# Methodology for Energy-Efficient Design of Digital Circuits

## Vojin G. Oklobdzija,

Advanced Computer Systems Engineering Laboratory
University of California / University of Texas
TxACE: Center of Excellence for Analog Circuits
http://www. acsel-lab.com

IEEE SSCS, Distinguished Lecture
Santa Clara Valley Chapter
April 15, 2010

# Summary of the Presentation

Energy Efficiency of Digital CMOS Circuits

- o Problems

- o Energy-Delay Relationship

- o Minimizing Energy for a given delay

- o Methodology

- o Determining the best structures for high-performance system
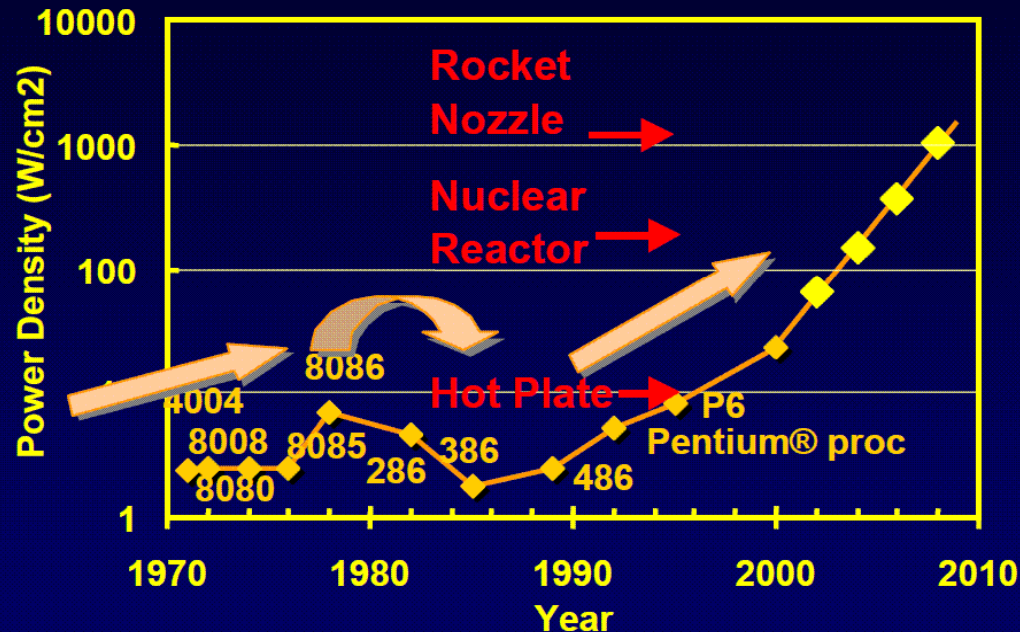
- o Implications on the architecture

# Challenges in High-Performance Design

❑ Optimizing for **power** – *not speed*! (or maximizing speed under the power budget)

❑ Logical Effort (LE) optimizes for speed, regardless of power i.e. brings us in the worse energy spot.

❑ Our method optimizes for:
  *power @ given speed* or *speed @ given power*

❑ We are developing new approaches for power efficiency (overlooked by delay optimization) applicable to:
  - Circuit structures
  - Design techniques
  - Energy-Delay Space
  - Creation of optimal Standard cell (ASIC) libraries

# Motivation for Energy Efficient Design
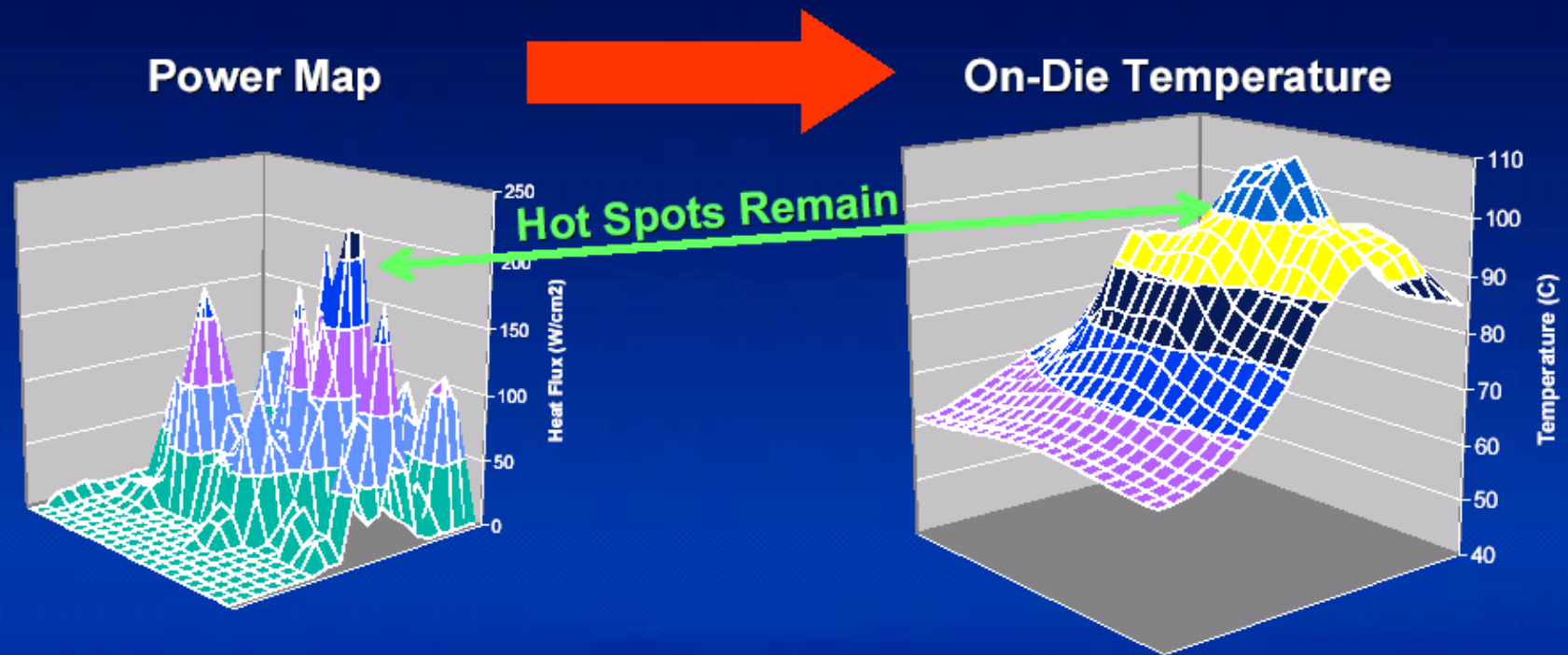


Shekhar Borkar

## Power density will increase

Power density too high to keep junctions at low temp

- Power density passed the level found in the Nuclear Reactor !
- Power density degrades the reliability and speed.

# Power Density: The Future

**Power Map**

**On-Die Temperature**
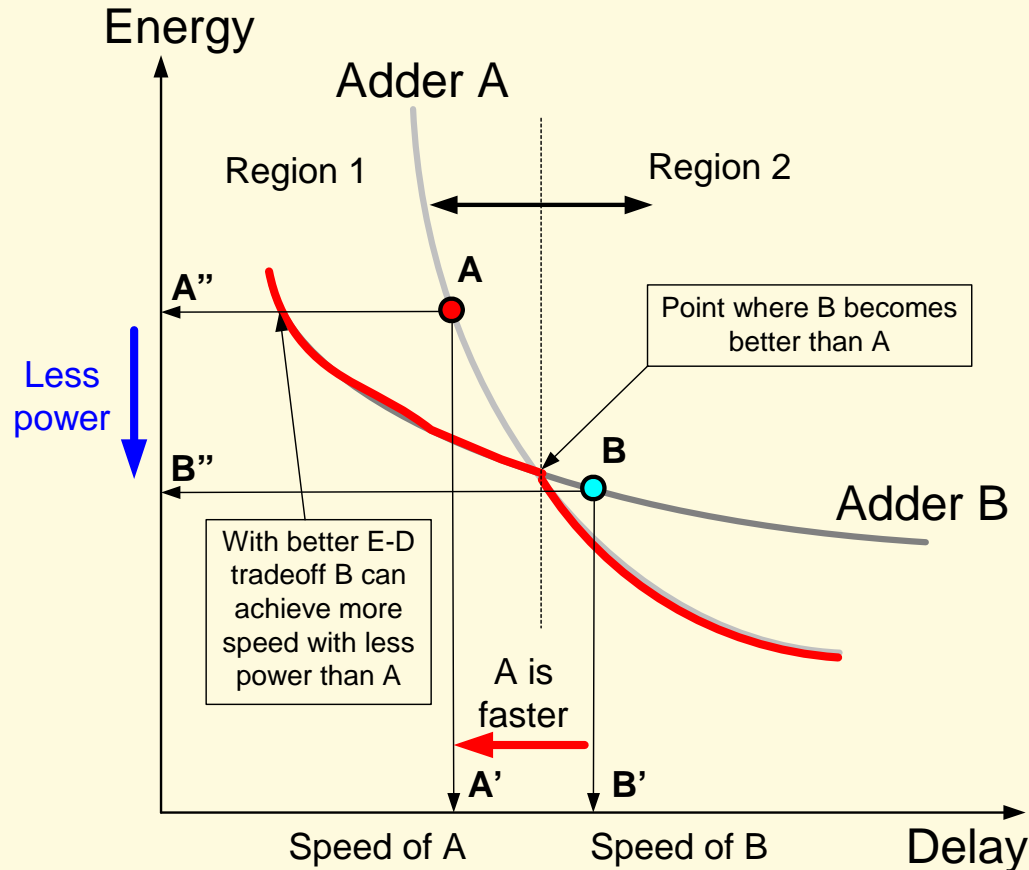
Hot Spots Remain

- **With high power density, cannot assume uniformity**
  - As die temperature increases, CMOS logic slows down
  - At high die temp., long-term reliability can be compromised
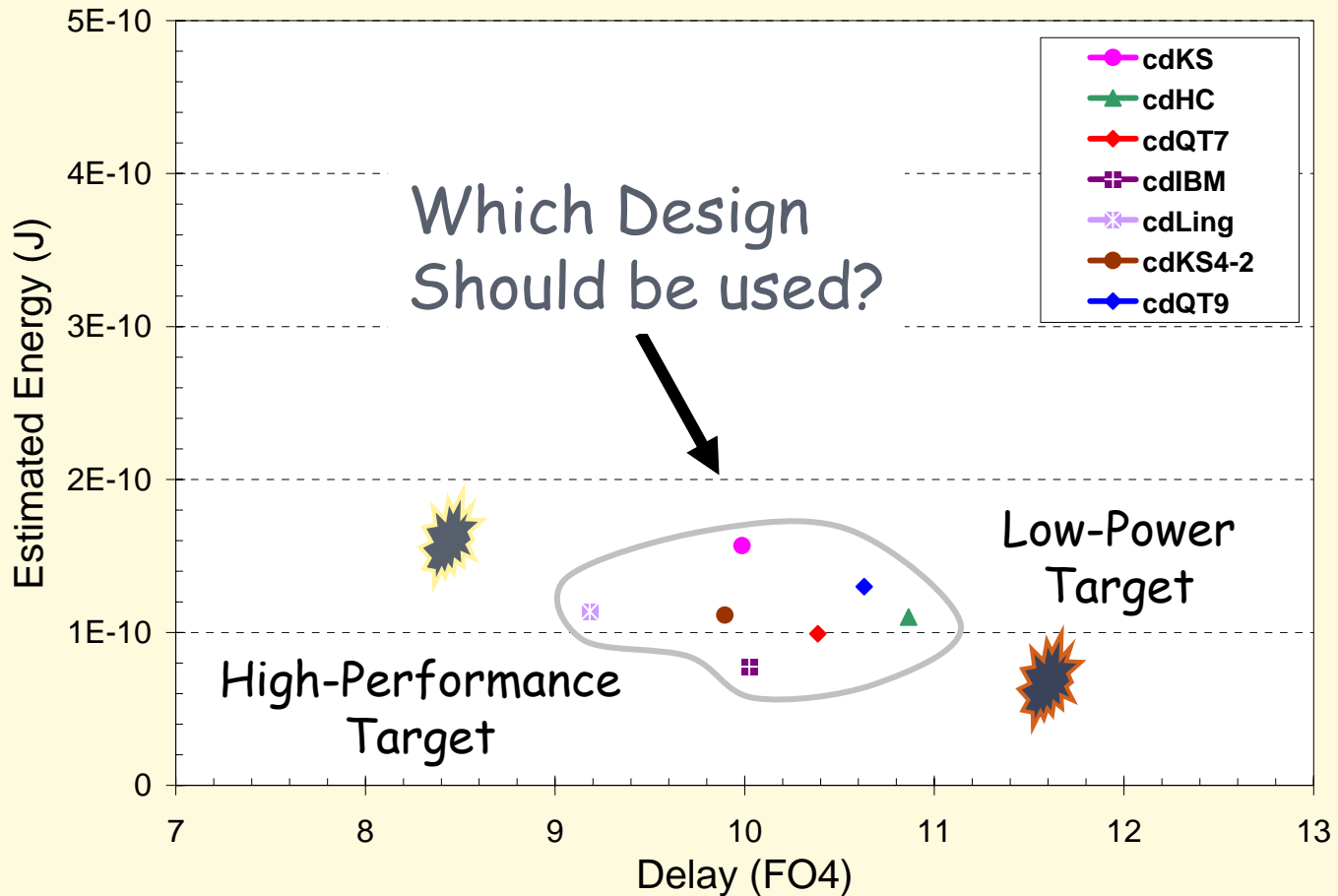
intel®

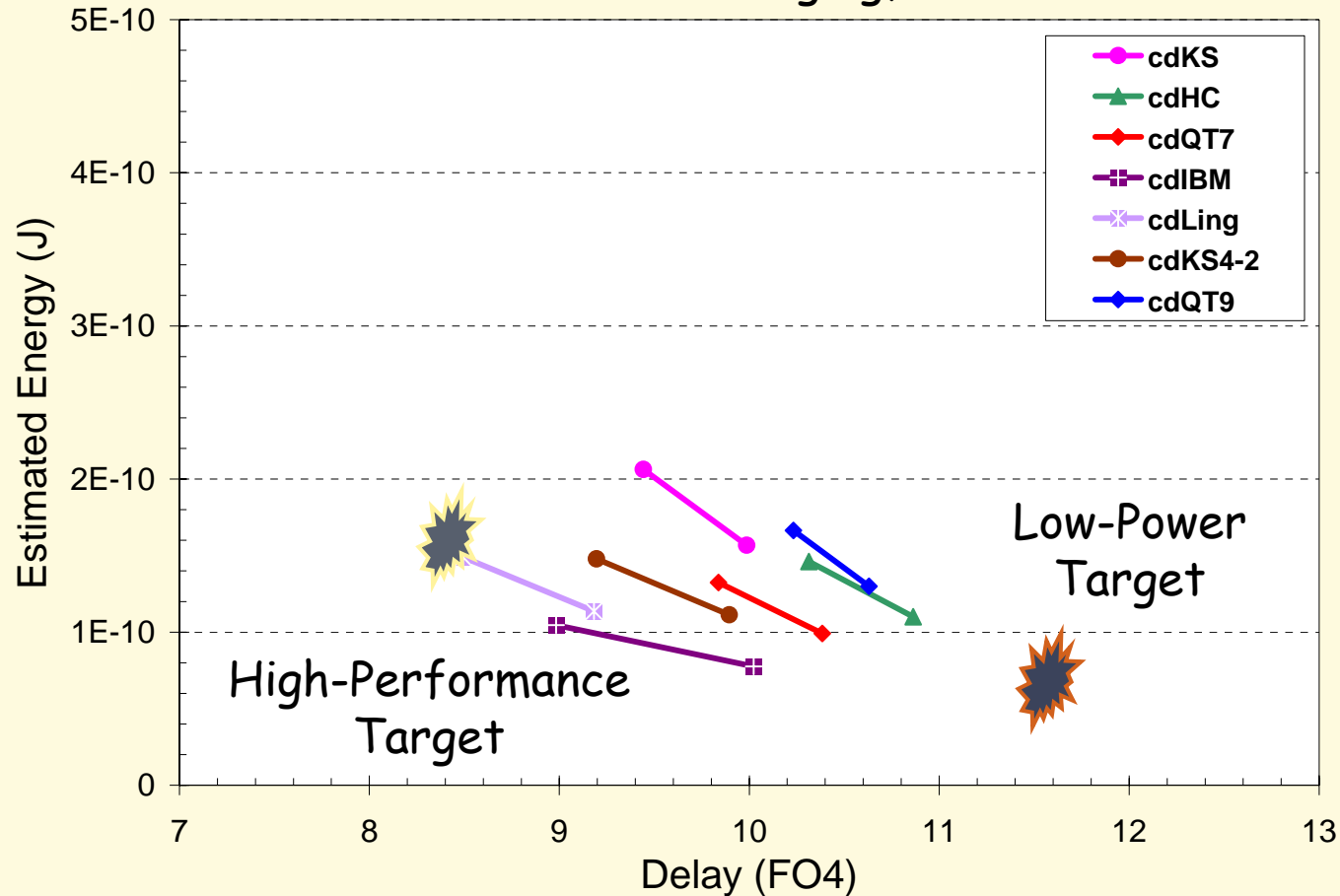# Energy-Delay Relationship

# Energy-Delay Space View



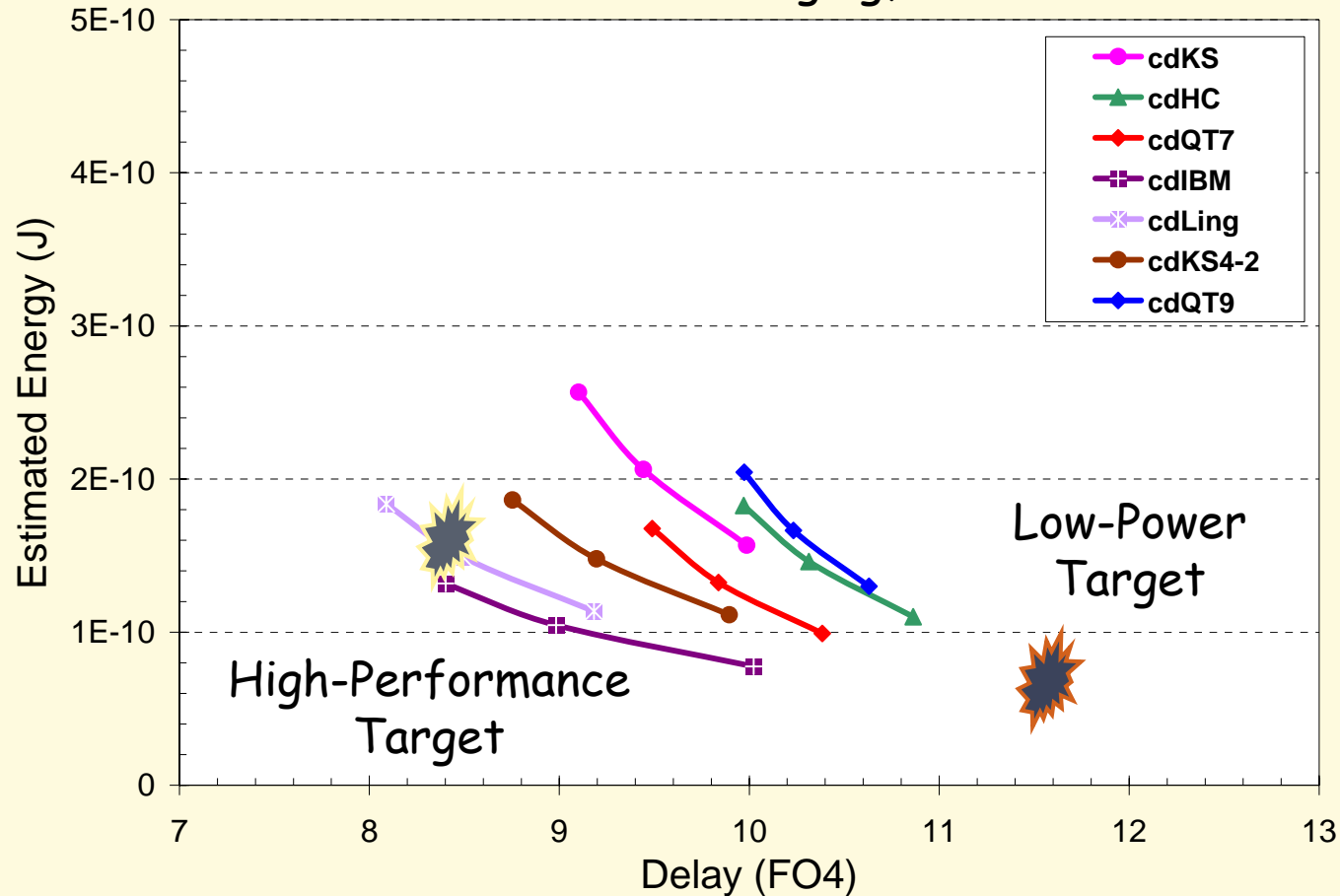- Must look at Energy-Delay Space of designs

