# POWER – PERFORMANCE OPTIMIZATION METHODS FOR DIGITAL CIRCUITS

## Radu Zlatanovici
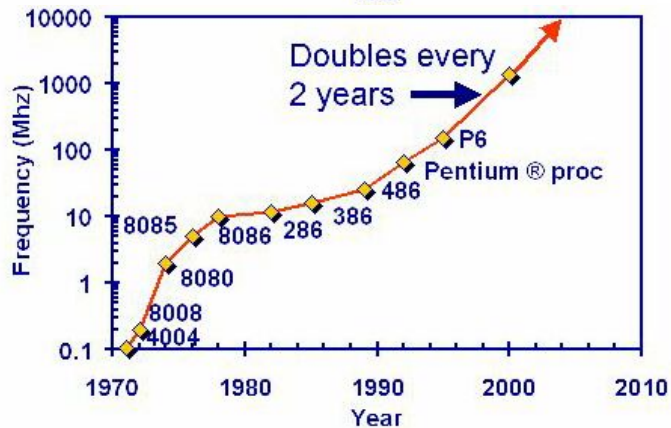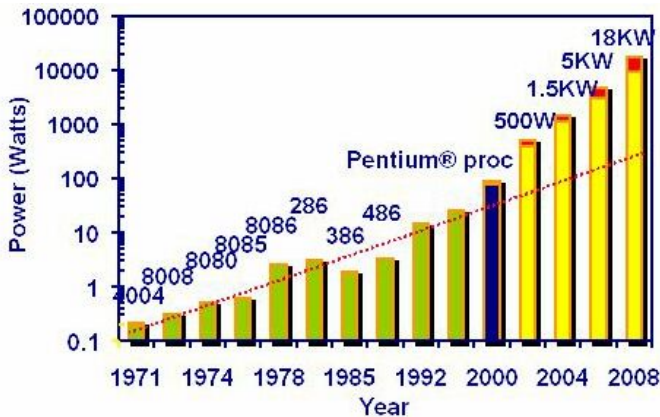
zradu@eecs.berkeley.edu
http://www.eecs.berkeley.edu/~zradu

Department of Electrical Engineering and Computer Sciences
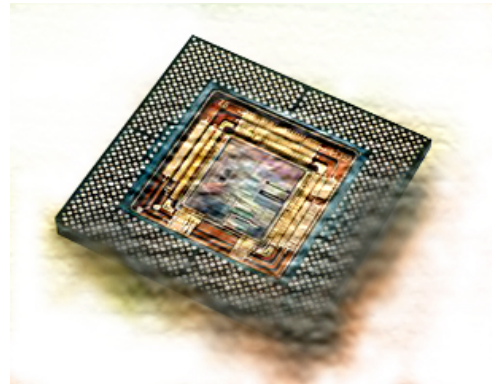**University of California, Berkeley**

# Power and Performance with Scaling

➢ **If we continue doing business as usual, both dynamic and leakage power will be a problem…**



…chips are getting hot…

…and phones leaky!

• Need to deliver maximum performance under power constraints

From S. Borkar, Intel

# A Common Problem



- Need to reduce power by 30%, while willing to give up 3% of performance

- What to do:
    - Decrease supply?
    - Increase thresholds?
    - Downsize?
    - Downsize latches or logic?
    - Use dual supplies?
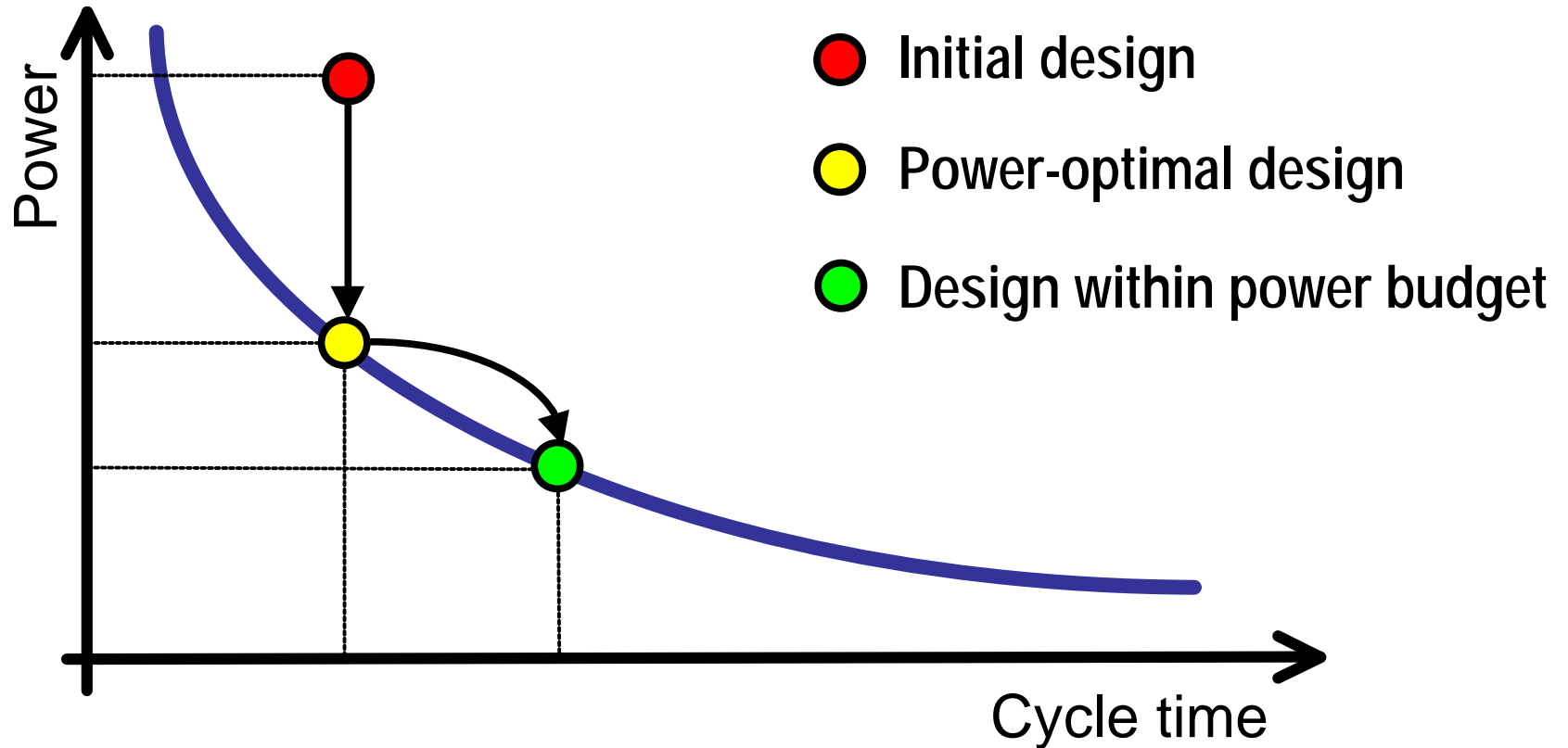    - Re-pipeline?
    - Parallelize?

