing slacks with high voltage starting corner
- Y-axis: High voltage timing slacks and low voltage timing slacks from separate runs combined
- Conclusion: STR matches separate high and low voltage timing runs



- X-axis: STR timing slacks with low voltage starting corner
- Y-axis: STR timing slacks with high voltage starting corner
- Conclusion: Starting corner selection gives same results

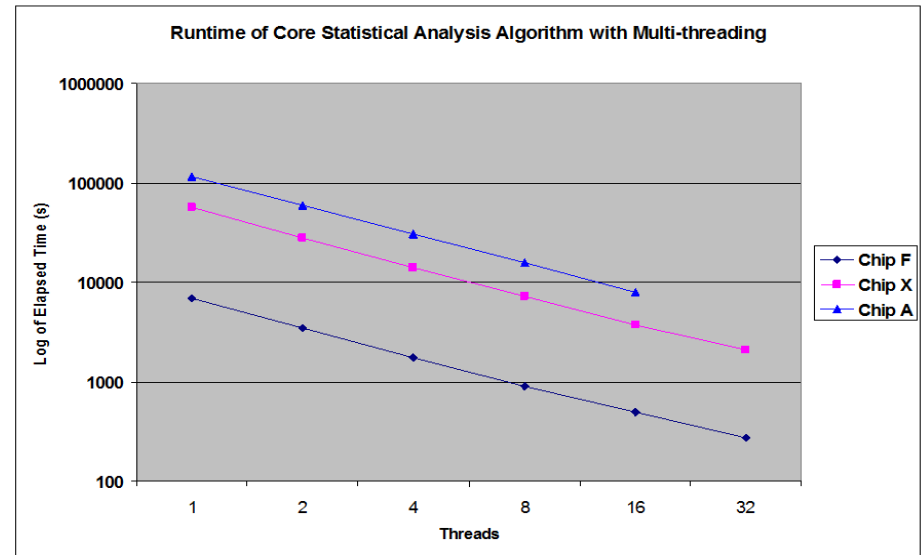
Runtime / Memory Comparison

■ Runtime

- Analysis is fully-threaded
- Runtime decreases linearly with increase in # threads
- For given X cpu machine:

$$STR[Xcpu] \prec \max\left(TR_{VH}\left[\frac{X}{2}cpu\right], TR_{VL}\left[\frac{X}{2}cpu\right]\right)$$

- Simulation results:
 - Show **16% Runtime Savings**



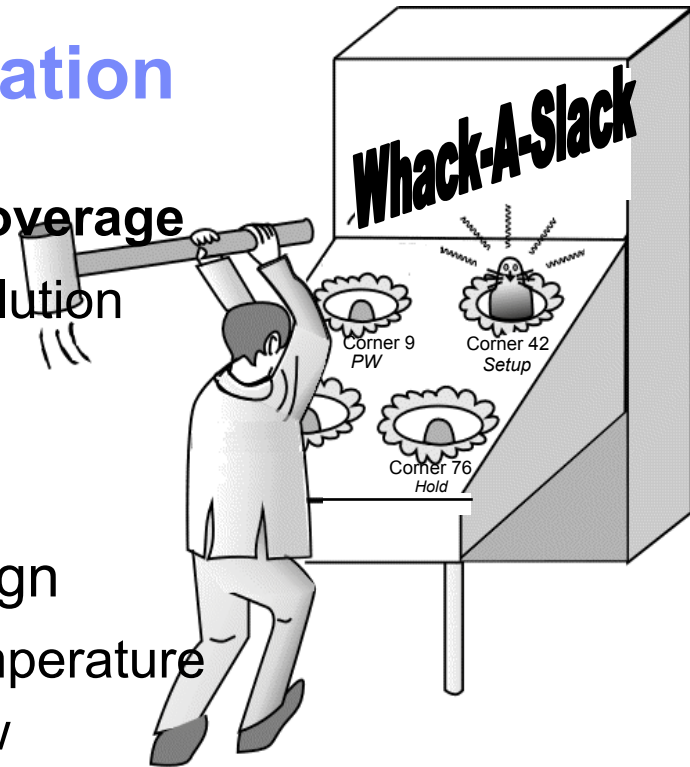
■ Memory

- STR loads timing models, parasitics, and netlist once
- STR memory footprint is less than separate high and low voltage timing runs
- Simulation results:
 - Show **30% Memory Savings**