# Energy-Delay Space View



April 19, 2010                    Energy-Efficient CMOS Circuit Design

# Energy-Delay Space View



H is changing, w/ Cout=constant

Legend: cdKS, cdHC, cdQT7, cdIBM, cdLing, cdKS4-2, cdQT9

Estimated Energy (J) vs Delay (FO4)

Low-Power Target

High-Performance Target

• Begin to see characteristics of designs

# Energy-Delay Space View



H is changing, w/ Cout=constant

Legend:
- cdKS
- cdHC
- cdQT7
- cdIBM
- cdLing
- cdKS4-2
- cdQT9

Y-axis: Estimated Energy (J)
X-axis: Delay (FO4)

Low-Power Target

High-Performance Target
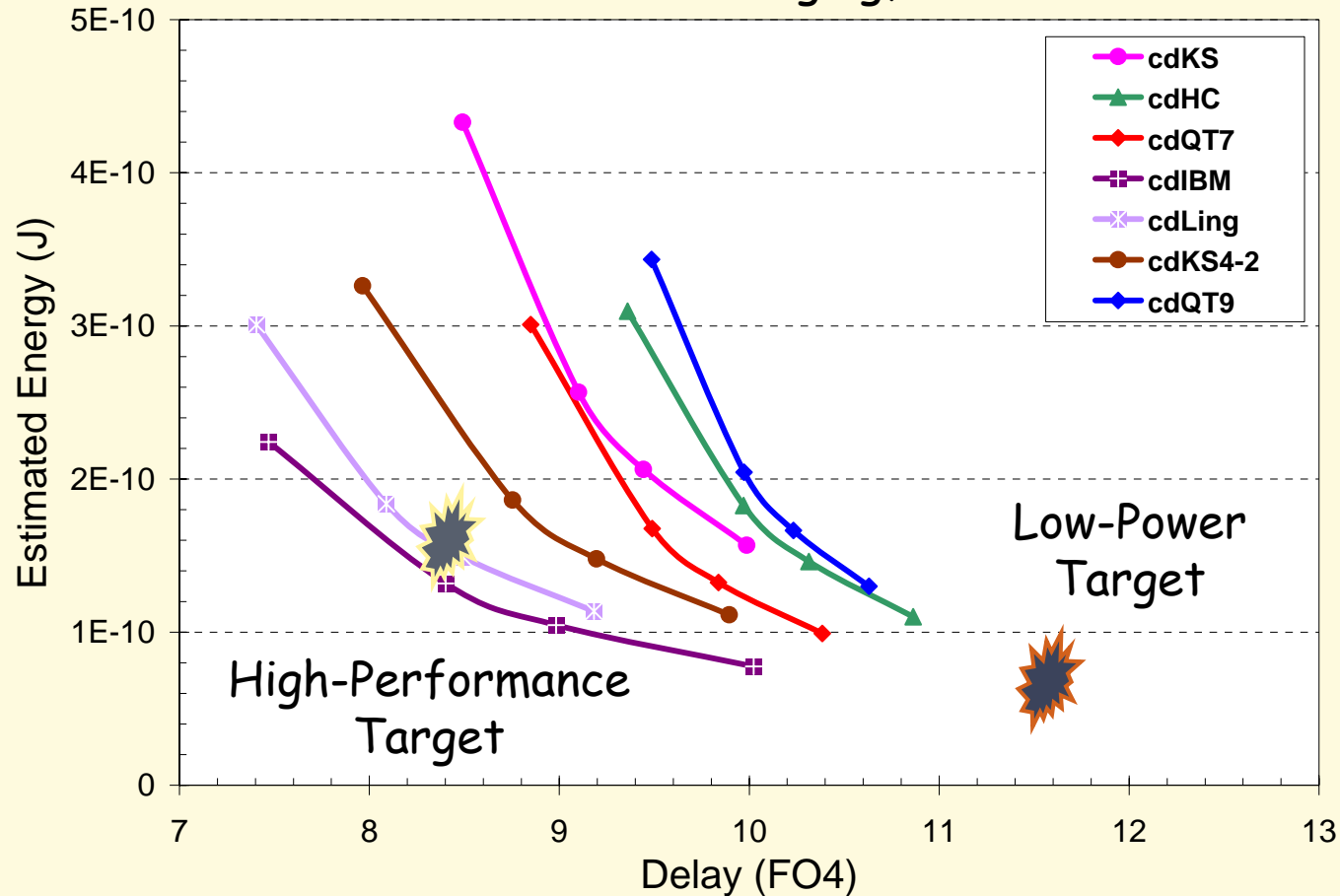
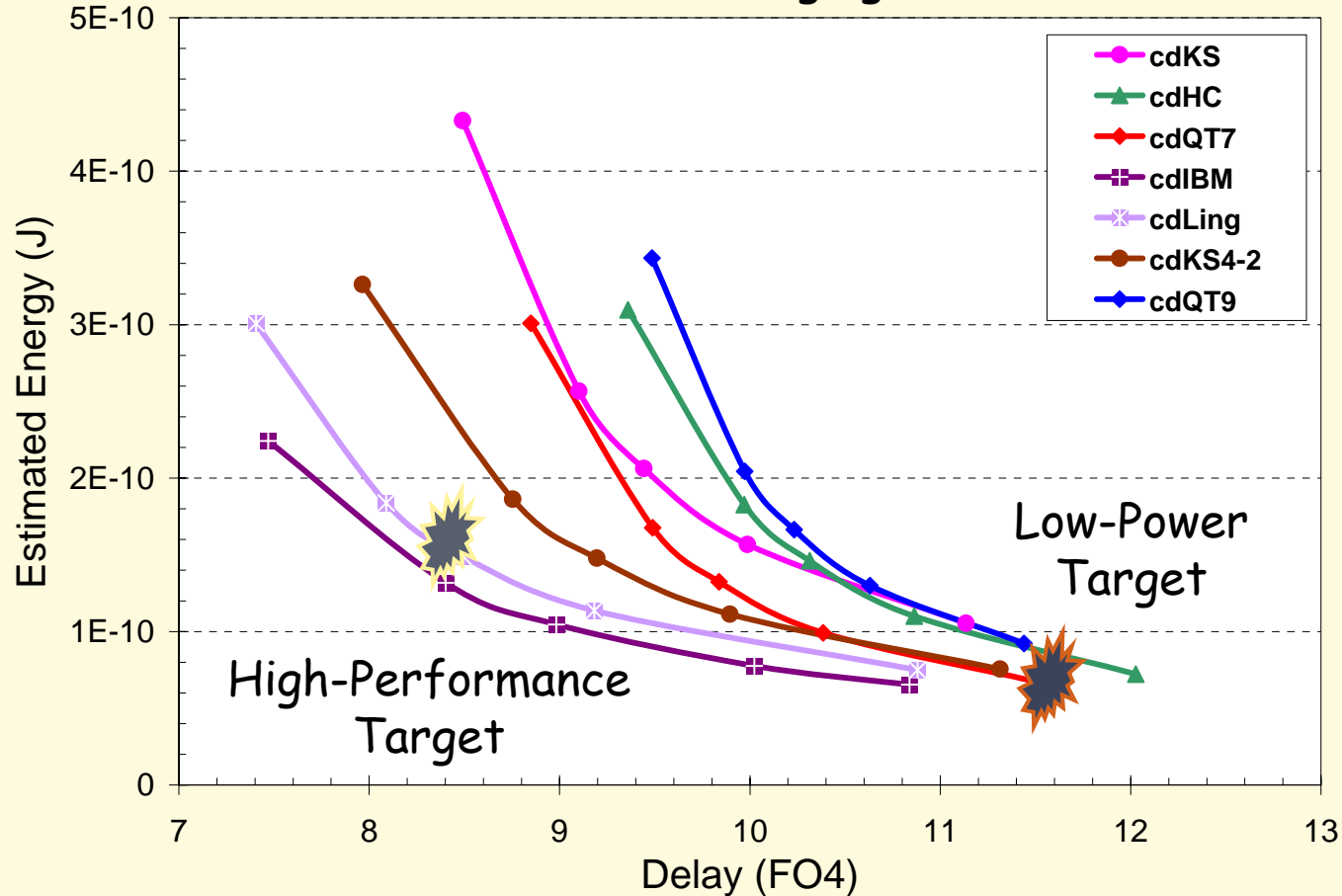- Begin to see characteristics of designs

# Energy-Delay Space View

H is changing, w/ Cout=constant



- Best High-Performance designs are clearly seen
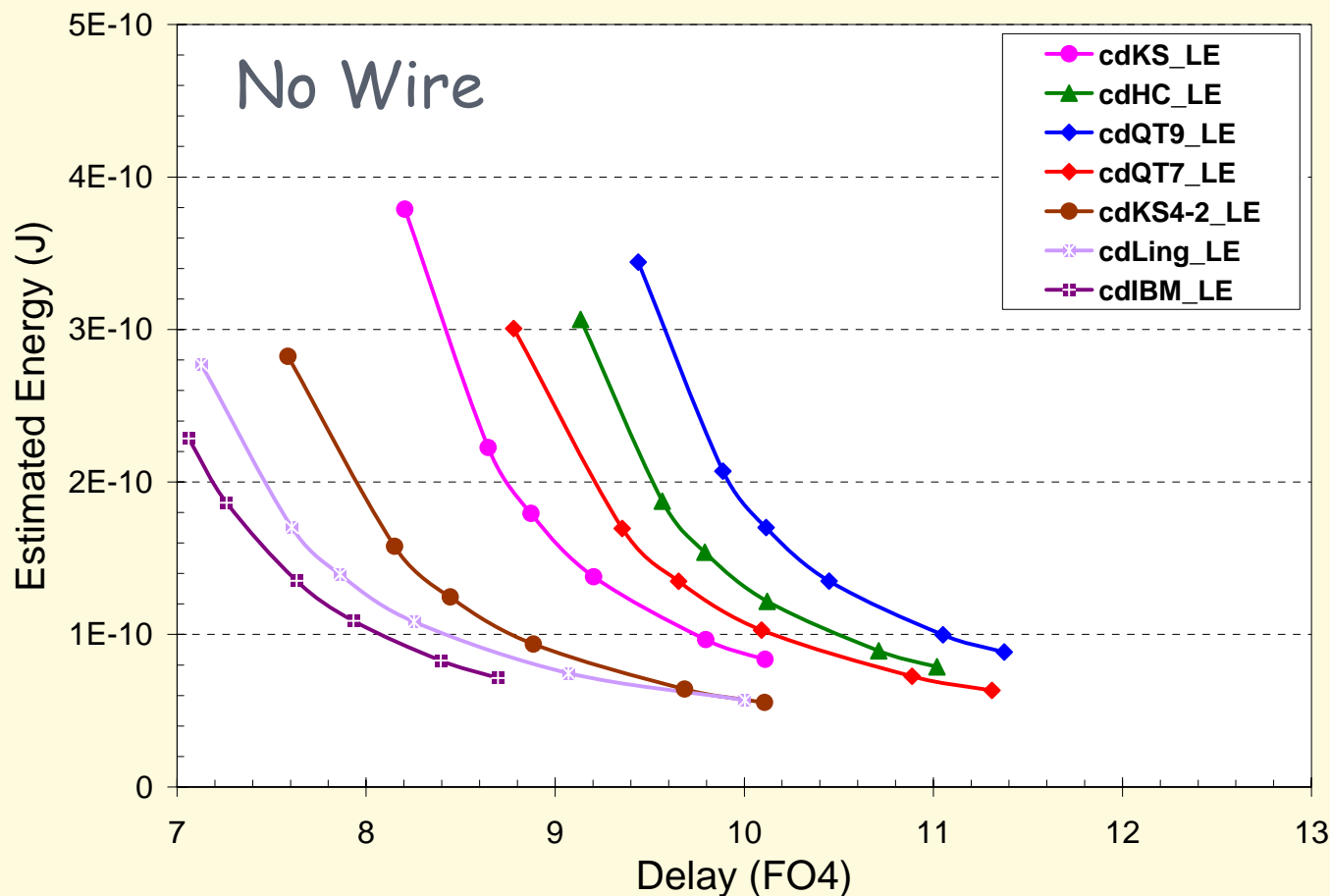- Different than what would be chosen from single point

Energy-Efficient CMOS Circuit Design

# Energy-Delay Space View
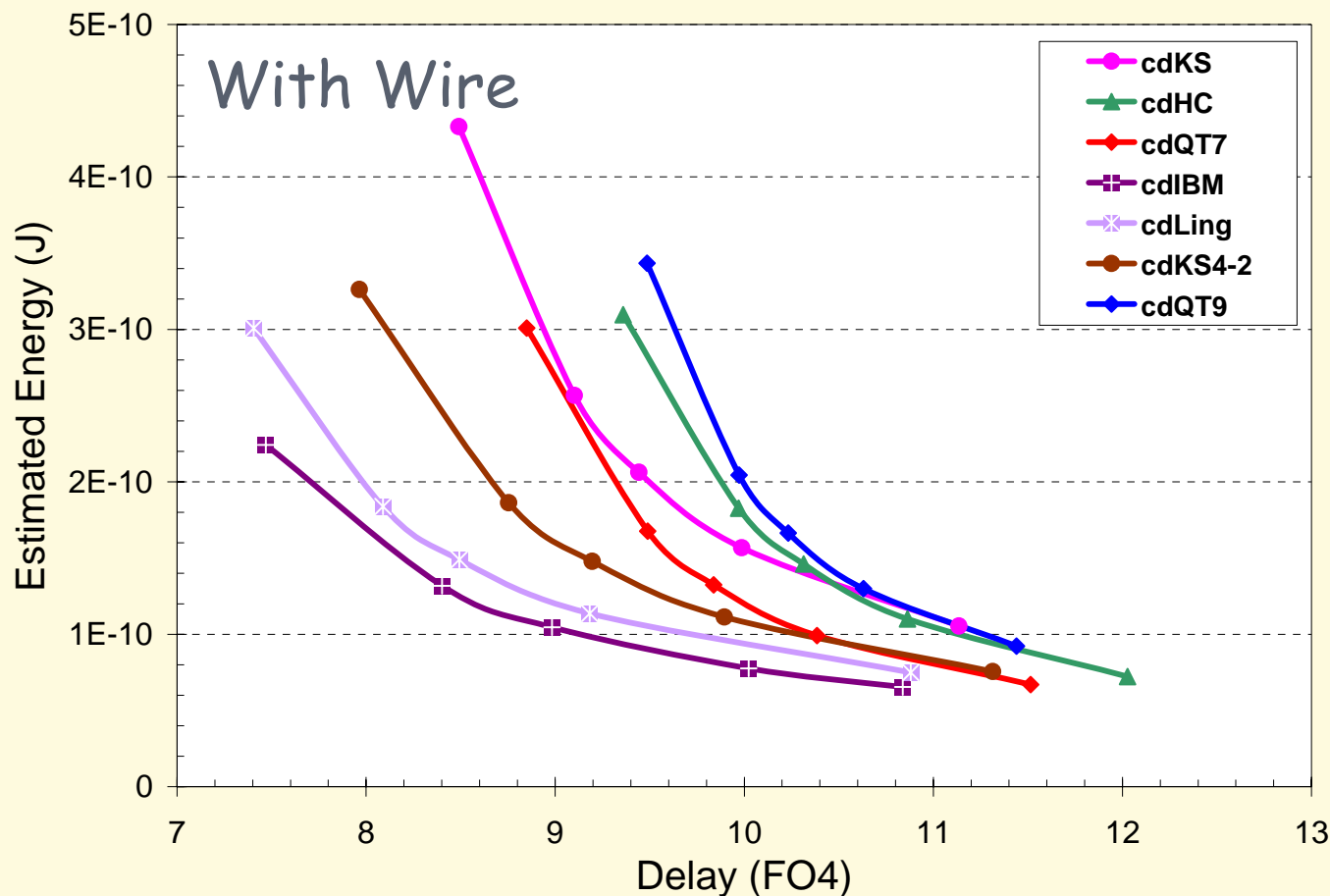


H is changing, w/ Cout=constant

- Also determines best design for Low-Power Target

# Contribution of Wire to Delay and Energy should be examined too



- Without wire, differences appear large

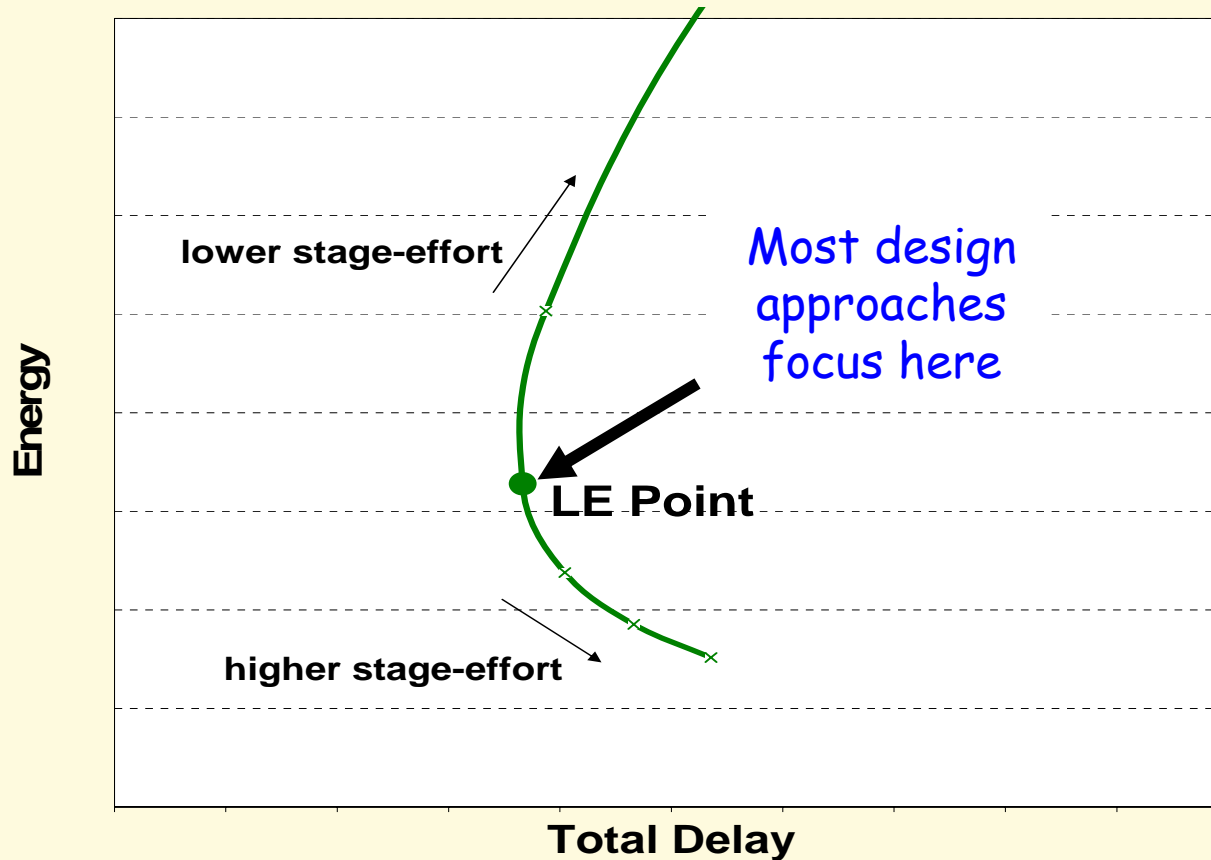# Contribution of Wire to Delay and Energy should be examined too