# Outline

- Design as a power – performance optimization problem
- Fundamentals of circuit optimization
- Design examples
- Dealing with variations
- Conclusions
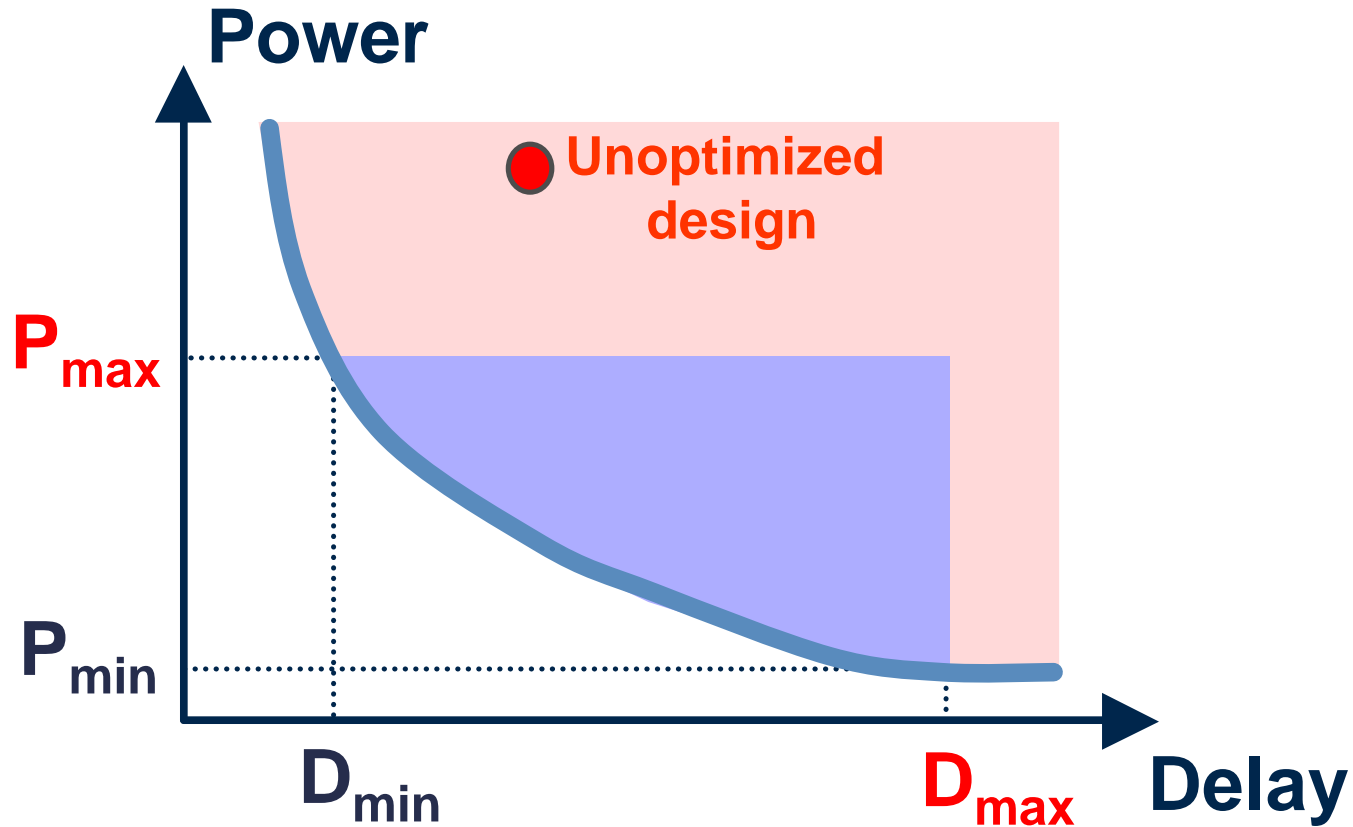
- Collaborative effort:
  - Students: R. Zlatanovici, D. Markovic, L.-T. Pang, J. Garrett, S. Kao
  - Faculty: B. Nikolic, R. Brodersen