Construction Flows & Optimization

■ SSTA: Comprehensive process space coverage

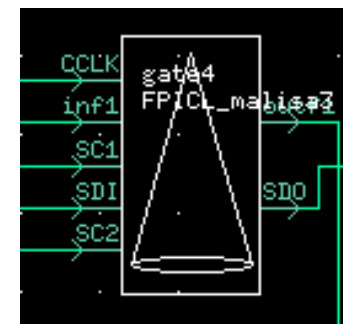
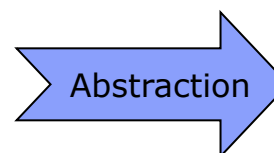
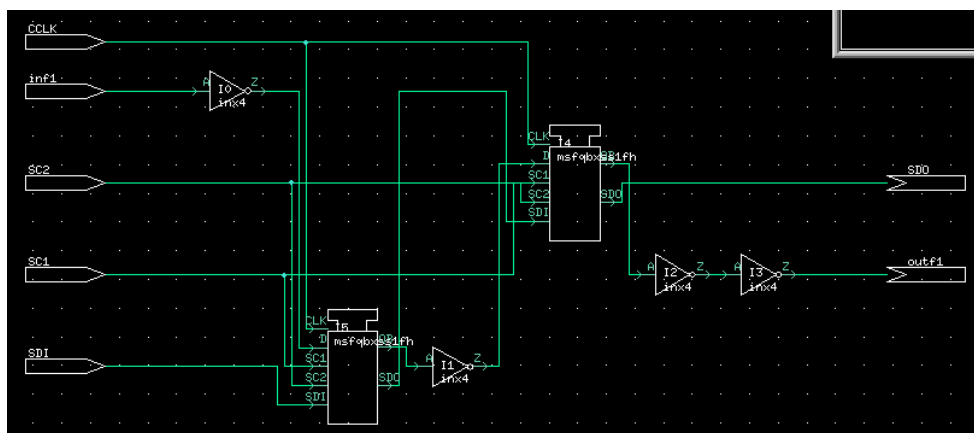
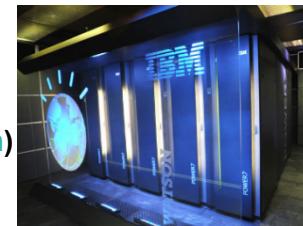
- Analyzes entire space, returns limiting solution
 - Avoids ‘whack-a-mole’ failing test iterations
 - Avoids or reduces any required margining
- SSTA flow can evolve along with design
 - Initially model Si-Process, Voltage, Temperature
 - Introduce metal parameters later in flow
 - Incremental parameter introduction can apply to other parms as well
- Other methods for managing runtime
 - Many of the performance filters have variable accuracy
 - Use low accuracy/high performance early in flow (large number of large fails)
 - > Can use other more aggressive techniques early on (e.g. constant pincap)
 - Increase accuracy as design matures (small number of small fails)



Statistical Abstraction: Enabling hierarchical timing with variability

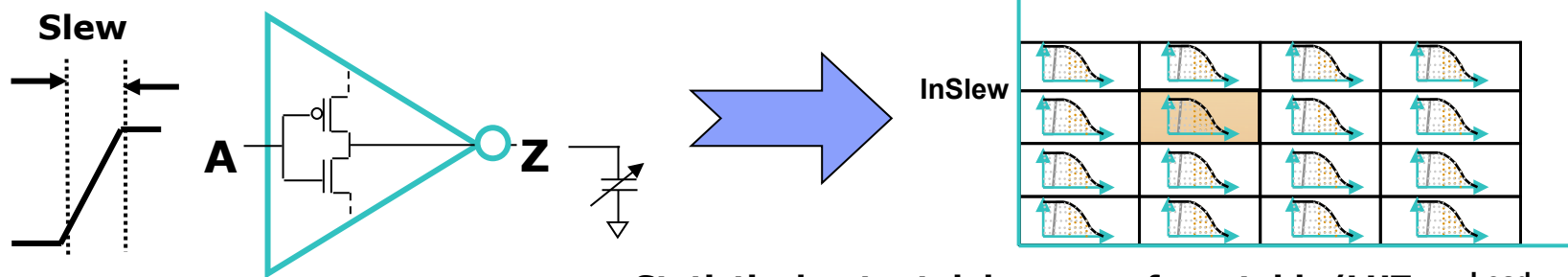
- Hierarchical designs: Partitioned (cores, units, macros)

(e.g. IBM POWER7-based Watson)



- Statistical abstracts (macro-models) of partitions used at parent level

- Creates Look-Up-Table (LUT) for each arc's **statistical**-delay and -output-waveform as function of **input-waveform's slew** and/or arc **capacitive-loading**



Statistical output delay, waveform table/LUT

Summary

- ✓ Statistical timing allows us to cover an exponential # process corners in a single analysis framework
- ✓ 2nd order cross-terms are supported
- ✓ Accuracy has been validated against exhaustive corner based timing
- ✓ SSTA is fully threadable and amenable to incremental optimization flows