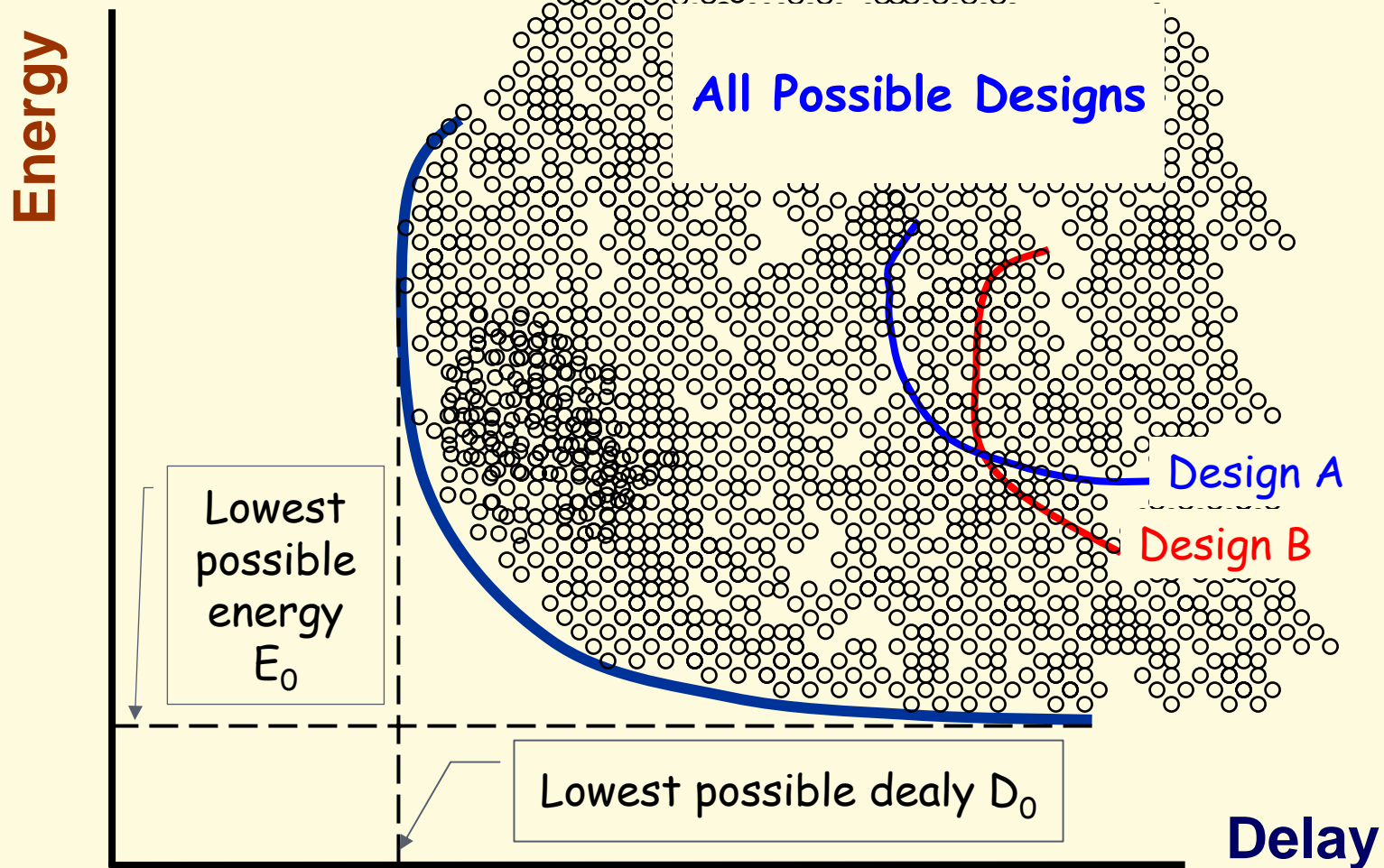- Wire strongly impacts selection of "best adders"

# Energy

# Where does Logical Effort lead us?



- It is possible to lower energy by trading delay? or ...
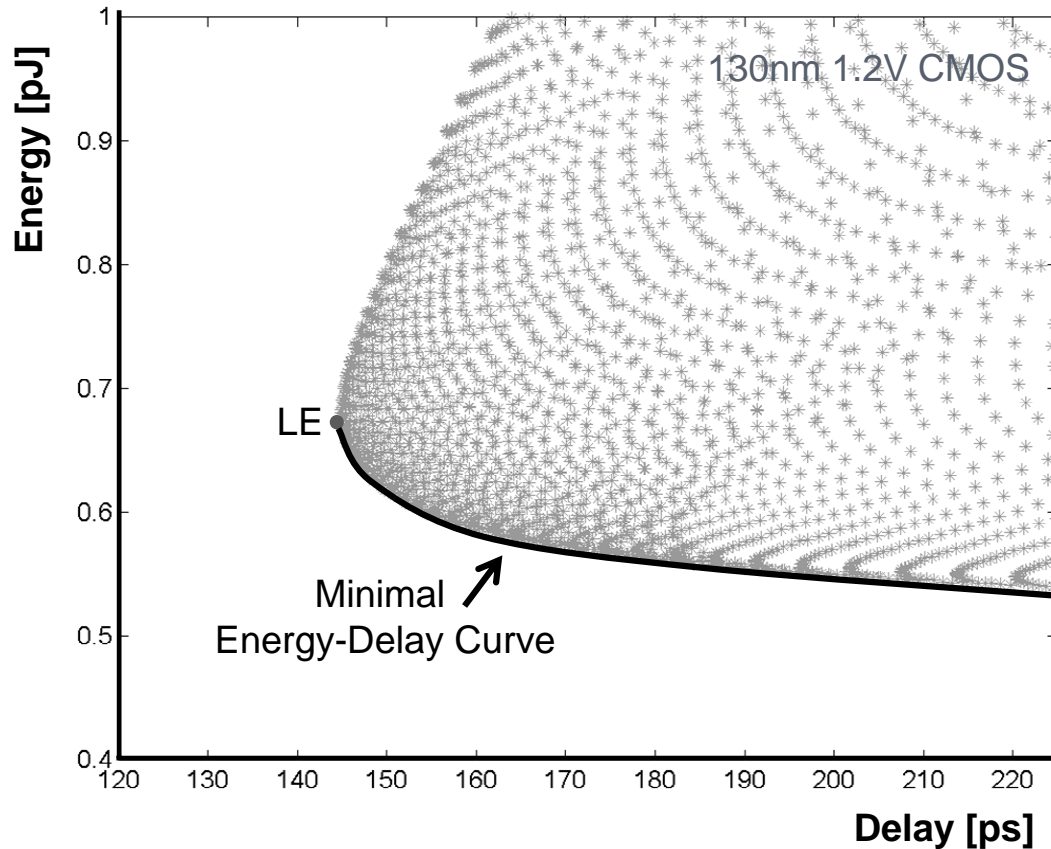
# Design in Energy-Delay Space



All Possible Designs

Design A
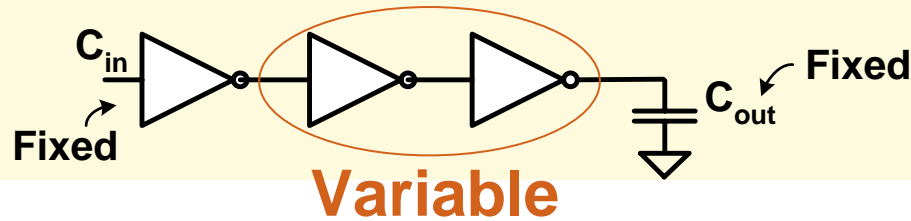
Design B

Energy

Lowest possible energy $E_0$

Lowest possible dealy $D_0$

Delay

$$(E-E_0)(D-D_0)=0.2 \times E_0 D_0$$
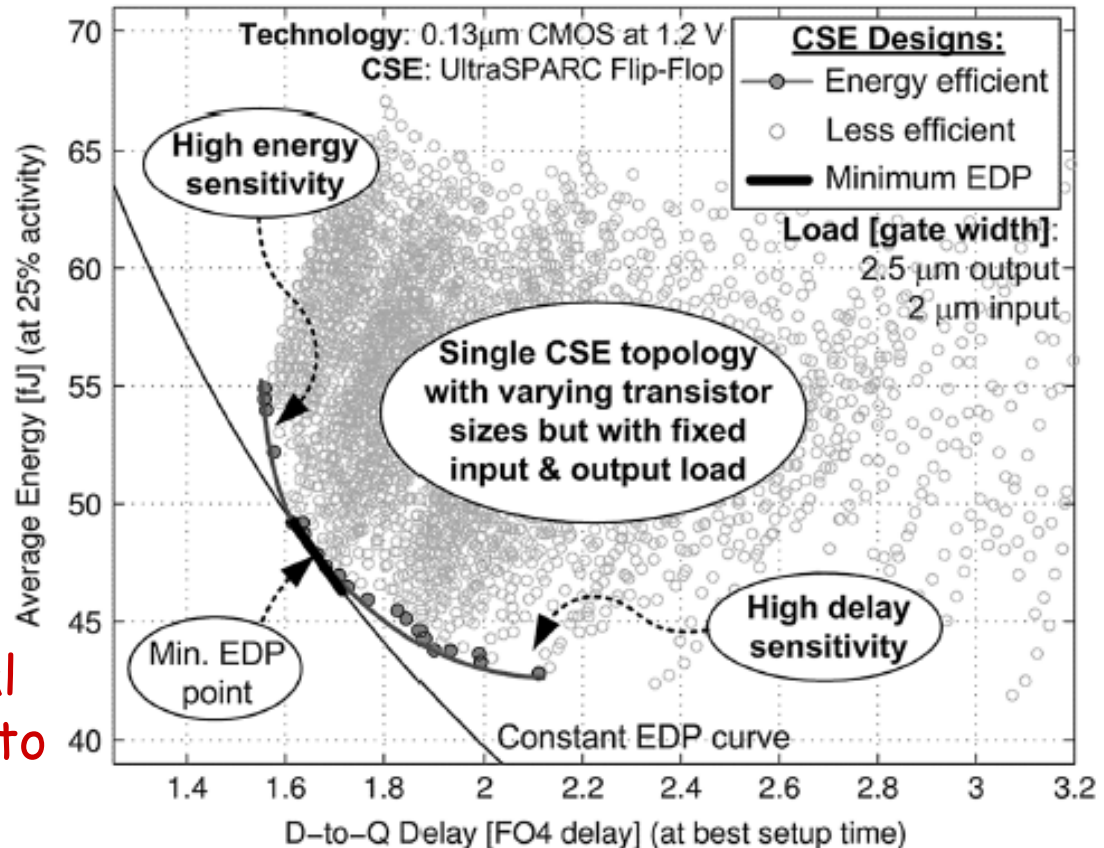
*P. Penzes, Caltech, PhD 2002, V. Zyuban, ISLPED 2002

# Energy-Delay Space of a Circuit
## (Fixed Input Size and Output Load)



$C_{in}$ Fixed

Variable

$C_{out}$ Fixed

Energy [pJ]

130nm 1.2V CMOS

LE

Minimal
Energy-Delay Curve

Delay [ps]

# Exhaustive search

A circuit with **10** transistors and **10** possible size for each transistor requires to check **$10^{10}$** possible solutions!



To obtain the minimal E-D curve do I have to check all possible solutions?
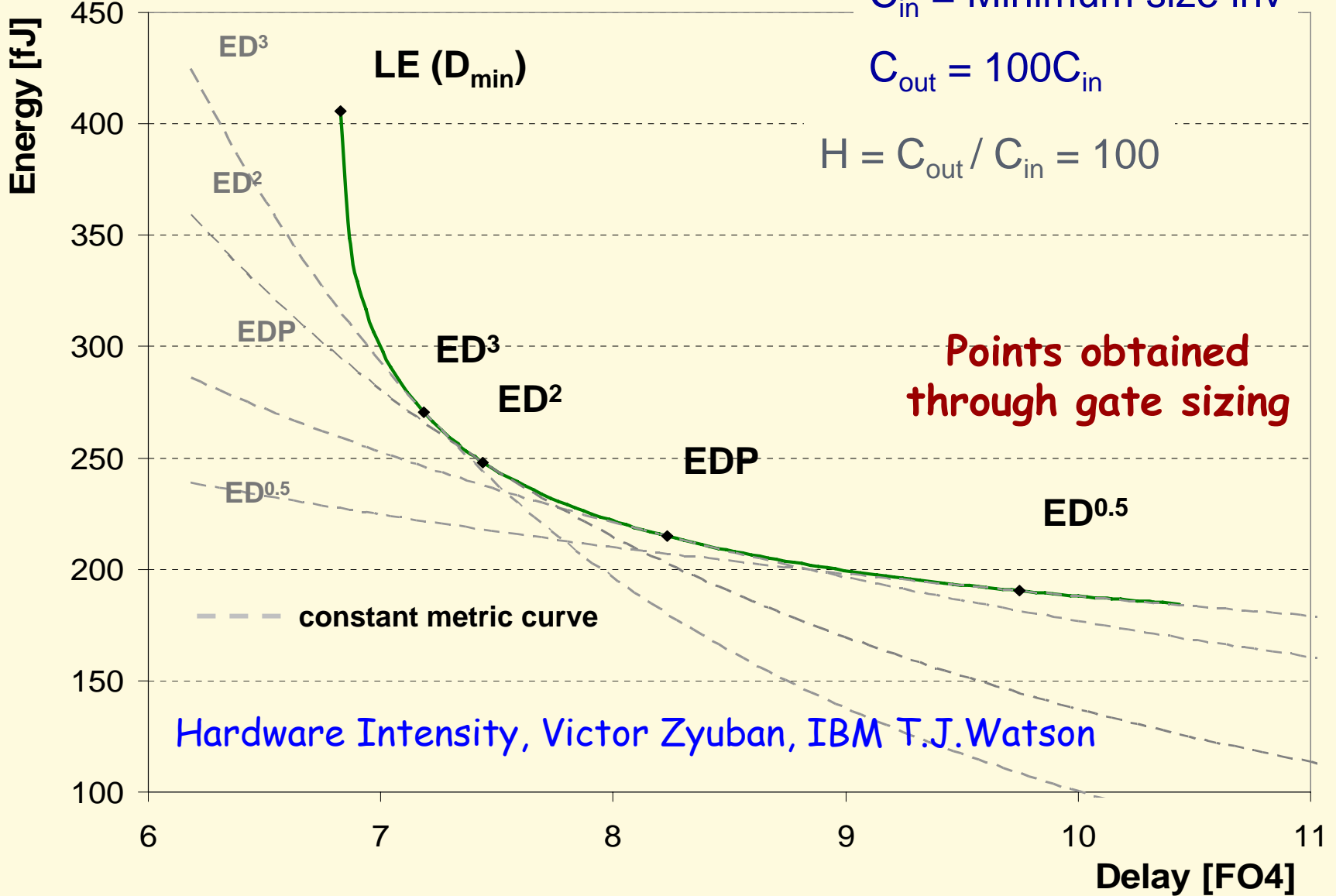
C. Giacomotto, N. Nedovic, and V. G. Oklobdzija, "*The Effect of the System Specification on the Optimal Selection of Clocked Storage Elements*", IEEE JoSSC, vol. 42, no. 6, June 2007.

# Hardware Intensity: V. Zyuban, IBM

$C_{in}$ = Minimum size inv

$C_{out}$ = 100$C_{in}$
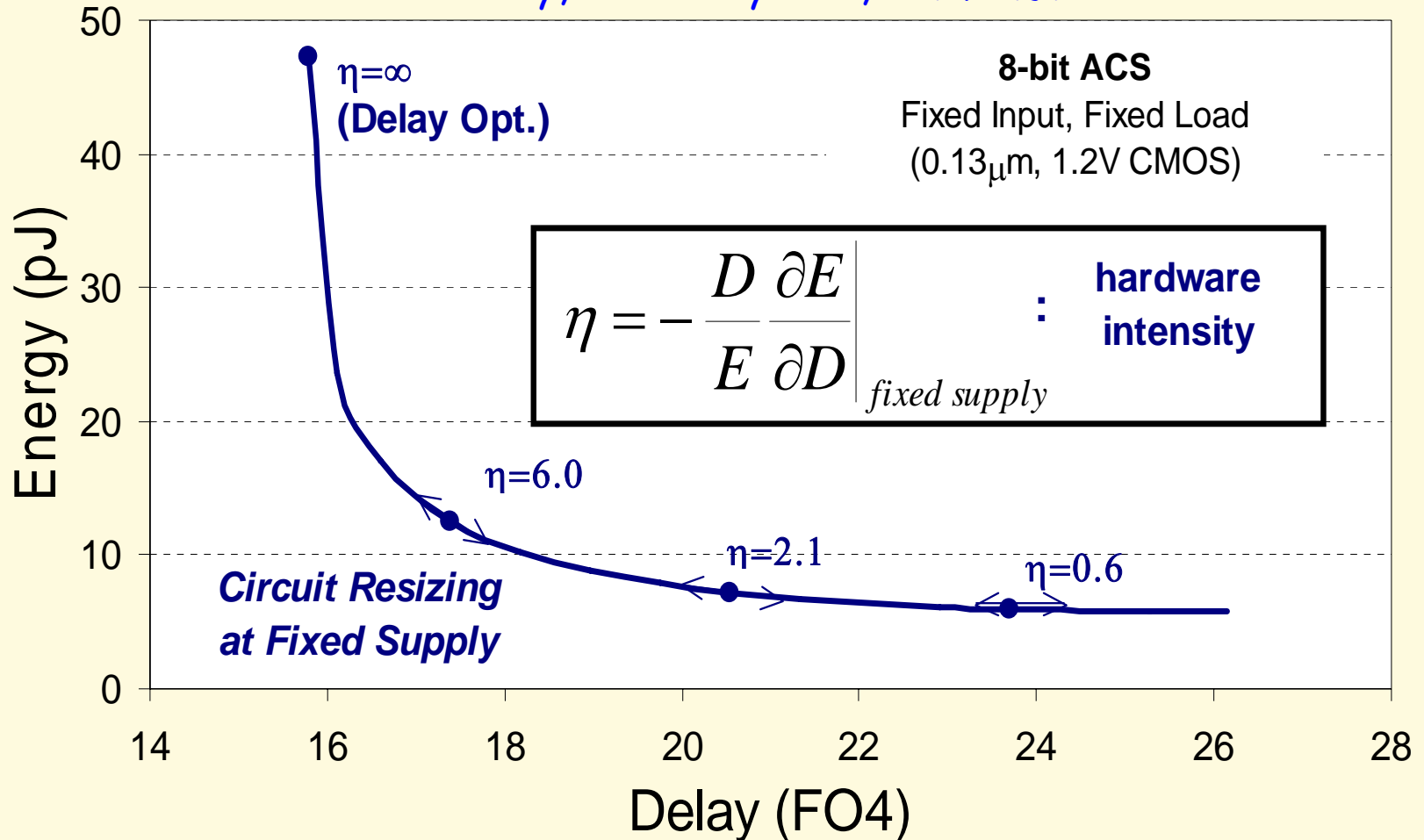
$H = C_{out} / C_{in} = 100$