# Power – Performance Optimization



OPTIMAL POWER – PERFORMANCE TRADEOFF CURVE

# Power Limited Operation



Power

$P_{max}$

$P_{min}$

Unoptimized design

$D_{min}$

$D_{max}$

Delay

**Achieve the highest performance under the power cap**

# Power Limited Operation

**Power**

**P**$_{max}$

**P**$_{min}$

**Unoptimized design**

**Var1**

**Design optimization curves**

**D**$_{min}$

**D**$_{max}$

**Delay**

## Achieve the highest performance under the power cap
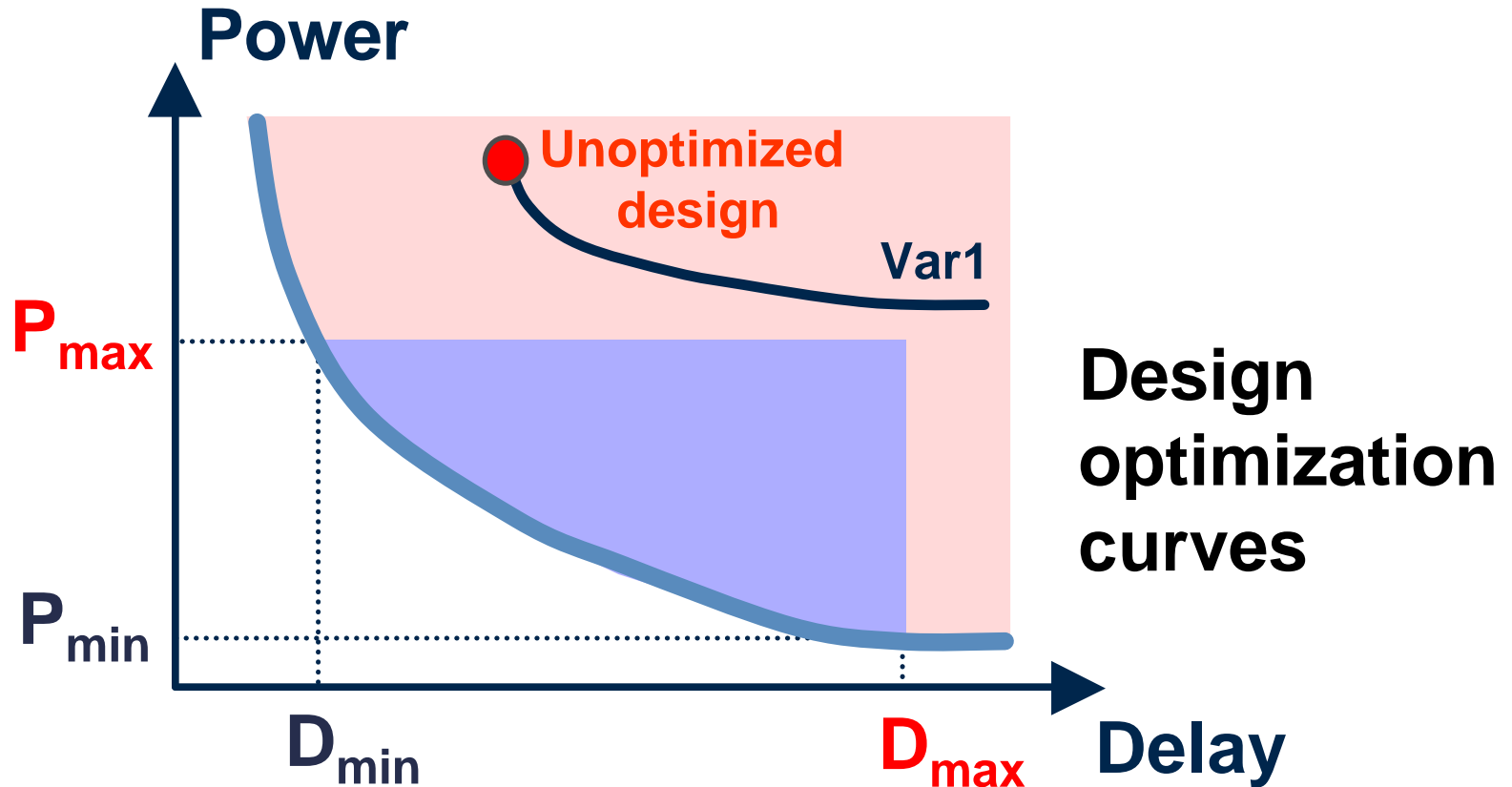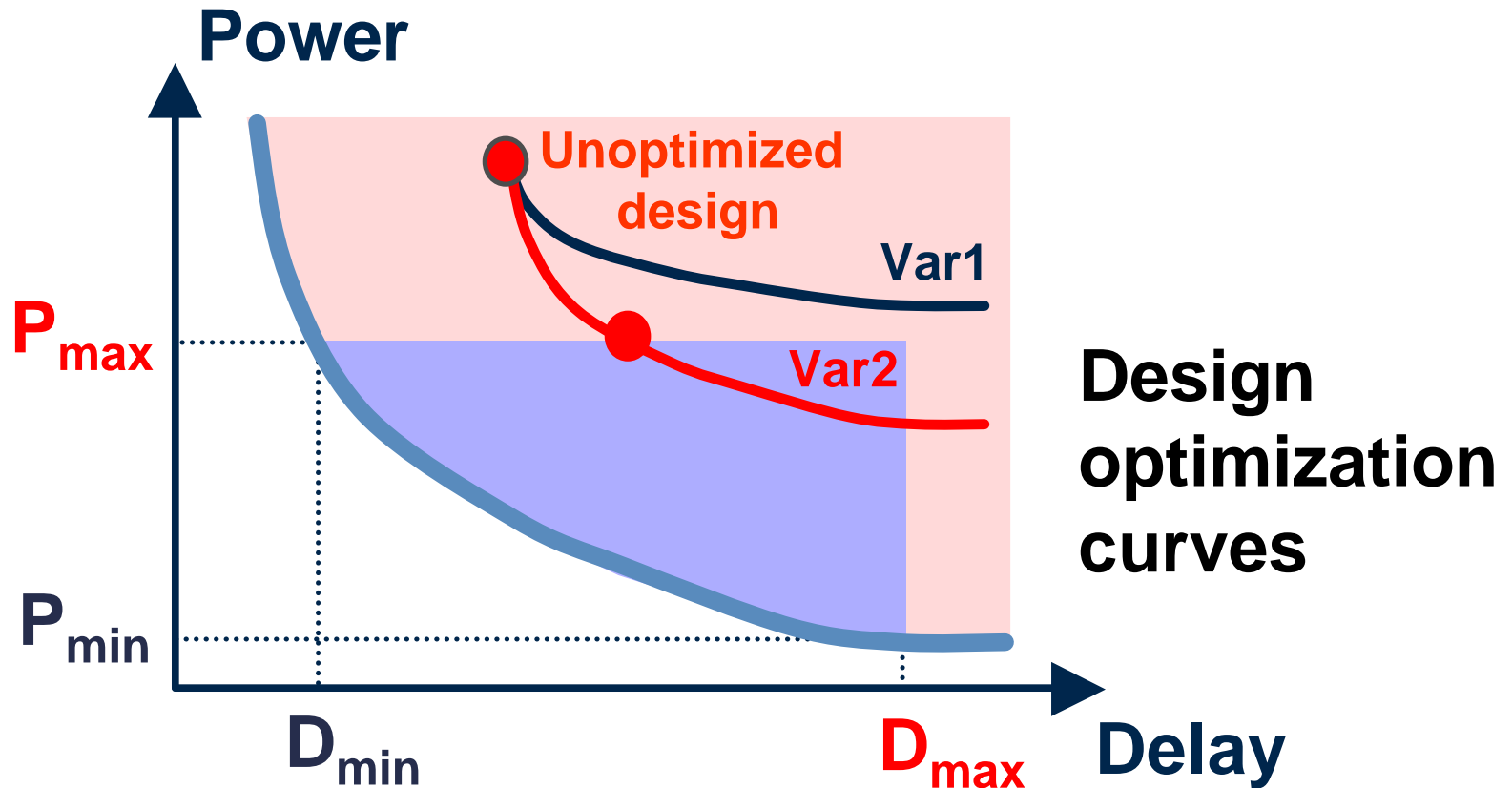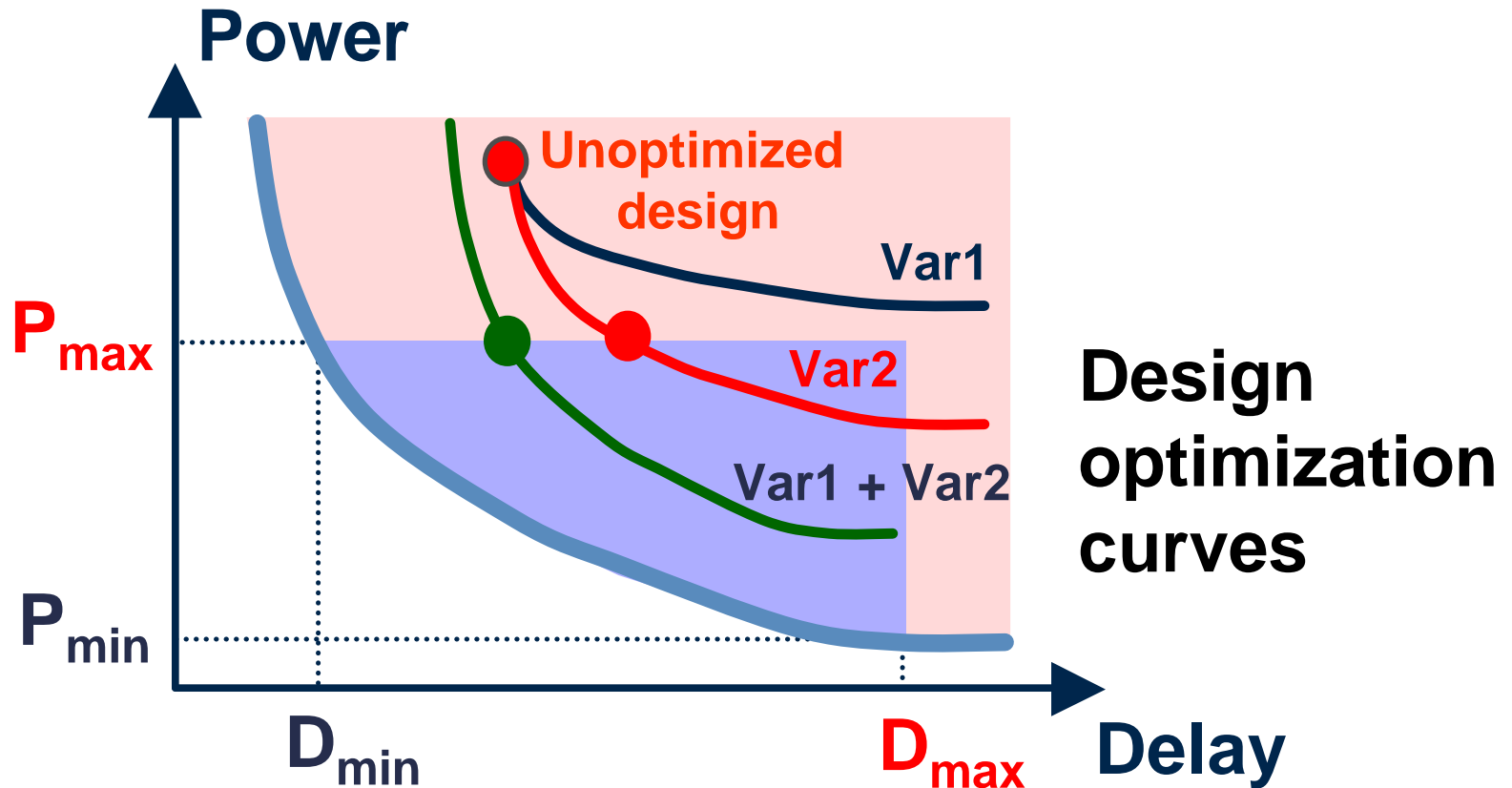
# Power Limited Operation