**Points obtained through gate sizing**

Hardware Intensity, Victor Zyuban, IBM T.J.Watson

ED³ · ED² · EDP · ED⁰·⁵ · LE (D_min) · ED³ · ED² · EDP · ED⁰·⁵ · constant metric curve

Energy [fJ] · Delay [FO4]

# Design Objective at Nomial $V_{DD}$

## Hardware Intensity, Victor Zyuban, IBM T.J.Watson



**8-bit ACS**
Fixed Input, Fixed Load
(0.13$\mu$m, 1.2V CMOS)

$$\eta = -\frac{D}{E}\frac{\partial E}{\partial D}\bigg|_{fixed\ supply}$$

: **hardware intensity**

$\eta=\infty$
**(Delay Opt.)**

$\eta=6.0$

$\eta=2.1$

$\eta=0.6$

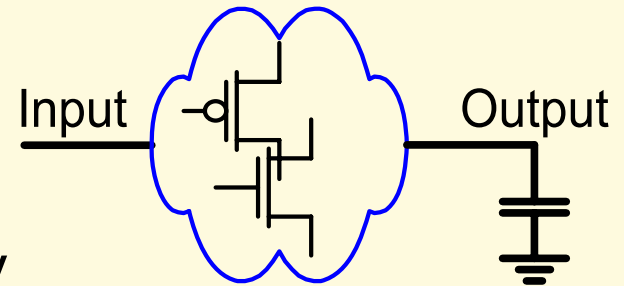*Circuit Resizing at Fixed Supply*

Energy (pJ) — Delay (FO4)

- Design choice depends on (E,D) requirements
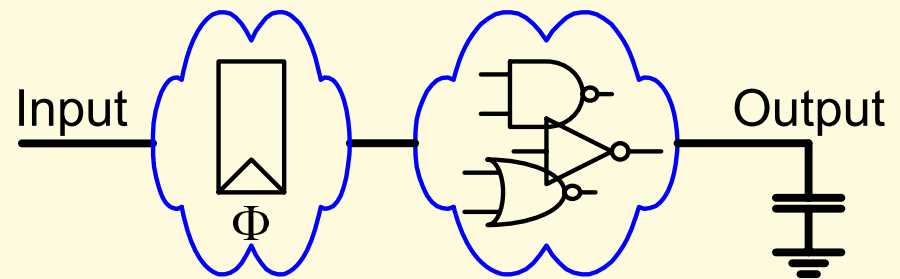
# Prior Work on Design Optimization

- **Transistor-based** [TILOS]
  - Sizing individual transistors
  - Growing complexity
  - Applicable to small blocks only



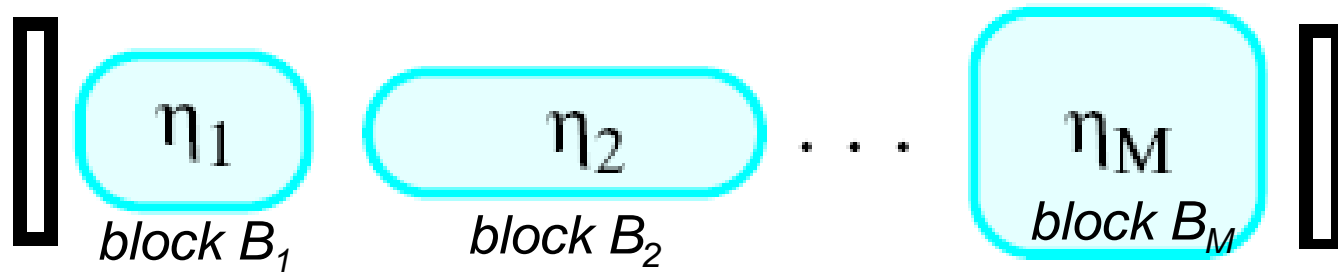- **Block-based** [Zuyban & Strenski]
  - Blocks: latch & logic
  - Trading {energy,delay} of blocks
  - CAD tools
  - Fixed interface

# Transistor-Based Approach

- Optimization problem: (i=1..M)
  - Minimize: $\quad$ Area$(W_1,...,W_M) \approx \Sigma W_i$
  - Constraint: $\quad$ $D_{worst}(W_1,...,W_M) = T$

- Delay modeling
  - Linear (RC-like): TILOS
  - Look-up table: AMPS (Synopsys)

- A convex problem $\Rightarrow$ minimal solution exists
  - Different polynomial algorithms developed
  - May have issues with convergence
  - Long run time with increasing design complexity

# Block-Based: Zyuban (IBM)



$$\eta_1 \quad \eta_2 \quad \cdots \quad \eta_M$$

block $B_1$      block $B_2$      block $B_M$

*[Zyuban, IBM T.J. Watson Research]*

- Optimization problem for a pipelined stage:
  - Minimize energy:     $E(B_1,...,B_M) = \Sigma E_{Bi}$
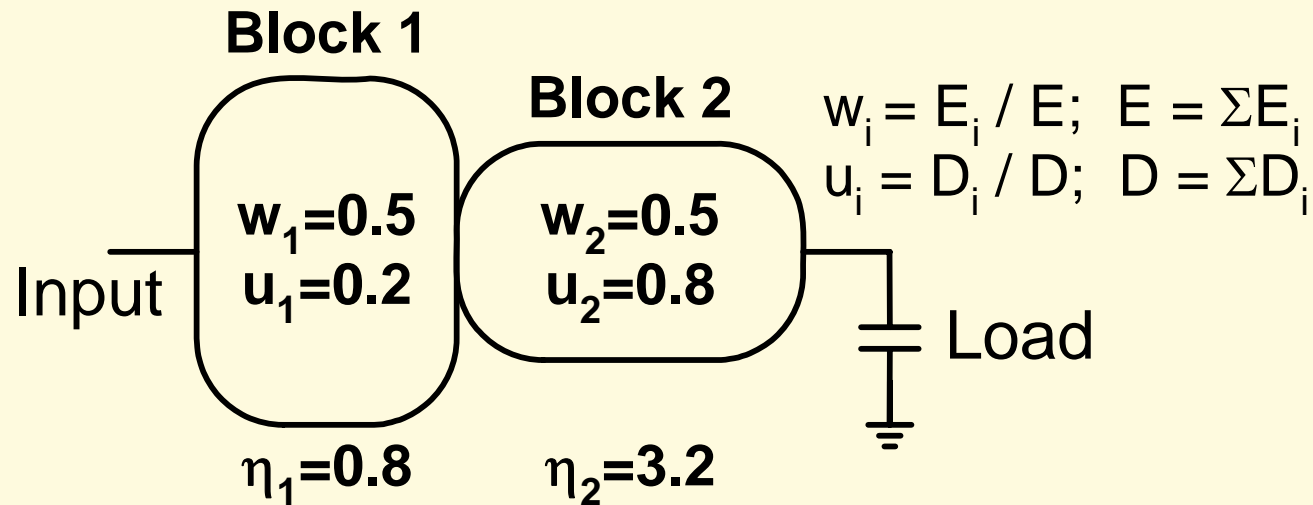  - Delay constraint:     $D(B_1,...,B_M) = \Sigma D_{Bi} = T$

- Optimal criteria:
  - $\boxed{(w_i / u_i) \cdot \eta_i = (w_k / u_k) \cdot \eta_k}$    for any (i,k)

    with   $w_x = E_{Bx}/E, \ u_x = D_{Bx}/T, \ \eta_x = (\partial E_{Bx}/E_{Bx})/(\partial D_{Bx}/D_{Bx})$
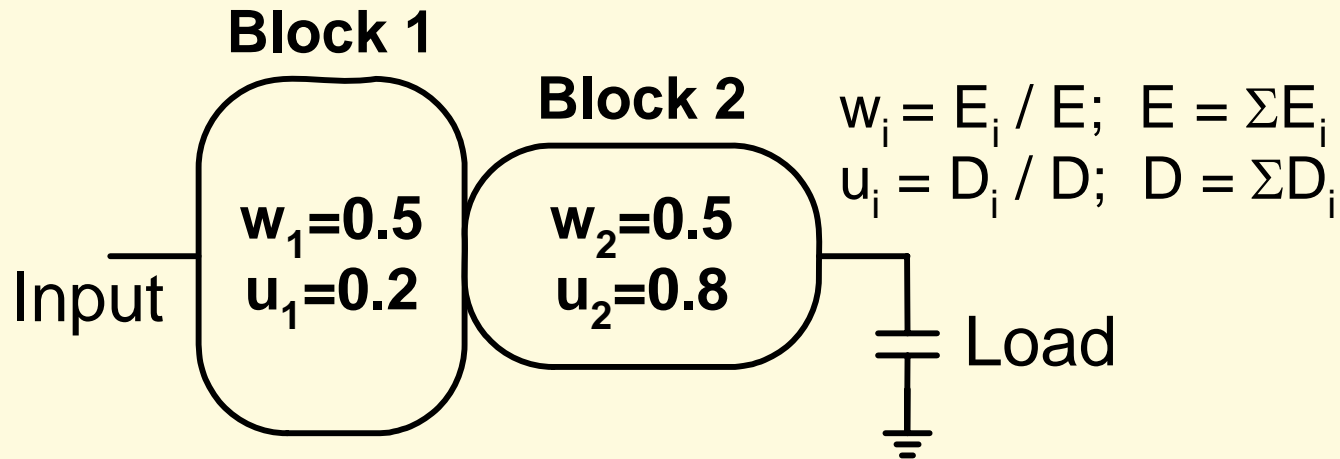
- Similar criteria for a simple pipelined system

# Application: Solution Verification

**Block 1**

**Block 2**

**w₁=0.5**
**u₁=0.2**

**w₂=0.5**
**u₂=0.8**

Input

Load

$w_i = E_i / E; \quad E = \Sigma E_i$
$u_i = D_i / D; \quad D = \Sigma D_i$

$\eta_1 = 0.8$ $\quad$ $\eta_2 = 3.2$
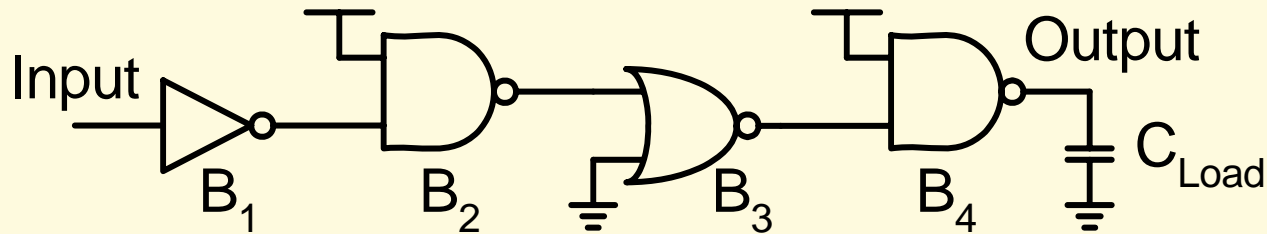
*[Zyuban, IBM T.J. Watson Research]*

- Verify optimality of solution:
  - Block 1: $(w_1/u_1) \cdot \eta_1 = 2.0$
  - Block 2: $(w_2/u_2) \cdot \eta_2 = 2.0$

  } *Equal $\Rightarrow$ optimal!*

# Application: Solution Verification

Block 1

Block 2

$w_i = E_i / E;\ \ E = \Sigma E_i$
$u_i = D_i / D;\ \ D = \Sigma D_i$

Input

$w_1 = 0.5$
$u_1 = 0.2$

$w_2 = 0.5$
$u_2 = 0.8$

Load

- If $\eta_1 = 3.2$ and $\eta_2 = 0.8$
    - Block 1: $(w_1/u_1)\cdot\eta_1 = 8.0$
    - Block 2: $(w_2/u_2)\cdot\eta_2 = 0.5$    *Unequal $\Rightarrow$ not optimal*
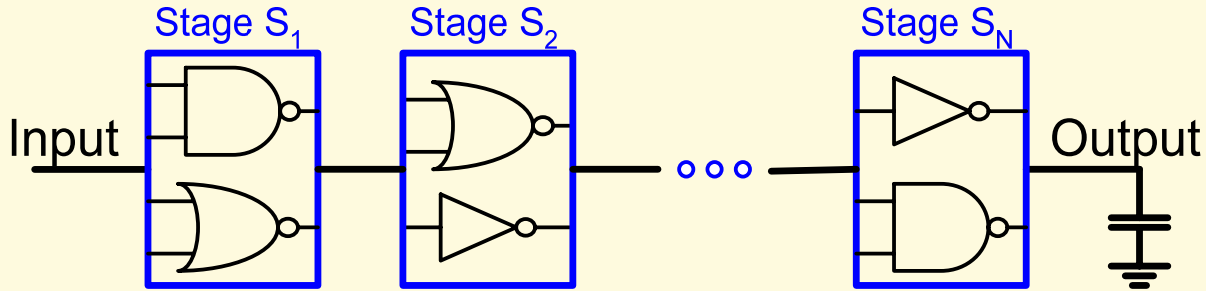    - Better solution? Relax $\eta_1$ and increase $\eta_2$

# Major Limitation



- Zyuban's assumption:
  - Delay & energy independence of each block $B_i$
- Single path: block $\equiv$ gate

$$\begin{cases} Delay\ T_d \propto \{\ C_{out}, C_{in}\ \} \\ Energy\ E \propto \{\ C_{out}, C_{in}\ \} \\ C_{in} : current\ gate\ cap \\ C_{out} : next\ gate\ cap \end{cases}$$

*: energy, delay dependency of adjacent gates*

- Similar dependency between blocks and pipelines
- No analytical solution if accounting dependency

# Proposed Stage-Based Approach



- Stage ≈ logic depth

- Gates → stage
  - Based on maximal distances to input and output
  - Stage delay: $d_{stage} = max\{d_{gate}\}$
  - Stage energy: $E_{stage} = \Sigma E_{gate}$
  - Estimated from gate energy & delay models

# Delay and Energy Modeling of Gates

- Logical Effort delay model:

reference delay $\tau$
*(technology dependent)*

$$T_{d,gate} = \left[ \frac{\kappa R_{gate} C_{gate}}{\kappa R_{inv} C_{inv}} \circ \frac{C_{out}}{C_{gate}} + \frac{\kappa R_{gate} C_{par}}{\kappa R_{inv} C_{inv}} \right] \left( \kappa R_{inv} C_{inv} \right) = \left( gh + p \right) \tau$$
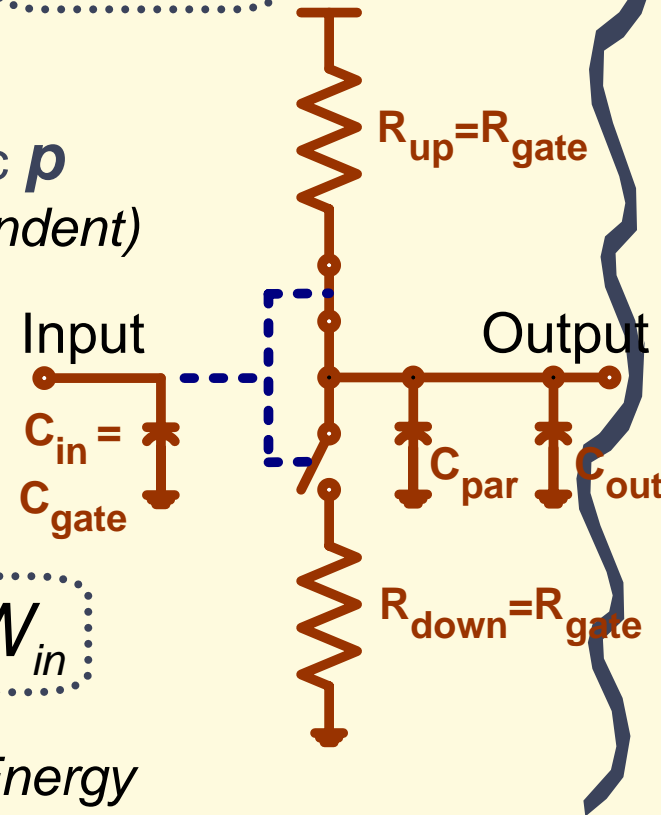
logical effort $g$
*(gate dependent)*

electrical effort $h$
*(load dependent)*

parasitic $p$
*(gate dependent)*

stage effort $f = g \bullet h$

- Energy model:

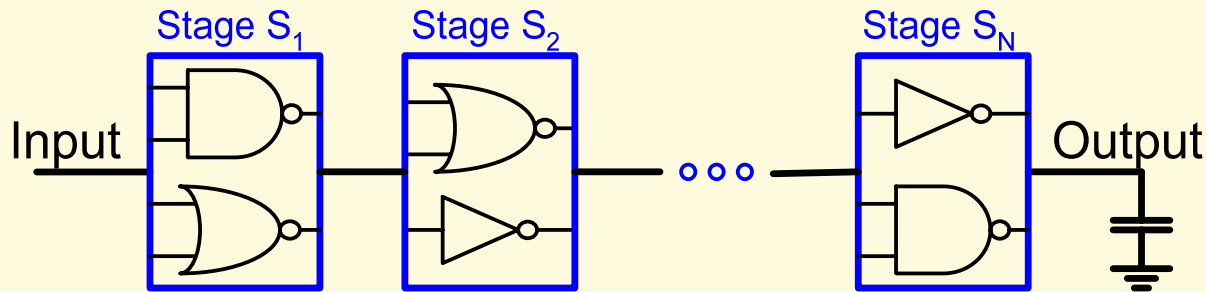$$E_{gate}(W) = E_g W_{out} + E_p W_{in} + P_l T_{CK} W_{in}$$
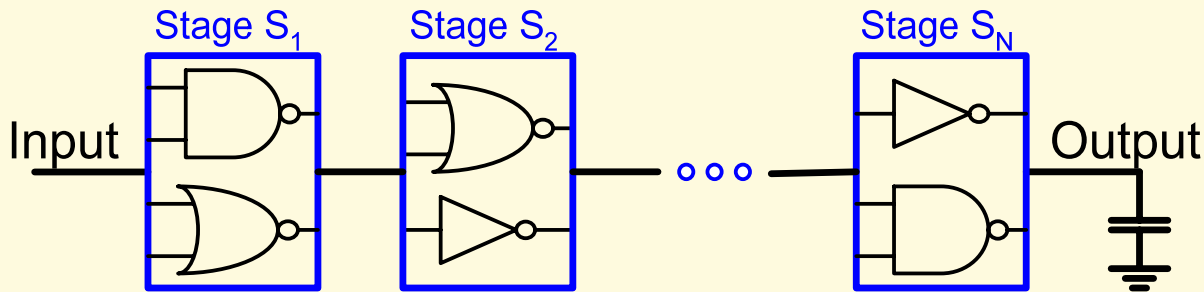
*Active Energy*      *Leakage Energy*

Input          Output

$R_{up}=R_{gate}$

$C_{in} =$
$C_{gate}$

$C_{par}$   $C_{out}$

$R_{down}=R_{gate}$

April 19, 2010

# Pipelined Stage Optimization

# Stage-Based Optimization



- **Optimization functions**
  - Delay: $\qquad D = \Sigma D_{Stage(i)}$
  - Energy: $\qquad E = \Sigma E_{Stage(i)}$
- **Possible design constraints**
  - Delay target, D
  - Input size, $W_{input}$
  - Output load, $C_{load}$
- **Posynomial Problems**
  - Solvable with polynomial algorithms

# Problem A: Delay Optimization



- Optimization problem
  - Minimize:     $D = \Sigma D_{Si}$
  - Constraint:   {Input, Load} = const.

- Objectives
  - Obtain minimally achievable delay, $D_{min}$
  - Wanted in performance-critical designs
  - Disregard energy consumption (*actually, $\partial E_i / \partial D_i = \infty$*)

# Single Path: Logical Effort



$$D = \sum_{i=1}^{N} \left( g_i \frac{C_{i+1}}{C_i} + p_i \right)$$

$C_i$ = input cap of gate i
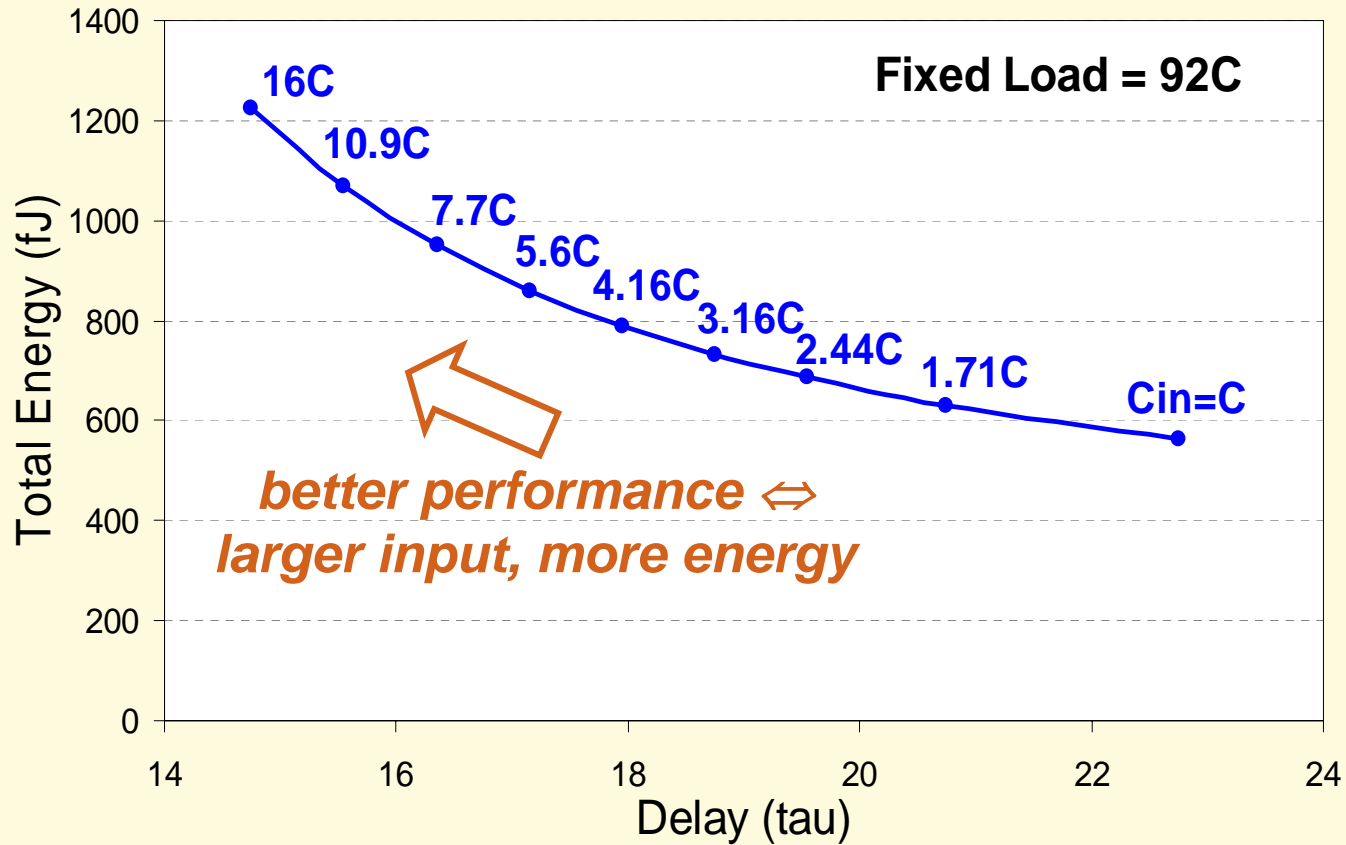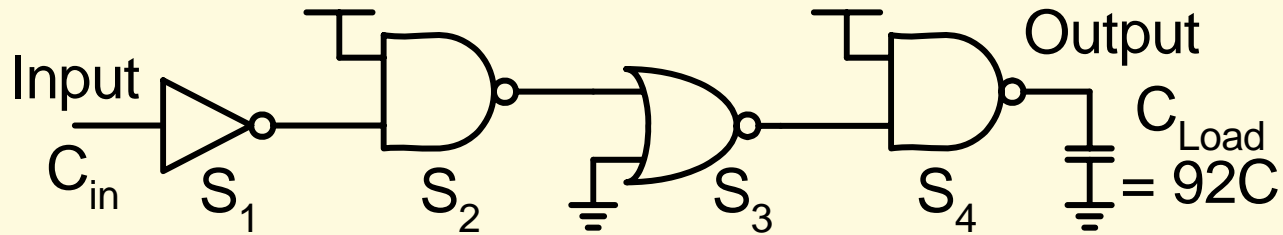$C_{in} = C_1$ : fixed
$C_{N+1} = C_{Load}$ : fixed

$dD/dC_i = 0 \ (\forall i = 1..N):$

$$f_i = f_{opt} = \left[ \left( \prod_{i=1}^{N} g_i \right) \left( \frac{C_{Load}}{C_{in}} \right) \right]^{1/N}$$
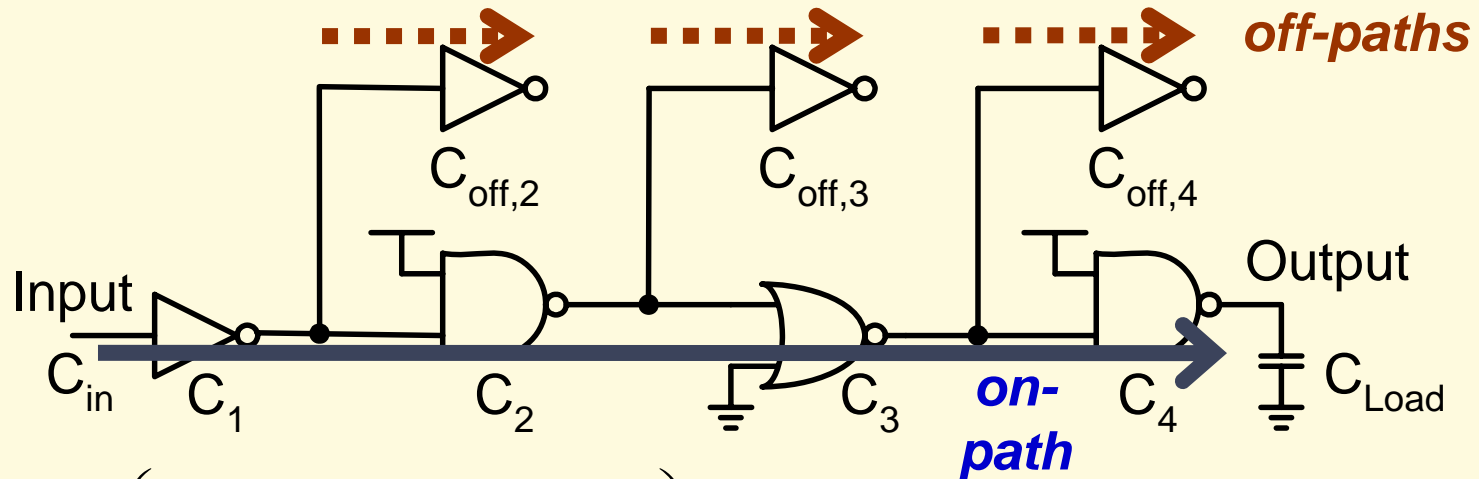
$$D_{min} = N \, f_{opt} + \sum_{i=1}^{N} p_i$$

⇧

*Path Gain, H*

- Solution = equal stage effort f  (i.e. fan-out)

# Energy Cost vs. Total Delay



April 19, 2010

# Multi-Path Circuits



$$D = \sum_{i=1}^{N} \left( g_i \frac{C_{i+1} + C_{off,i+1}}{C_i} + p_i \right)$$

$$= \sum_{i=1}^{N} \left( g_i \frac{C_{i+1} + C_{off,i+1}}{C_{i+1}} \frac{C_{i+1}}{C_i} + p_i \right)$$
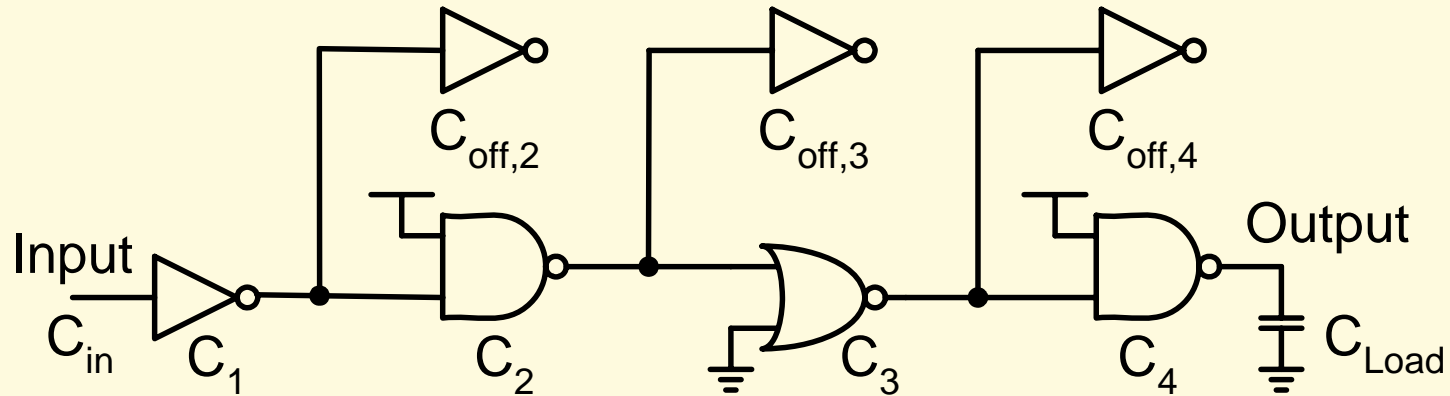
$C_i$ = *input cap of gate i*

$C_{off,i}$ = *input cap of i<sup>th</sup> off-path gate*

$b_{i+1}$ = *branching factor at (i+1)<sup>th</sup>-gate input*

• Optimal delay depends on off-path load

# Linear Branching



- Linear branching: $C_{off,i} / C_i = const.$

$dD/dC_i = 0 \ (\forall i = 1..N): \quad f_i = f_{opt} = \left[\left(\prod_{i=1}^{N} g_i\right)\left(\prod_{i=1}^{N} b_i\right)\left(\frac{C_{Load}}{C_{in}}\right)\right]^{1/N}$
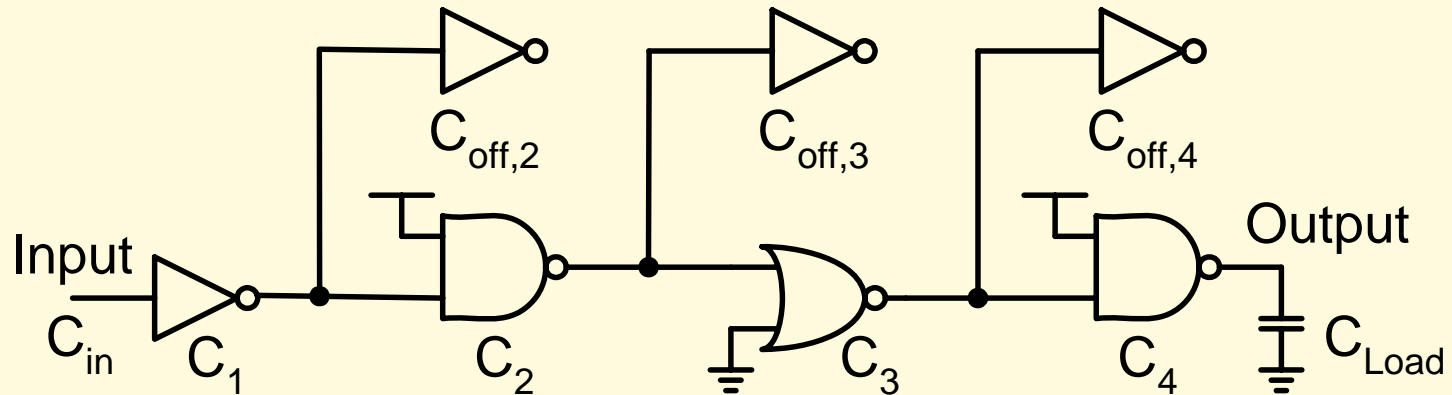
$$\boxed{D_{min} = N \, f_{opt} + \sum_{i=1}^{N} p_i}$$
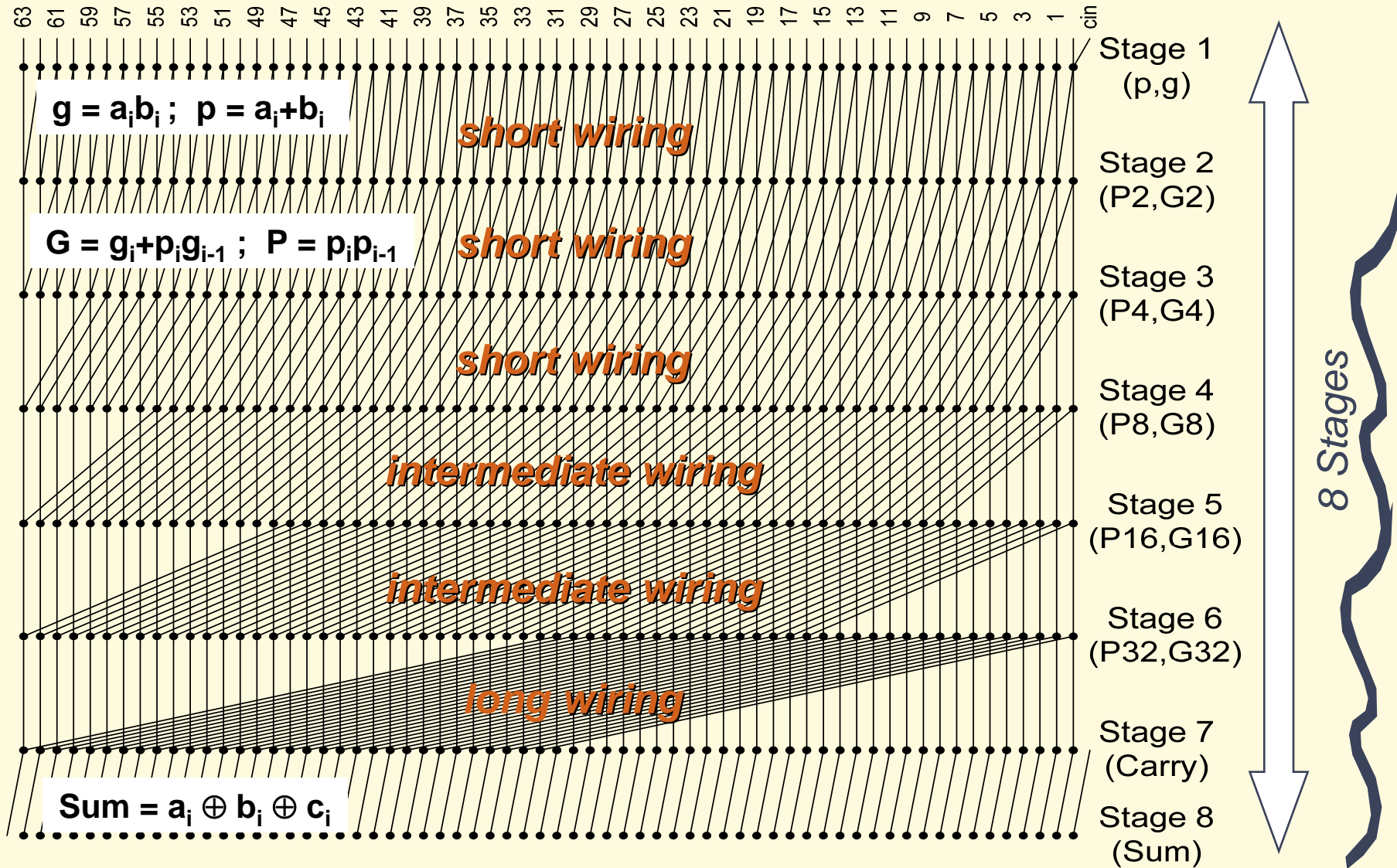
⇧ *Branching Factor, B*

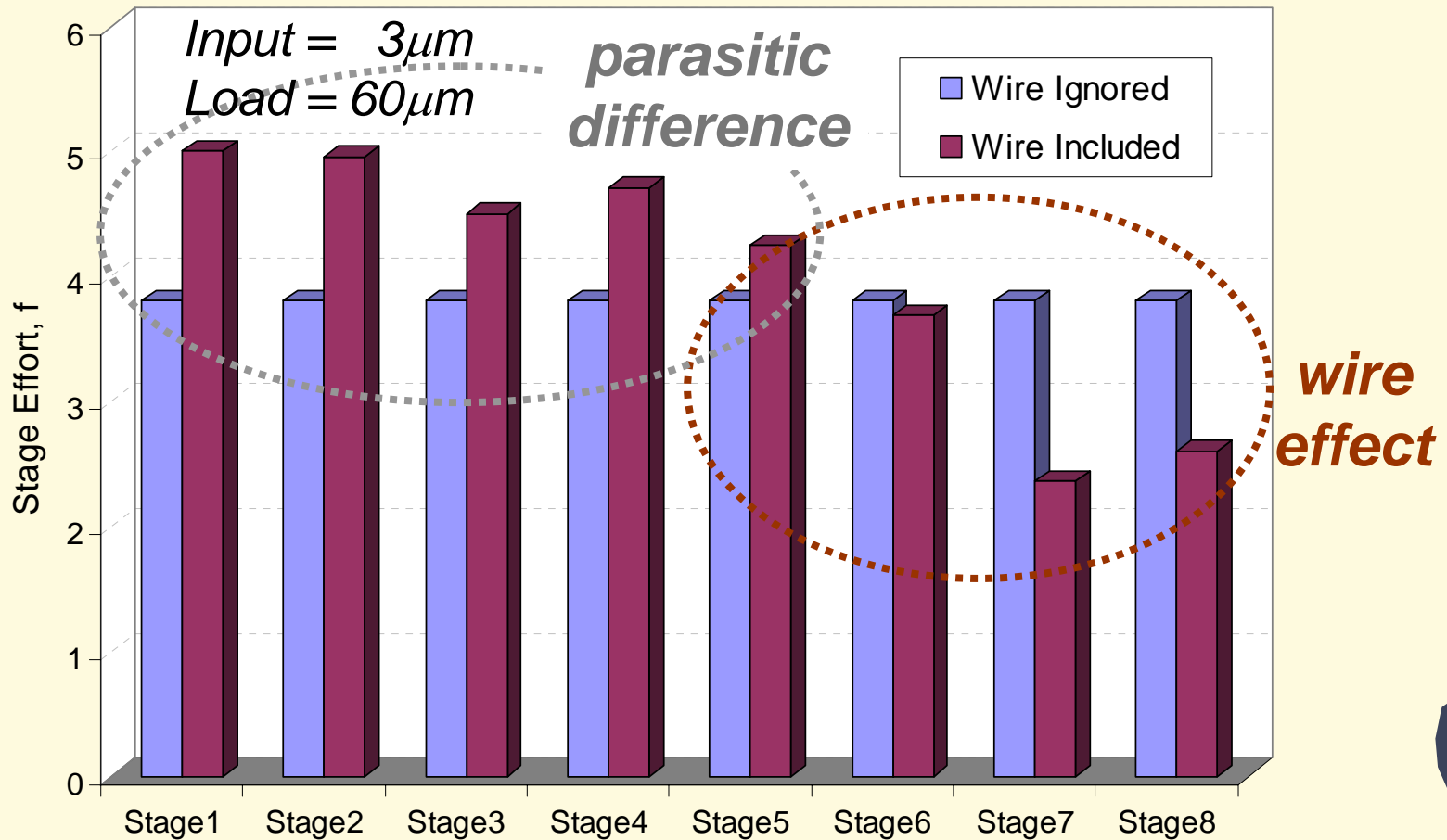- Similar analytical form for solution

# Non-Linear Branching



- Nonlinearity due to:
  - Constant off-path load (wire cap, min-size gates)
  - Unequal path lengths
  - Parasitic delay difference of gates