**Achieve the highest performance under the power cap**

# Power Limited Operation



**How far away are we from the optimal solution?**

# Power Limited Operation



**Global optimum – best performance**

# Design Optimization

- There are many sets of parameters to adjust:
  - Circuit
    (sizing, supply, threshold)
  - Logic style
    (domino, static, pass-gate, …)
  - Block topology
    (adder: CLA, CSA, RCA,…)
  - Micro-architecture
    (parallel, pipelined)

# Design Optimization

- There are many sets of parameters to adjust:
  - Circuit
    (sizing, supply, threshold)
  - Logic style
    (domino, static, pass-gate, …)
  - Block topology
    (adder: CLA, CSA, RCA,…)
  - Micro-architecture
    (parallel, pipelined)



Globally optimal boundary curve:
pieces of E-D curves for
different topologies

# Outline

- Design as a power – performance optimization problem
- **Fundamentals of circuit optimization**
- Design examples
- Dealing with variations
- Conclusions

# Optimization Problem

$$\min_{x \in R^n} f(x)$$

subject to

$$g_i(x) \leq 0 \quad i = 1..m$$

$$h_j(x) = 0 \quad j = 1..p$$

- Very difficult in the general case
- Optimality not guaranteed

## Convex Optimization

- $f$, $g_i$ – convex, $h_j$ – linear
- Key property: every local minimum is a global minimum
- Optimality guaranteed



$f(x)$

$ax_1 + (1 - a)x_2$

$x_1 \quad Y \quad x_2$

$x$

# Circuit Optimization

# Power Constrained Optimization Problem

**Minimize DELAY**

subject to

**Maximum POWER**

## Constraints:

- Maximum output slew
- Maximum internal slew
- Maximum input capacitance
- Minimum sizes

**Basic Result:**

- Power - Performance tradeoff curve

# PD, PD$^2$, P$^2$D ... : Obsolete



$P^0 D^1$

$P^a D^b$

$P^1 D^0$

Power

Delay

**ALL P-D METRICS ARE INCLUDED IN THE TRADEOFF CURVE**

# Choice of Models

## ANALYTICAL

- − Limited accuracy
- + Fast parameter extraction
- + Provide insight in the operation of the circuit
- + Can exploit their mathematical properties to help optimization
  - ≻ Target: convex optimization

## TABULATED

- + Very accurate
- − Slow to generate
- − No insight in the operation of the circuit
- − Can't guarantee convexity
- − Optimization is "blind"
- ☞ If convex models are any good, the optimization problem is not very "non-convex"

$$t_D = p + g \cdot \frac{C_L}{C_{in}}$$

$$C_{in} = g \cdot C_{inv} \cdot W$$

$$b = \frac{C_{on-path} + C_{off-path}}{C_{on-path}}$$



$g_1, p_1$     $g_2, p_2$     $g_3, p_3$     $g_4, p_4$

$$D = \sum_{i=1}^{N} p_i + g_i \cdot \frac{C_{load,i}}{C_{in,i}} = \sum_{i=1}^{N} p_i + \sum_{i=1}^{N} g_i \cdot \frac{b_i g_{i+1} W_{i+1}}{g_i W_i}$$