- *No analytical form*
  - Recursive solving
  - Solution: unequal stage efforts

# 64-bit Static Kogge-Stone Adder

63 61 59 57 55 53 51 49 47 45 43 41 39 37 35 33 31 29 27 25 23 21 19 17 15 13 11 9 7 5 3 1 cin

**Stage 1 (p,g)**

$g = a_i b_i$ ;  $p = a_i + b_i$

*short wiring*

**Stage 2 (P2,G2)**

$G = g_i + p_i g_{i-1}$ ;  $P = p_i p_{i-1}$

*short wiring*

**Stage 3 (P4,G4)**

*short wiring*

**Stage 4 (P8,G8)**

*intermediate wiring*

**Stage 5 (P16,G16)**

*intermediate wiring*

**Stage 6 (P32,G32)**

*long wiring*

**Stage 7 (Carry)**

$Sum = a_i \oplus b_i \oplus c_i$
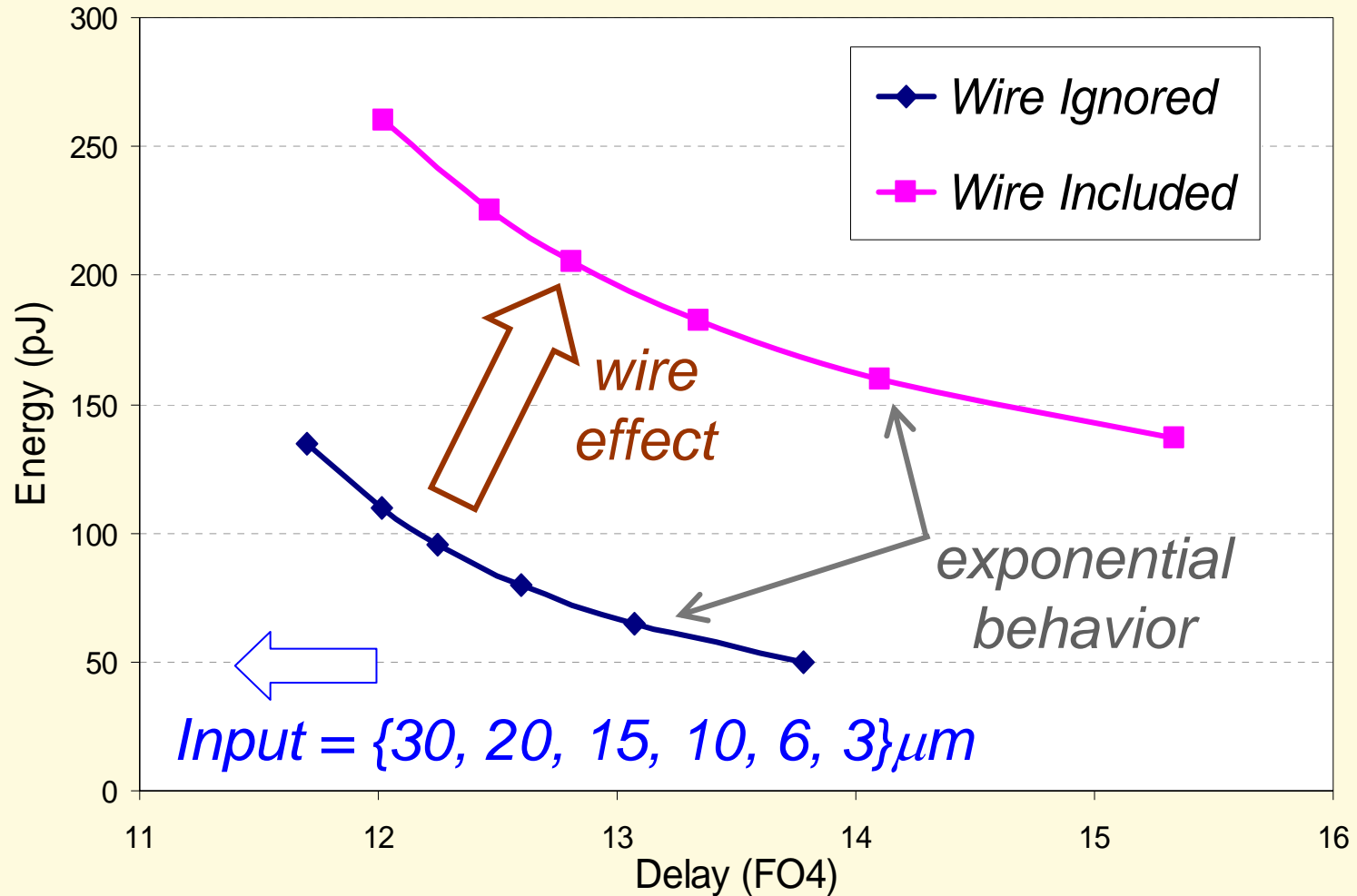
**Stage 8 (Sum)**

*8 Stages*
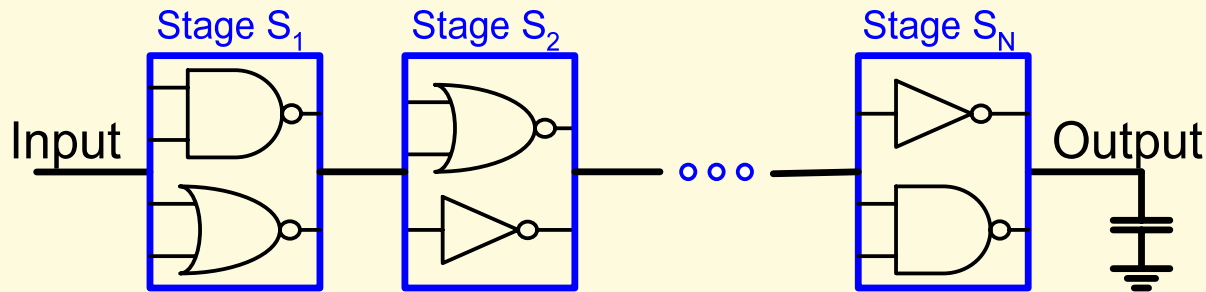
# 64-b KS: Stage Effort Distribution



- Nonlinearity causes unequal stage efforts
  - Nonlinear factors: wire load, parasitic delay diff.

# 64-b KS: Energy versus Delay
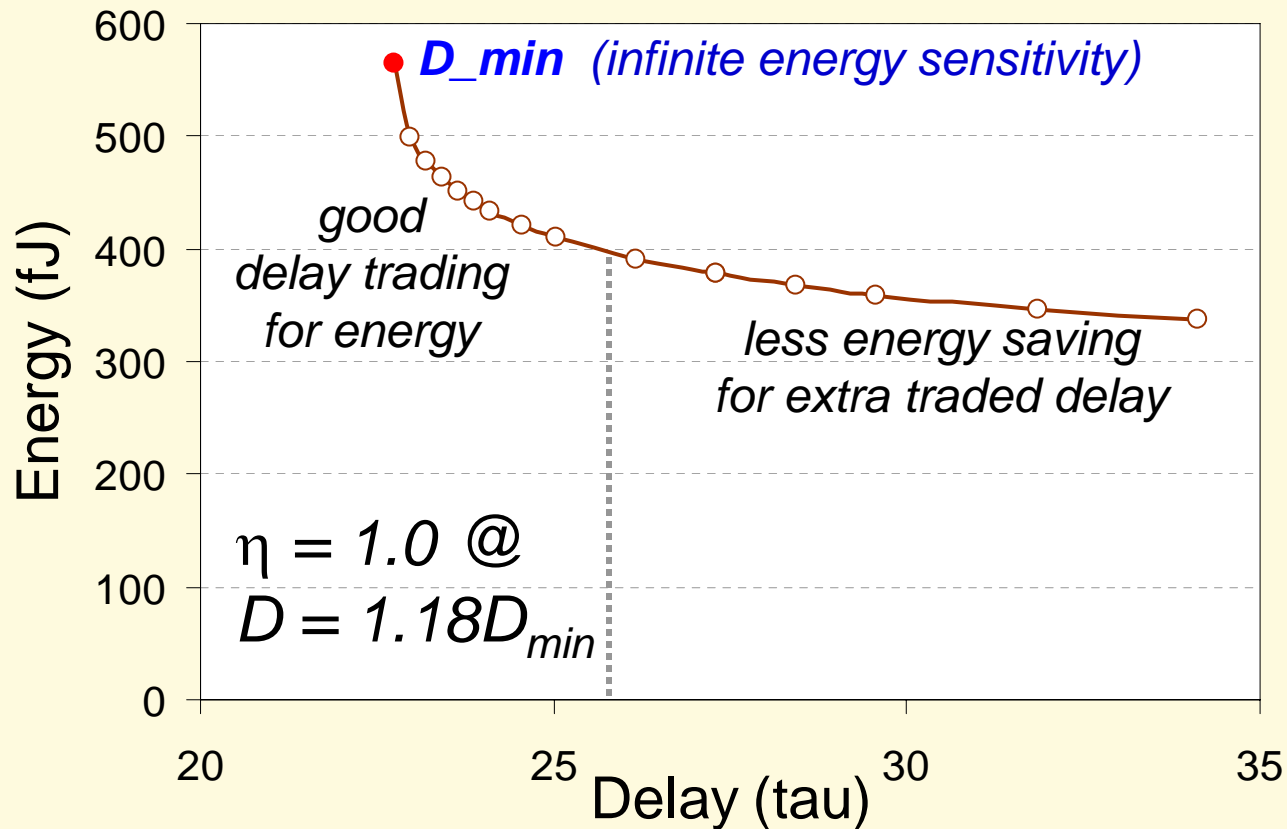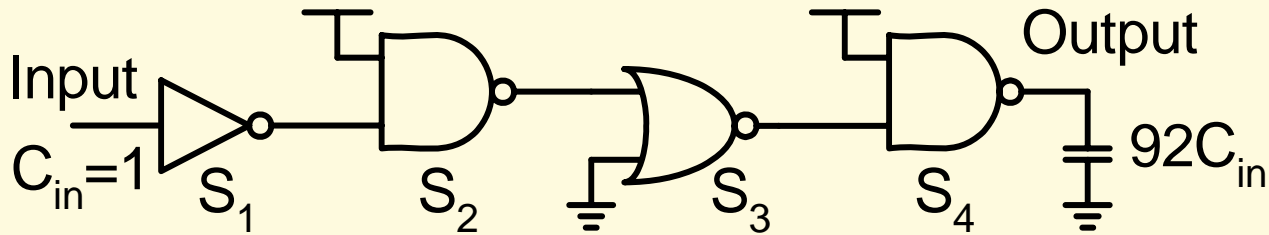


- **Significant wire effect on delay & energy**

April 19, 2010

# Problem B: Energy-Delay Tradeoff



- Optimization Problem
  - Minimize: $E = \Sigma E_{Stage(i)}$
  - *Constraint:* *{Input, Load} = const.*
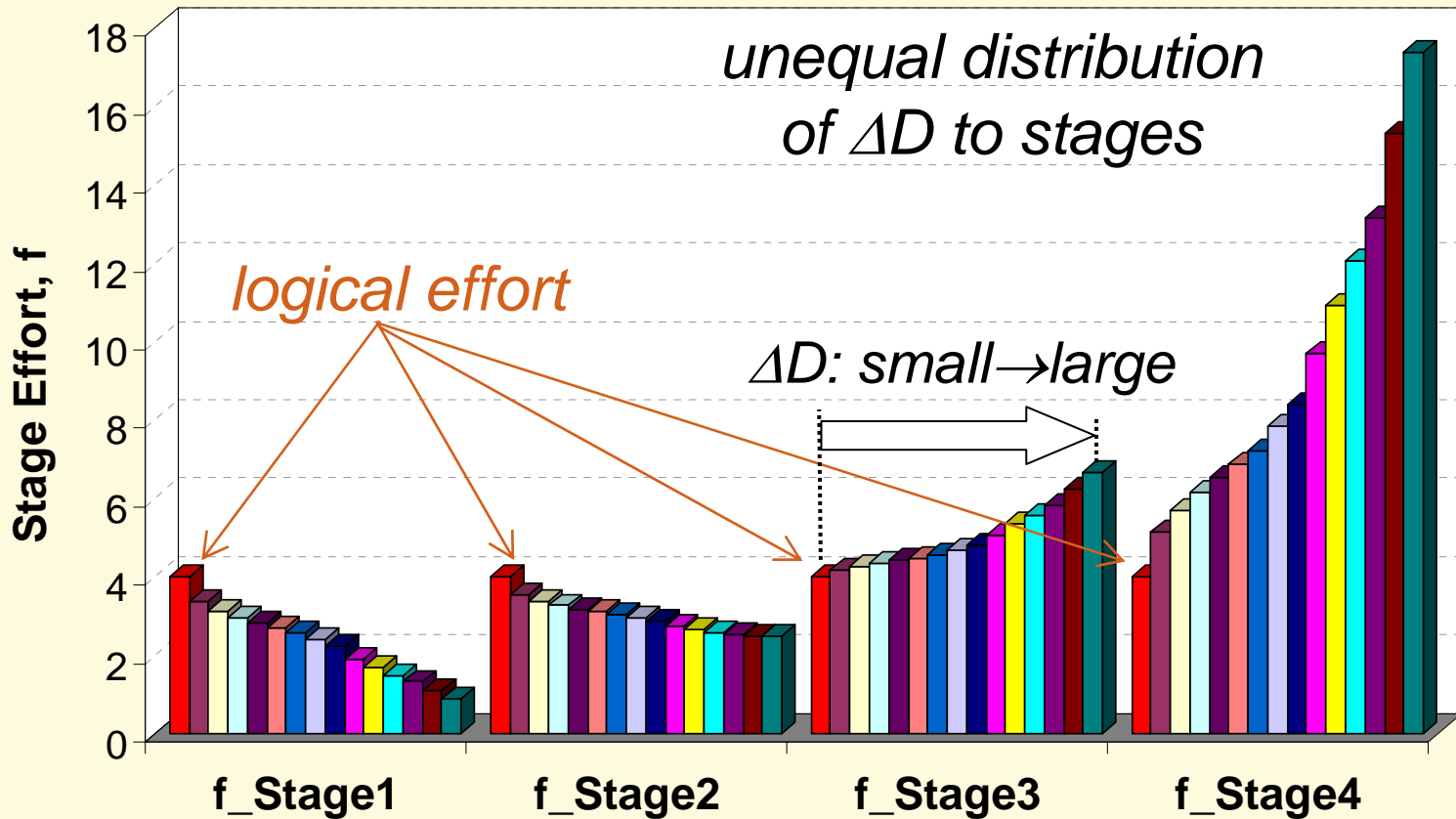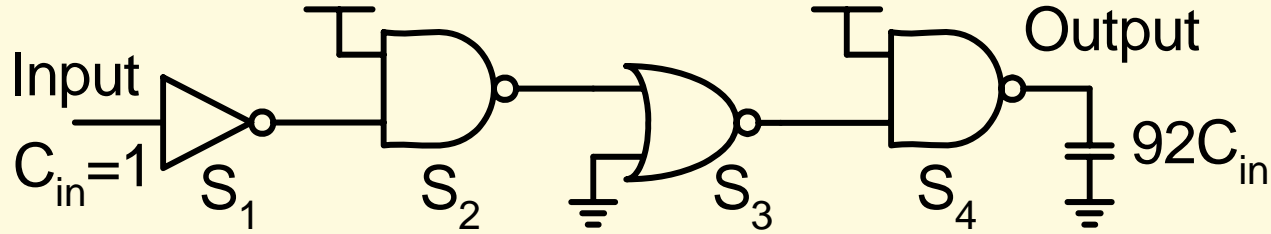    $$\Sigma D_{Stage(i)} = D_{min} + \Delta D$$

- Objectives
  - Avoid infinite energy sensitivity at $D_{min}$
  - Equalize energy-delay sens.: $(\partial E/\partial D)_{Si} = (\partial E/\partial D)_{Sk}$
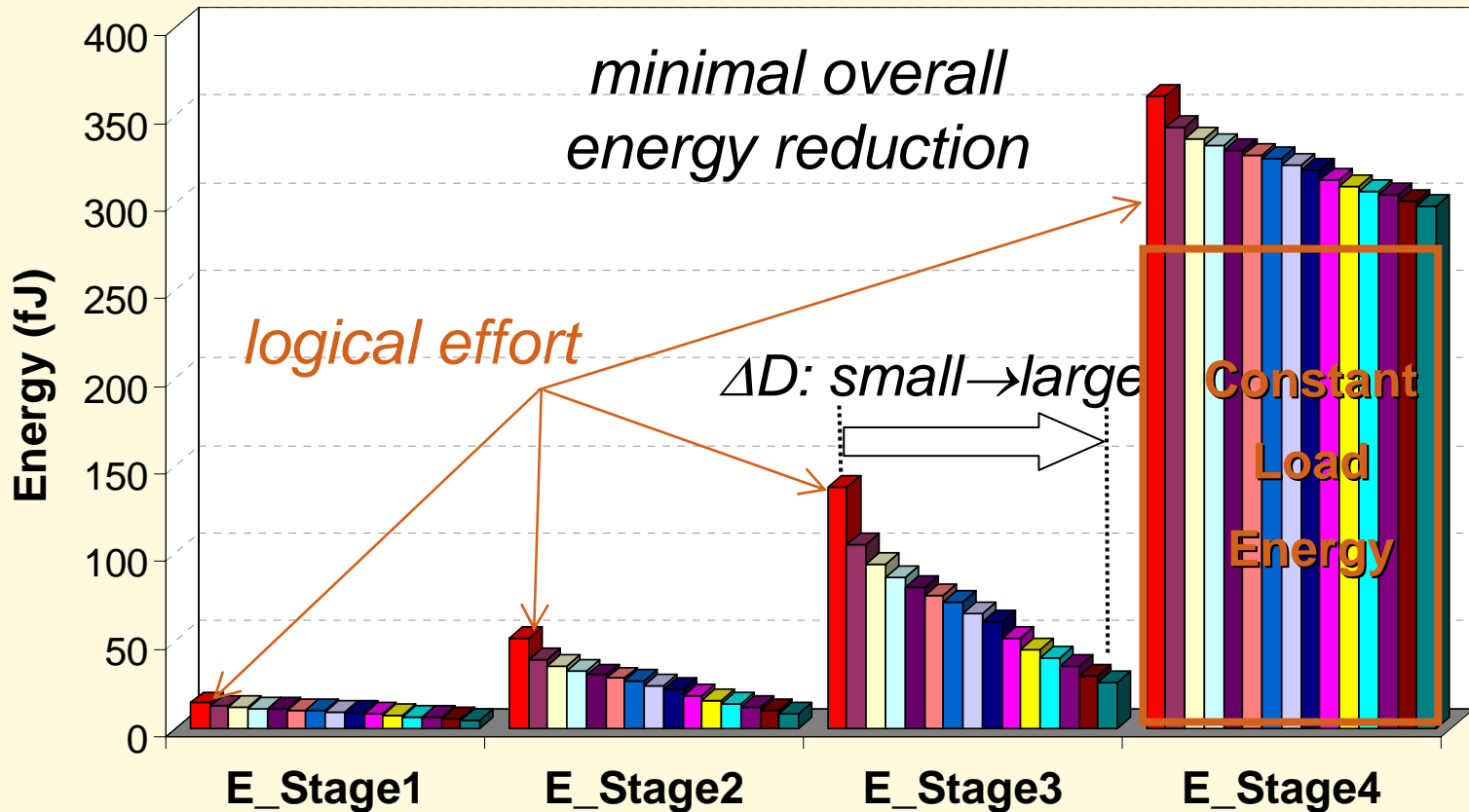  - Trade delay for energy (*traditional approach*)

# E-D Trade-off: Single Path

# E-D Trade-off: Stage Effort Distrib.
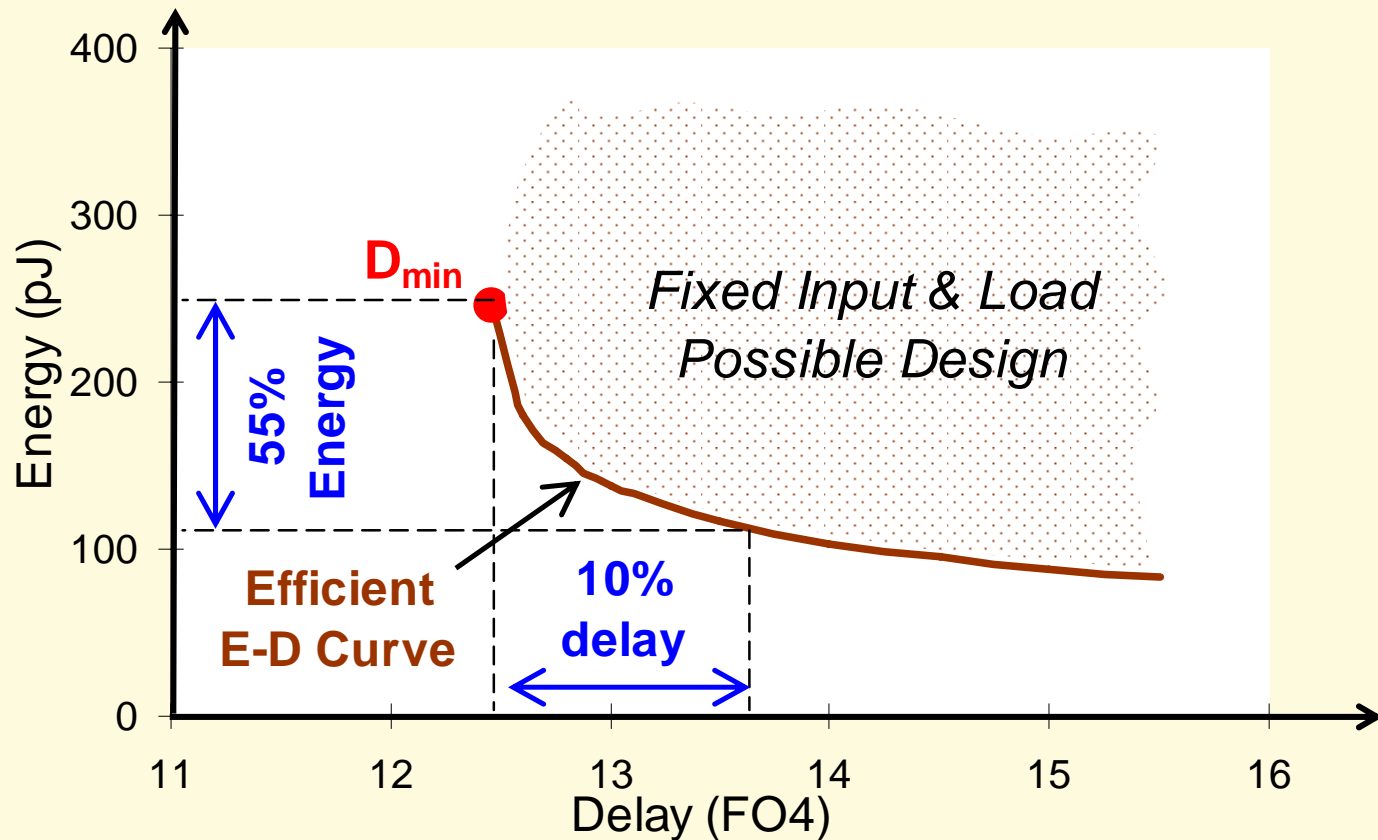
# E-D Trade-off: Energy Distrib.



April 19, 2010
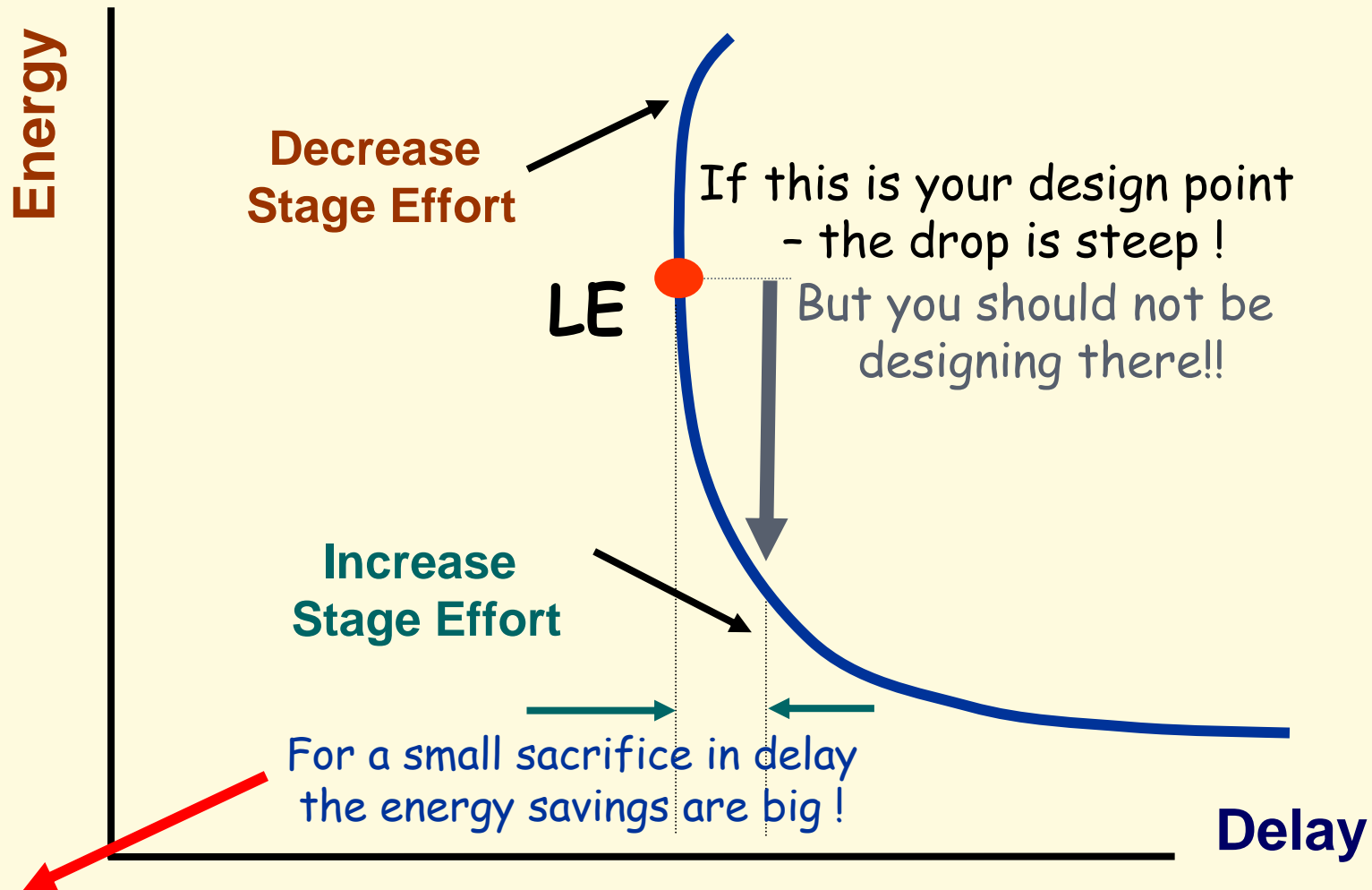
# 64-b KS: Energy Delay Tradeoff



- $D_{min}$ solution is very inefficient in energy
- 55% energy saving with 10% delay traded
- Solution = equal stage energy-delay sensitivity

# Design in Energy-Delay Space

**Energy**

**Decrease Stage Effort**

If this is your design point – the drop is steep !

But you should not be designing there!!

**LE**

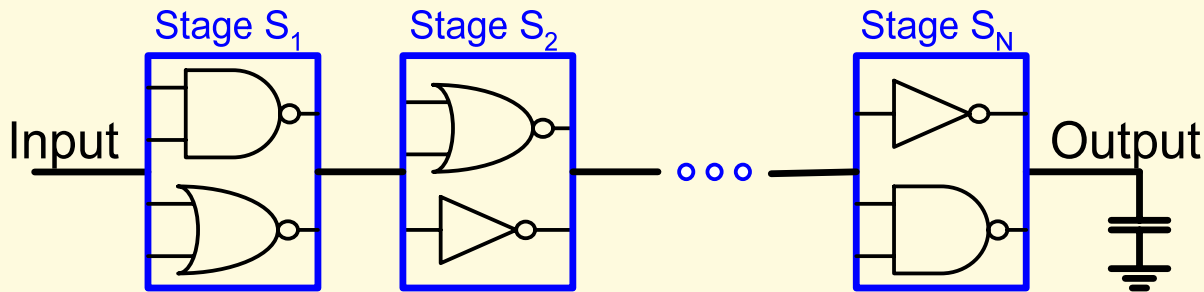**Increase Stage Effort**

For a small sacrifice in delay the energy savings are big !

**Delay**

From:
R.W. Brodersen, M.A. Horowitz, D. Markovic, B. Nikolic, V. Stojanovic, "Methods for True Power Minimization," International Conference on Computer-Aided Design, ICCAD-2002, Digest of Technical Papers, San Jose, CA, November 10-14, 2002, pp. 35-42.

# Problem C: Energy Minimization



- Problem:
  - Minimize:     $E = \Sigma E_{Si}$
  - *Constraint:   $D = \Sigma D_{Si} = const.$*

    *Load = const.*

- Objective:
  - Obtain absolute minimal energy @ given delay
  - Equalize energy-delay sens.: $(\partial E / \partial D)_{Si} = (\partial E / \partial D)_{Sk}$
  - Trade input size such that $(\partial E / \partial Input)_D = 0$

# Single Path: Stage Effort Redistrib.

# Single Path: Energy Distribution



April 19, 2010

# 64-b KS: Minimal Energy vs. Delay



- 30 - 50% energy saving @ same performance
- 1.6 - 3.6X input size

# Pipelined System Optimization

# Pipelined System Optimization

- Design constraints
  - Delay target
  - External I/O constraint

- How to obtain minimal-energy solution?
  - Pipelined stages
    - Minimized for energy
    - Sensitive to input and load variations
  - System level
    - Balancing energy sensitivities at pipelined boundaries
  - Recursive process

# Pipelined Stage: Efficient E-D Area



*Load = 60 μm gate width*

**Delay-Optimized Design Points (Logical Effort)**

30μm input

20μm input

10μm input

6μm input

4μm input

3μm input

Smaller H = $C_{out}/C_{in}$

Design Region for Possible Energy Reduction at Fixed Load

Energy Range for Designs Achieving the Same Delay

**Energy-Minimized Design Points (Lowest Energy Achievable)**

(64-bit static Kogge-Stone adder in 0.13μm CMOS at 1.2V)

Energy (pJ)

Delay (FO4)

• Efficient E-D area is upper- and lower-bounded

# Efficient Input-Delay Area



(64-bit static Kogge-Stone adder in $0.13\mu$m CMOS at 1.2V)

*Load = 60 $\mu$ m gate width*

- Energy-Minimized Design Points
- Delay-Optimized Design Points (Logical Effort)
- Range of Input Size for Achieving the Same Delay
- Possible Energy Reduction for Input Size of Design Region at Fixed Load

Input Size ($\mu$m) vs. Delay (FO4)

- • Larger/smaller input ⇔ less/more energy

# Efficient Energy-Input @ D = const.



Energy

Delay

Input

Delay

Energy

Tradable Energy

*Delay-Opt*

$$\sigma_{E,Input} = -\left.\frac{\partial E}{\partial Input}\right|_{Load=fixed}$$

**Delay = const**
*Load = const*

*Energy-Min.*

Tradable Input

Input Size

- Energy vs. input size: one-to-one relation
  – Energy sensitivity to input, $\sigma_{E,Input}$
- Flat energy in upper half of input range

# Sensitivity to Load @ D = const.

$$\sigma_{E,Load} = \frac{\partial E}{\partial Load}\bigg|_{Input=fixed}$$
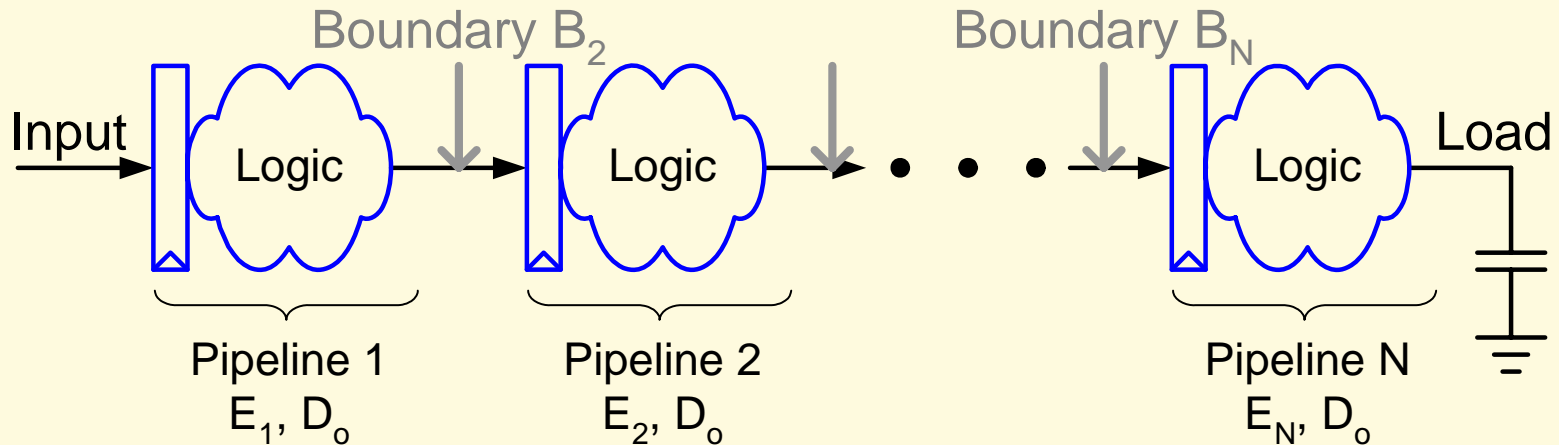


Exponential

Almost Linear

*(Fixed delay, fixed input)*

Energy

Load

$\sigma_{E,Load}$

$\sigma_{E,Input}$

Upper Boundary

Lower Boundary

Energy (pJ)

300

250

200

150

100

50

0

**Delay = const**

Load (μm)

100

80

60

40

Input (μm)

10    15    20    25    30    35

- Energy sens. to output load, $\sigma_{E,Load}$
  - More @ upper bound (delay-opt.)
  - Less @ lower bound (energy-min.)
- Energy components = energy + its sensitivities
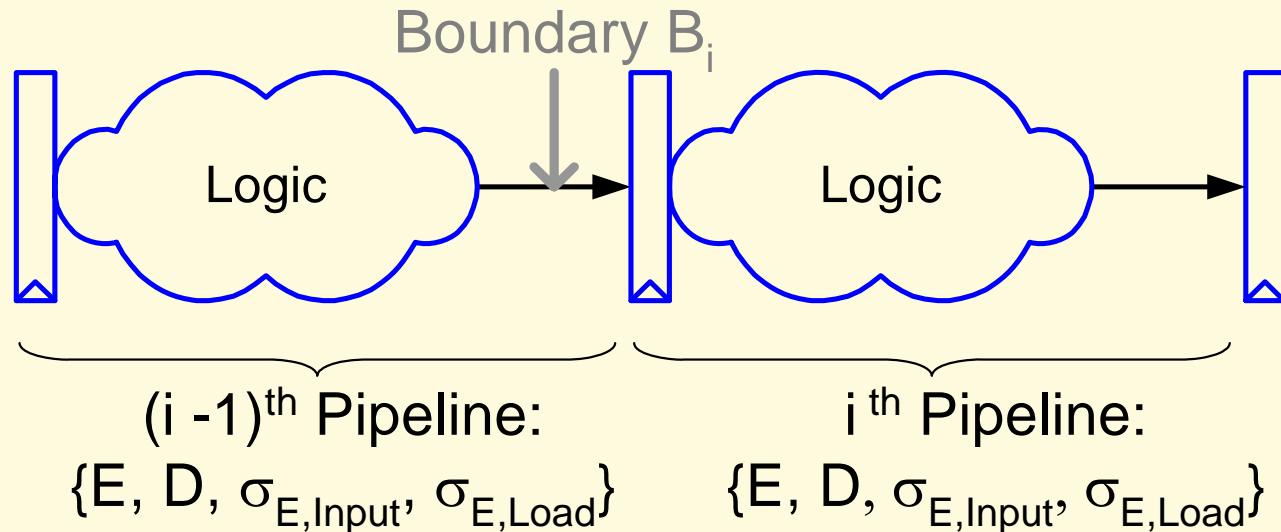
# System Energy Optimization



- **Energy minimization**
  - Pipeline: minimize energy @ given input & load
  - System: balance energy sensitivities @ boundaries

- **Trading elements: input size, output load**

- **Optimal criteria:**
$$\begin{cases} E_{Stage\,A_i} = minimal \\ \sigma_{E,Input\,A_i} = \sigma_{E,Load\,A_{i-1}} \end{cases}$$

# How to Achieve Less Total Energy

Boundary $B_i$



(i -1)$^{th}$ Pipeline:
{E, D, $\sigma_{E,Input}$, $\sigma_{E,Load}$}

i$^{th}$ Pipeline:
{E, D, $\sigma_{E,Input}$, $\sigma_{E,Load}$}

**Case A:**

□ $\sigma_{E,Input\,(i)} > \sigma_{E,Load\,(i-1)}$

- *Increase i$^{th}$ input*
  $\Rightarrow$ *less E$_{Total}$*

**Case B:**

■ $\sigma_{E,Input\,(i)} < \sigma_{E,Load\,(i-1)}$

- *Reduce i$^{th}$ input*
  $\Rightarrow$ *less E$_{Total}$*

$$E_{i-1} + E_i = min \iff \sigma_{E,Input(\,i\,)} = \sigma_{E,Load(\,i-1)}$$

# Case Study: Media Datapath

Maximal Input = $30\mu m$

| 16-b Reg X | 16-b Reg Y | 16-b Reg U | 16-b Reg V |
|---|---|---|---|

16x16 Multiplier

16x16 Multiplier

Target Delay = $17FO_4$

2 x 32-b Reg (C1,S1)

2 x 32-b Reg (C2,S2)

64-b ALU

Output Load = $60\mu m$ gate width

64-b Reg Z

- What is the minimal energy solution?