# Posynomial Functions

› Definition of a posynomial:

$$p(x) = \sum_j g_j \prod_{i=1}^{n} x_i^{a_{ij}} \quad a_{ij} \in R \quad g_j \in R^+$$

› A posynomial can be converted to a convex function using a simple change of variables: $x_i = e^{z_i}$ for $x_i > 0$

› Logical effort delay is posynomial

   › Fishburn, ICCAD '85: Elmore delay formula can be written as a posynomial

› Switching energy is linear in $W_i$ ➜ posynomial

› Convex optimization with posynomials is called <span style="color:red">geometric programming</span>

# V$_{DD}$- Dependent Analytical Delay Model

- Gate equivalent resistance can be computed from analytical saturation current models (a reduced form of the BSIM3v3 equation)

$$R_{EQ} = \frac{1}{V_{DD}/2} \int_{V_{DD}/2}^{V_{DD}} \frac{V_{DS} dV_{DS}}{I_{DSAT}} = \frac{3}{4} \frac{V_{DD}(\boldsymbol{b}_1 V_{DD} + \boldsymbol{b}_0 + V_{DD} - V_{TH})}{W \cdot K (V_{DD} - V_{TH})^2} (1 - \frac{7V_{DD}}{9V_A})$$

- Include supply and threshold dependency in the delay model:

$$d = c_2 R_{EQ}(W, V_{DD}, V_{TH}) + c_1 R_{EQ}(W, V_{DD}, V_{TH}) \cdot \frac{C_L}{C_{in}} + (\boldsymbol{h}_0 + \boldsymbol{h}_1 V_{DD}) t_{s,in}$$

- Accurate over a <span style="color:red">reasonable yet limited</span> range of fanouts (2.5-6), supplies and threshold (+/- 30%)
    - Most datapath blocks are within this range
- Compatible with convex optimization
    - Captures dependencies on V$_{DD}$ and V$_{TH}$ ➔ they can be optimization variables

*21*

Work by Joshua Garrett

# Vdd and W Optimization

# Vth,Vdd and W Optimization

# Outline

- Design as a power – performance optimization problem
- Fundamentals of circuit optimization
- **Design examples**
- Dealing with variations
- Conclusions

# Dual V$_{DD}$ ALU in Domino Logic



Y. Shimazaki, R. Zlatanovici, B. Nikolic: A Shared – Well Dual – Supply – Voltage 64-bit ALU

ISSCC'03, JSSC 03/2004

# Extending the Space: Dual Supply

# CLA Adders in E-D Space

- **Adders are common in critical paths**
- **CLA adders:**
  - Many designs, commonly used in practice
  - Recent interest in sparse adders
  - No fair comparison in energy – delay space
- **This work:**
  - Optimization of representative 64-bit adders in energy – delay space
  - Optimal 64-bit adder design

R. Zlatanovici, B. Nikolic: Power – Performance Optimal 64-Bit Carry-Lookahead Adders, ESSCIRC 2003

# 64-Bit CLA Adders

- Generic 64-bit adder block diagram

```
a [ 0:63 ]  ──────────►┌──────────┐   Carry[...]  ┌────────┐
                       │  Carry   │─────────────►│        │
b [ 0:63 ]  ──────┬───►│  Tree    │               │  Sum   │──► Sum [ 0:63 ]
                  │    └──────────┘               │ Select │
                  │    ┌──────────┐   S0[0:63]    │        │
             ├───►│  Sum     │─────────────►│        │
             └───►│Precompute│   S1[0:63]    └────────┘
                  └──────────┘
```

- Classical CLA, Ling equations
- Static, single-rail domino, compound domino logic
- Radix-2 and radix-4 carry trees
- Full and sparse trees (sparseness of 2 and 4)
- Use tabulated delay and analytical energy models (switching, leakage)

# Equation Set Comparison



Ling adders achieve shorter delays
Radix-4 are faster than radix-2

# Chosing a Logic Style



Energy [pJ] vs Delay [FO4]

Legend:
- Static R2
- Domino R4
- Domino R2
- Compound Domino R2

- Static adders are low power but slow
- Domino logic is the choice for short cycle times

# Full vs. Sparse Trees

## R2 Full

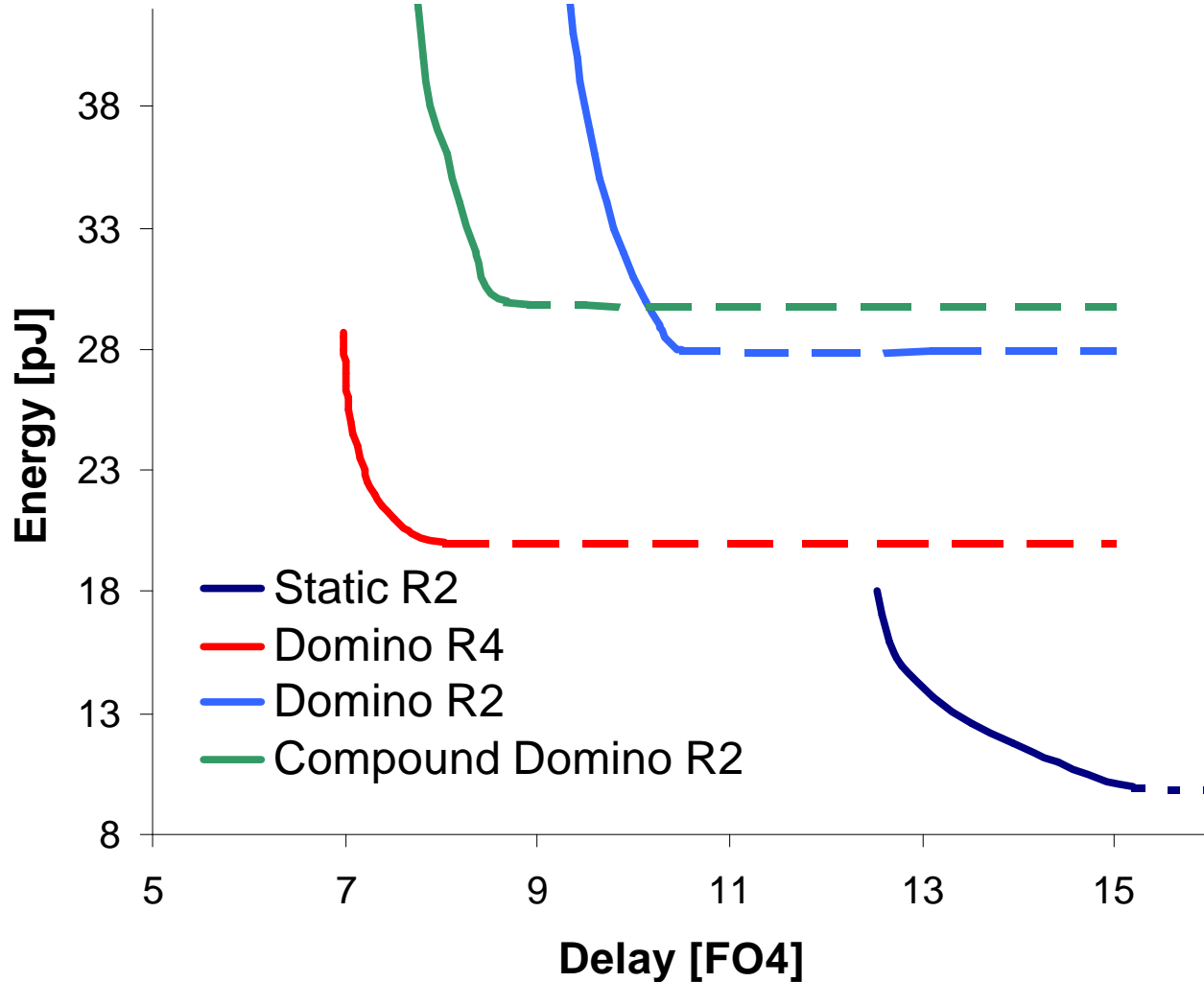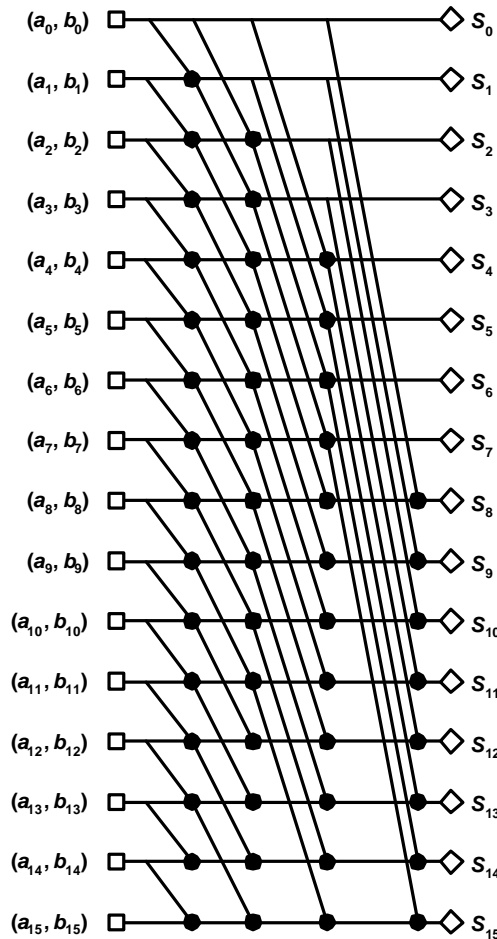| | |
|---|---|
| $(a_0, b_0)$ | $S_0$ |
| $(a_1, b_1)$ | $S_1$ |
| $(a_2, b_2)$ | $S_2$ |
| $(a_3, b_3)$ | $S_3$ |
| $(a_4, b_4)$ | $S_4$ |
| $(a_5, b_5)$ | $S_5$ |
| $(a_6, b_6)$ | $S_6$ |
| $(a_7, b_7)$ | $S_7$ |
| $(a_8, b_8)$ | $S_8$ |
| $(a_9, b_9)$ | $S_9$ |
| $(a_{10}, b_{10})$ | $S_{10}$ |
| $(a_{11}, b_{11})$ | $S_{11}$ |
| $(a_{12}, b_{12})$ | $S_{12}$ |
| $(a_{13}, b_{13})$ | $S_{13}$ |
| $(a_{14}, b_{14})$ | $S_{14}$ |
| $(a_{15}, b_{15})$ | $S_{15}$ |

## R2 SP2

| | |
|---|---|
| $(a_0, b_0)$ | $S_0$ |
| $(a_1, b_1)$ | $S_1$ |
| $(a_2, b_2)$ | $S_2$ |
| $(a_3, b_3)$ | $S_3$ |
| $(a_4, b_4)$ | $S_4$ |
| $(a_5, b_5)$ | $S_5$ |
| $(a_6, b_6)$ | $S_6$ |
| $(a_7, b_7)$ | $S_7$ |
| $(a_8, b_8)$ | $S_8$ |
| $(a_9, b_9)$ | $S_9$ |
| $(a_{10}, b_{10})$ | $S_{10}$ |
| $(a_{11}, b_{11})$ | $S_{11}$ |
| $(a_{12}, b_{12})$ | $S_{12}$ |
| $(a_{13}, b_{13})$ | $S_{13}$ |
| $(a_{14}, b_{14})$ | $S_{14}$ |
| $(a_{15}, b_{15})$ | $S_{15}$ |

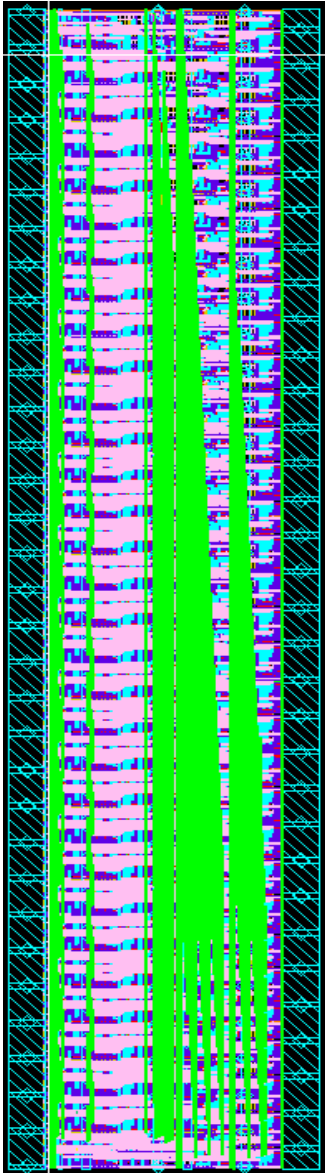## Sparse trees:

↘ Heavier load on the carry tree
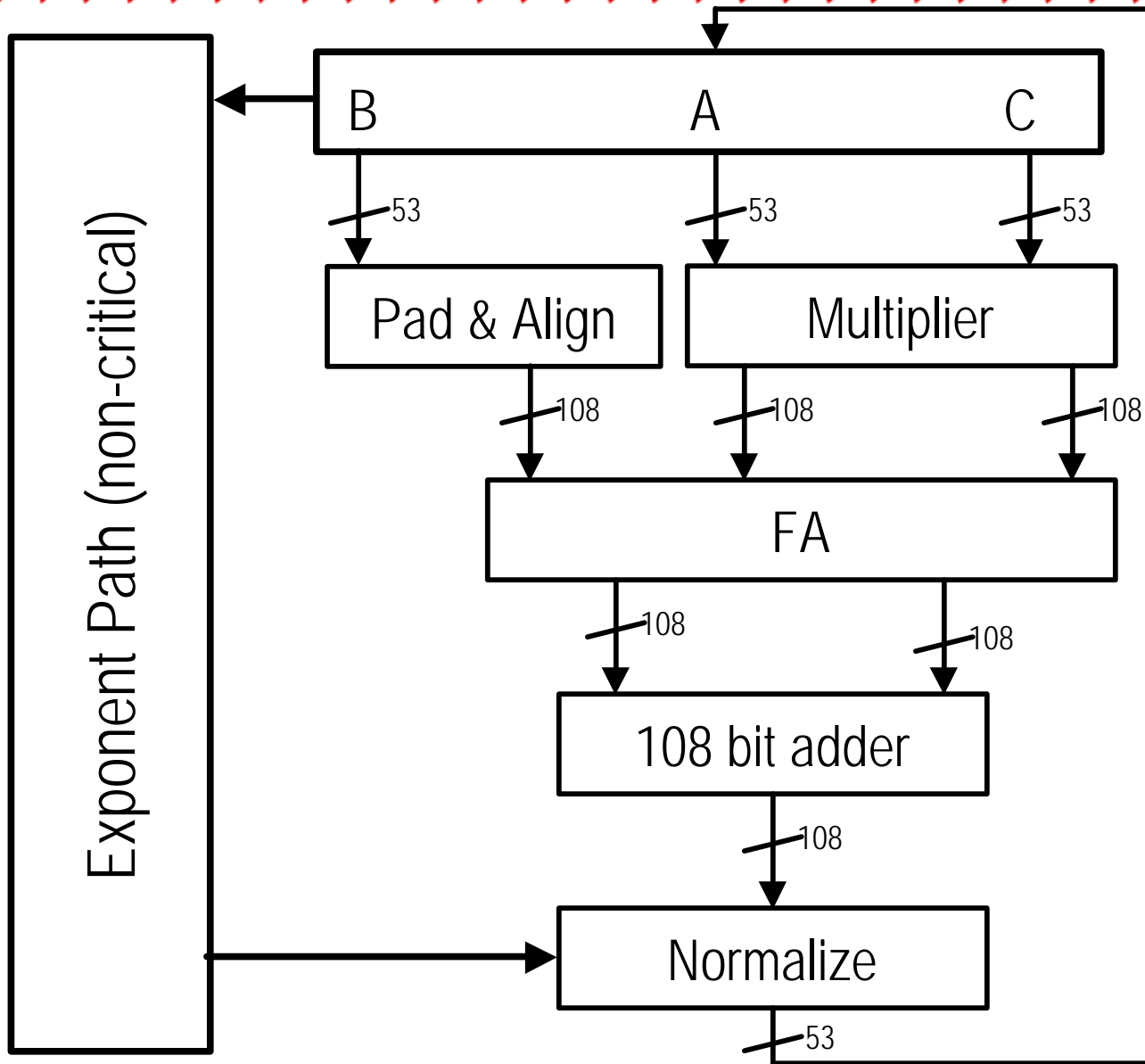
↗ Reduced input loading

↘ More complex sum precompute gates

# Proof of Existence: Fastest Adder

- Radix-4  sparse-2 domino Ling adder
- Technology:
  - 90 nm 1P 7M
  - $V_{DD}$ = 1 V
- Performance:
  - Delay: 210 ps (post – layout simulation)
  - Energy: 9.1 pJ / cycle (optimization tool)
- **Core dimensions: 417.3 mm x 75.3 mm**
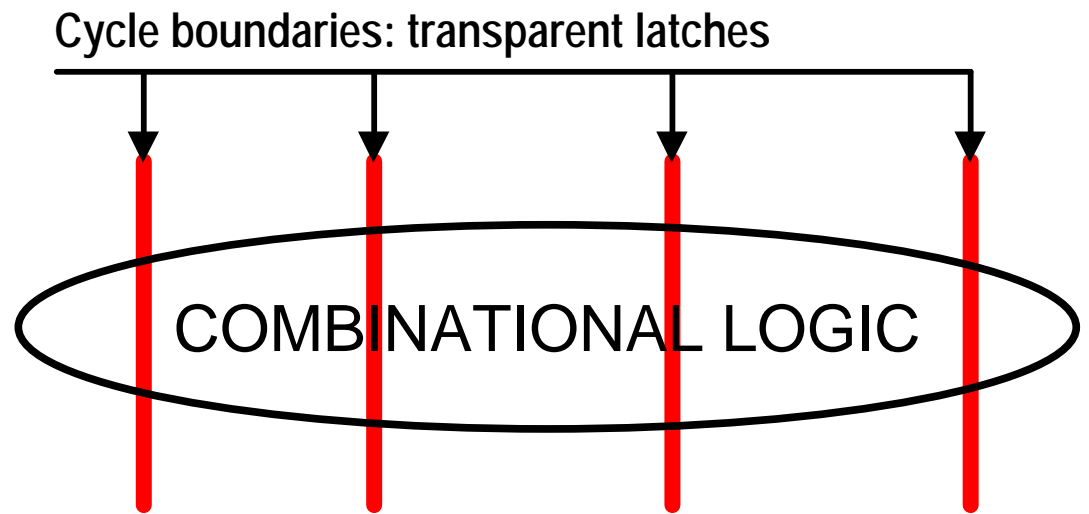- **Chip to be taped out 11/1/04**

*with Sean Kao

# Micro-Architecture Optimization: Power4 FPU



- 5 cycles
- 2-phase clocking
- Static CMOS

Exponent Path (non-critical)

B          A          C

53          53          53

Pad & Align          Multiplier

108          108          108

FA

108          108

108 bit adder

108

Normalize

53

# Optimizing Pipelined Circuits

| Models: Posy-nomials | Block Level Netlist | Minimize $T_{CYCLE}$ Subject to Maximum ENERGY | Gate Sizes Latch Positions |
|---|---|---|---|

Static timer

Optimizer

**Optimal Pipeline Configuration**

Cycle boundaries: transparent latches

COMBINATIONAL LOGIC

> **Fix pipeline depth**

> **Find shortest cycle time for fixed cutset**
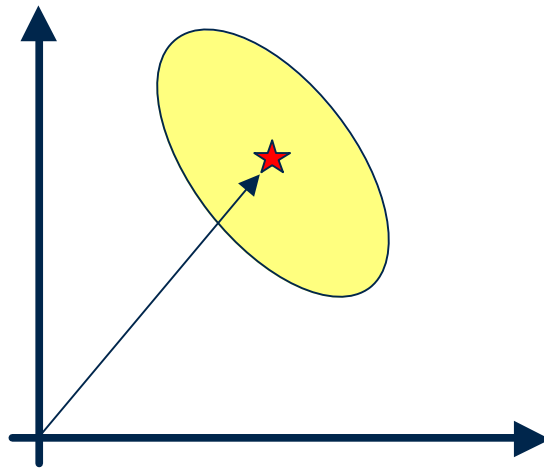
> **Search for optimum cutset**

Work in progress…

# Outline

- Design as a power – performance optimization problem
- Fundamentals of circuit optimization
- Design examples
- **Dealing with variations**
- Conclusions

# Robust Optimization

- Parameters are within an **ellipsoid** centered on the nominal values
- Optimize the **worst case**



- Compatible with convex optimization for the presented analytical models
- Problem: computing the ellipsoid

# Stochastic Optimization for Yield

- Parameters are random variables centered on the nominal values
- Optimize for a desired yield

$$\min_{x \in R^n} f(x)$$

$$\text{subject to}$$

$$g_i(x) \le 0 \quad i = 1..m$$

$\longrightarrow$

$$\min_{x \in R^n} f(x)$$

$$\text{subject to}$$

$$P(g_i(x) \le 0, \quad i = 1..m) \ge h$$

- Compatible with convex optimization under certain conditions
  - Convex analytical models
  - Jointly Gaussian parameters
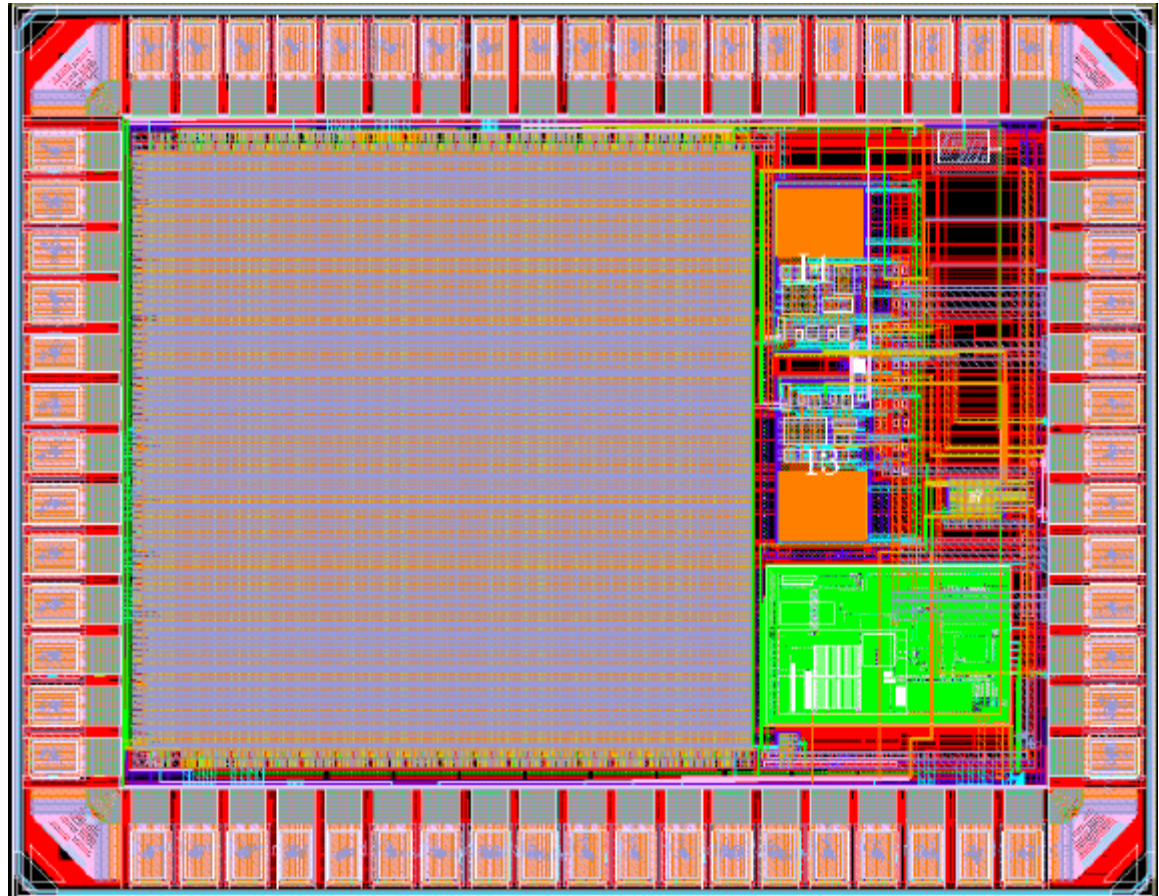- Problem: finding the distributions of the parameters, especially the correlations

# Impact of Layout on Variations

- Stacked gates vs. non-stacked gates (e.g. gates vs. buffers)
- Proximity effects, orientation of gates, metal layer above gate (annealing)



Stacked gates

Non-stacked gates

**Work by Liang – Teck Pang**

# Test Chip

- 90nm 1P 7M

- $V_{DD}$ = 1 V

- 1.55 x 1.17 mm$^2$

- Taped out 9/04

- To be packaged and tested

- Measurements will provide statistical data for the stochastic optimizer

# Conclusions

> Power and performance are the two sides of the same coin.
> The connection: the power – performance tradeoff curve.

- **Built a suite of stackable tools to design "power – performance optimal circuits"**

- **Design space explorations:**
  - Different optimization variables
  - Various levels of abstraction
  - Impact of variations

- **Experimented the tool suite on various designs**
  - Dual-supply ALU
  - CLA Adders in the E-D space