# Case Study: Optimal Criteria

*Simplified Datapath*

MAX{X,Y,U,V} = 30μm

{Z} = 60μm

Delay = 17FO$_4$

E{P1} = min.

E{P2} = min.

E{P3} = min.

$\sigma_{E,Load\ P1} + \sigma_{E,Load\ P2} = \sigma_{E,Input\ P3}$



*Boundary*

PIPELINE 1
X1, Y1
Multiplier 16x16

PIPELINE 2
U1, V1
Multiplier 16x16
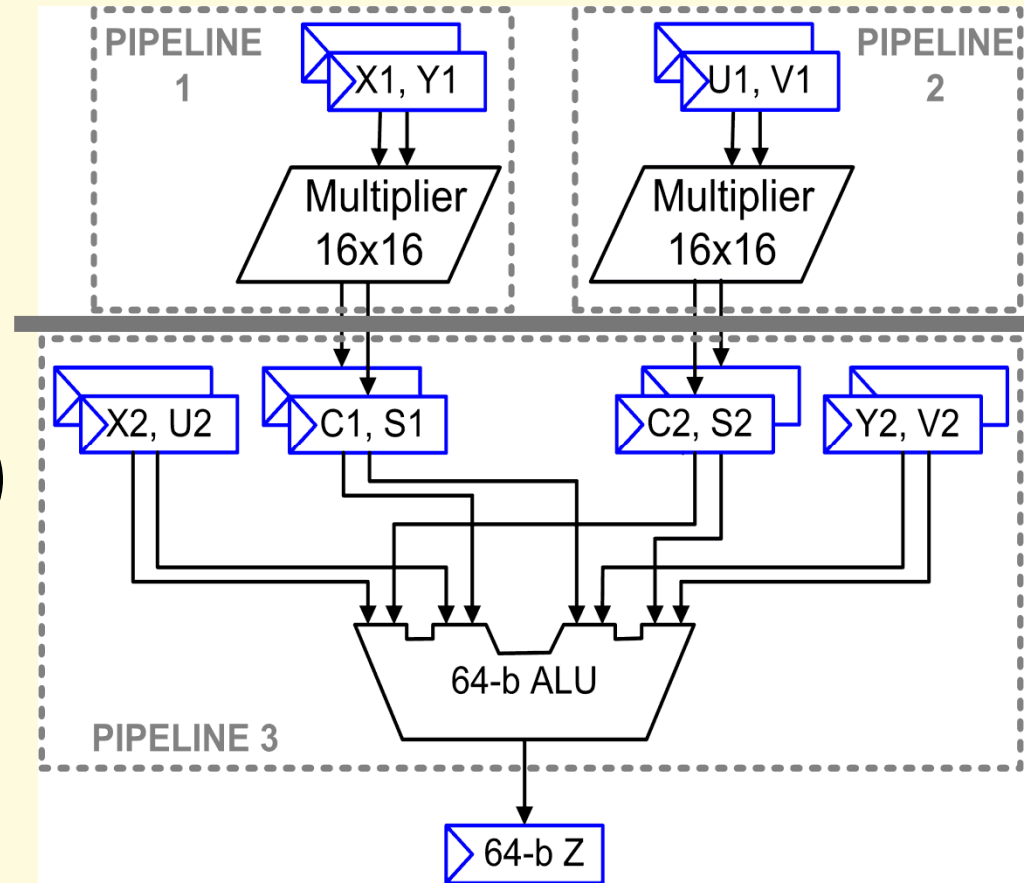
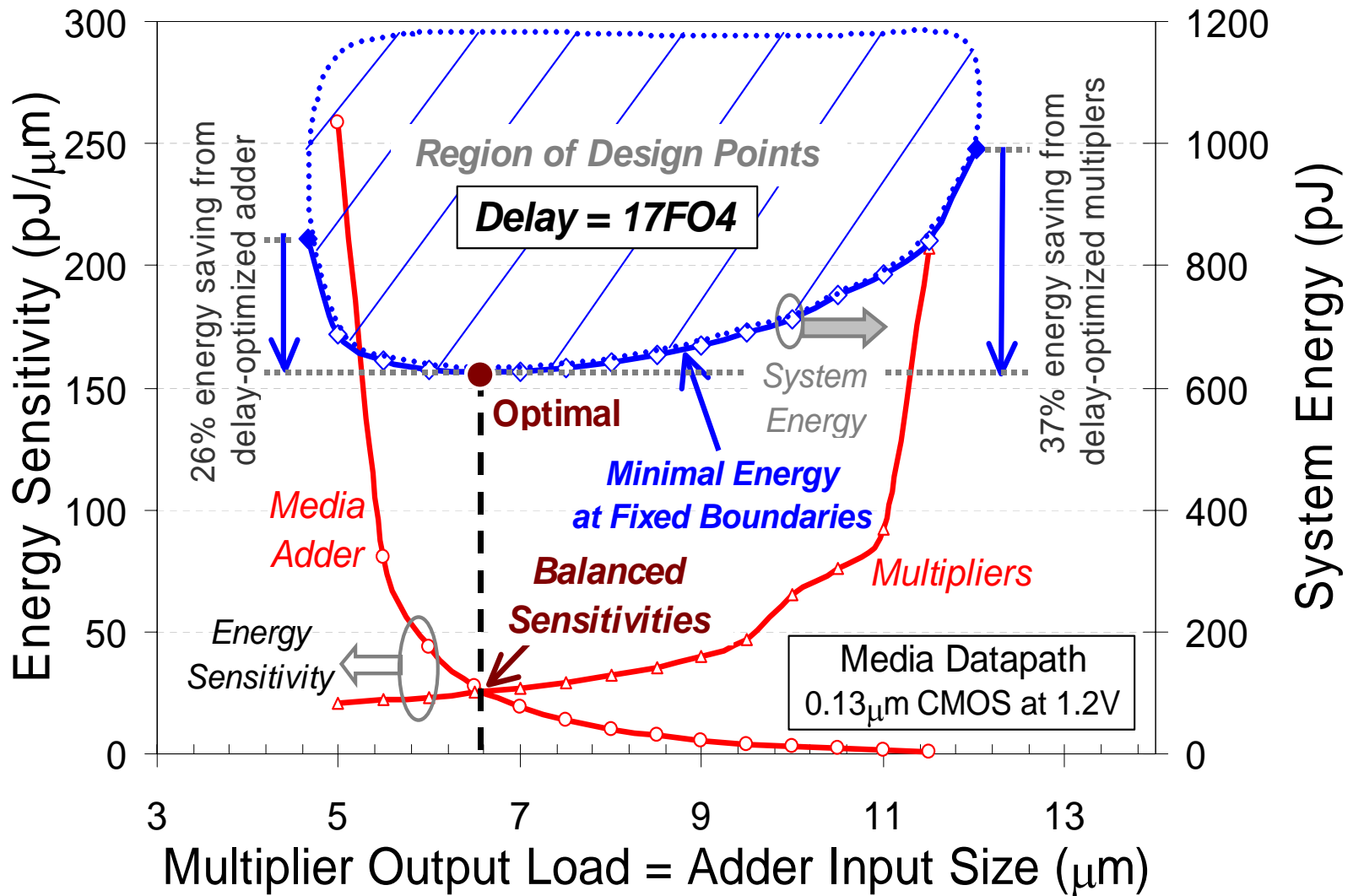X2, U2    C1, S1    C2, S2    Y2, V2

64-b ALU
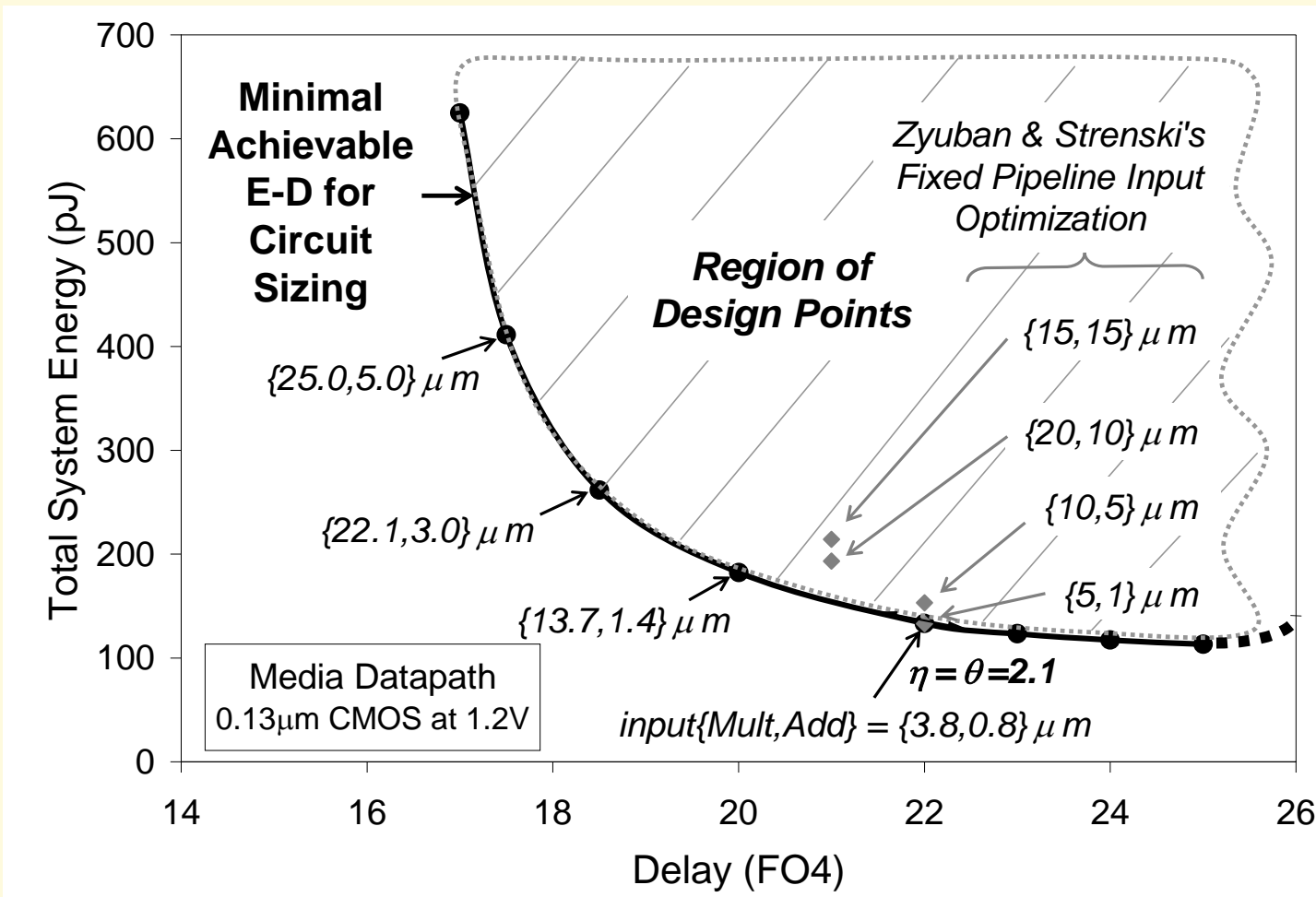
PIPELINE 3

64-b Z

# Case Study: Optimal Algorithm



Simplified Datapath

# Case Study: Media Datapath Solution

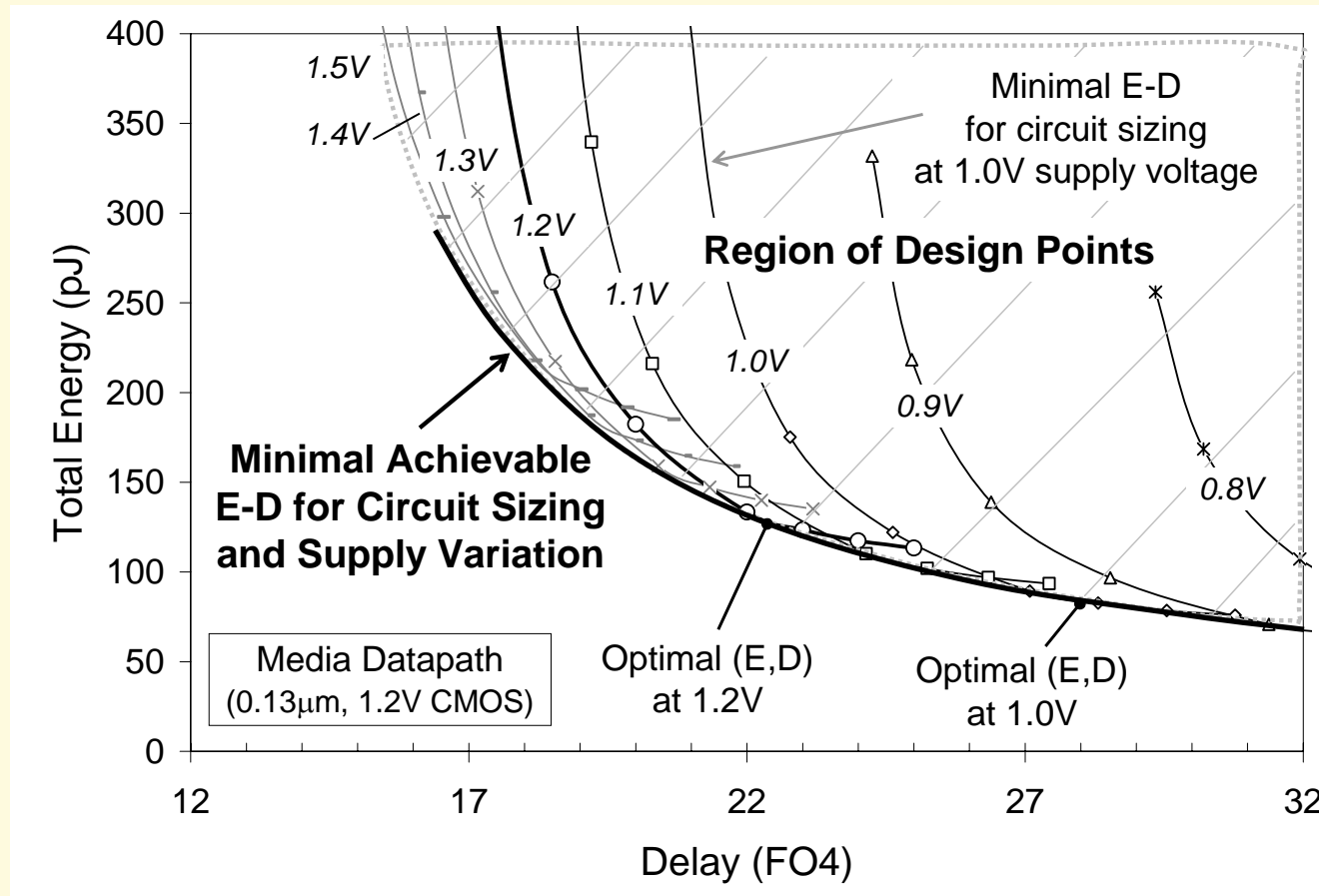# System Energy-Delay @ $V_{DD}$=const.



- Similar E-D characteristics as single stages
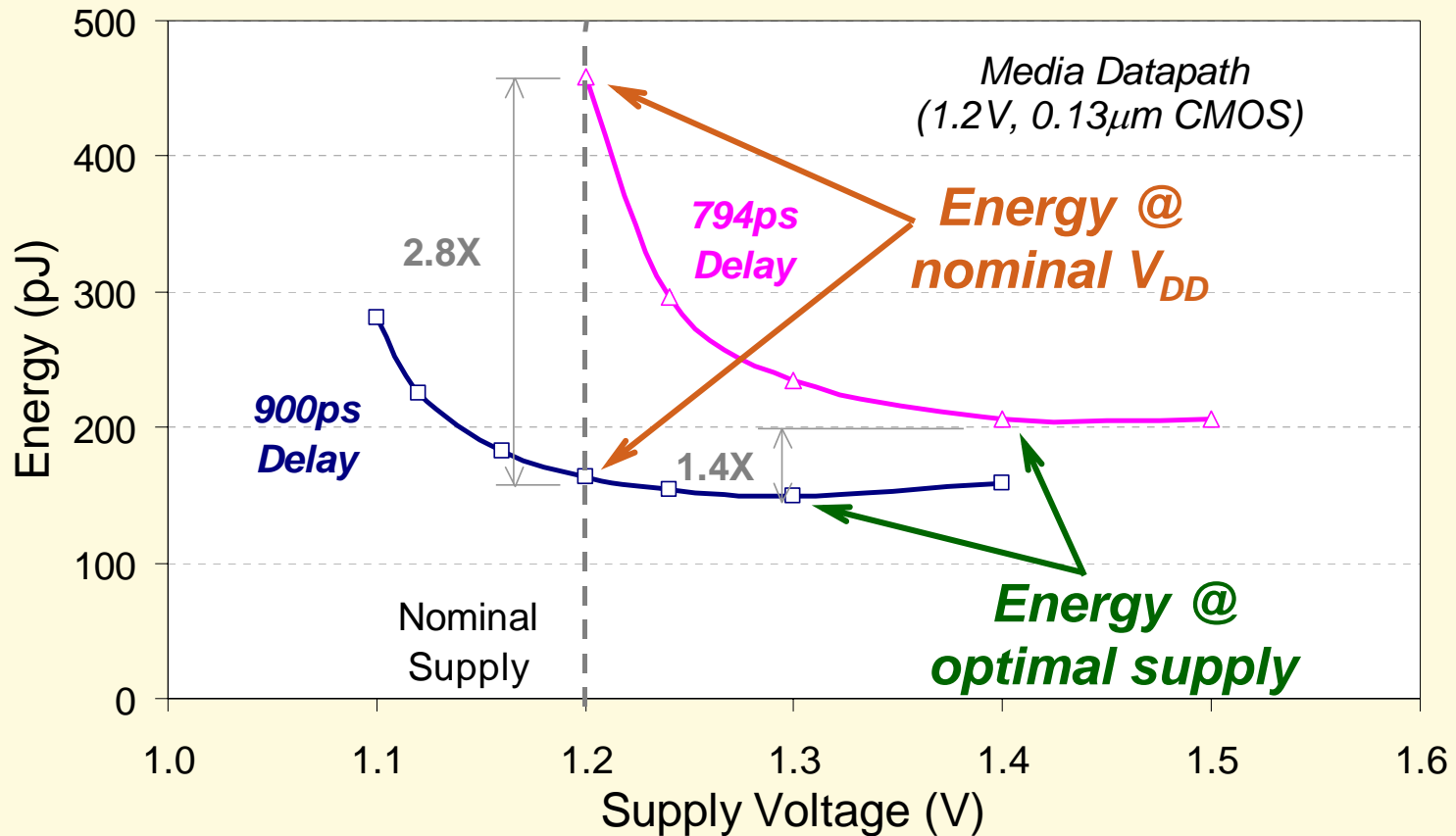- Possibly less input size @ lower delay

# Effect of Supply Scaling



- Efficient {E, D, $V_{DD}$} are dependent
- Optimal criterion:

$$\frac{D}{E}\frac{\partial E}{\partial D}\bigg|_{sizing} = \frac{D}{E}\frac{\partial E}{\partial D}\bigg|_{supply} \quad or \quad \eta_{system} = \theta$$
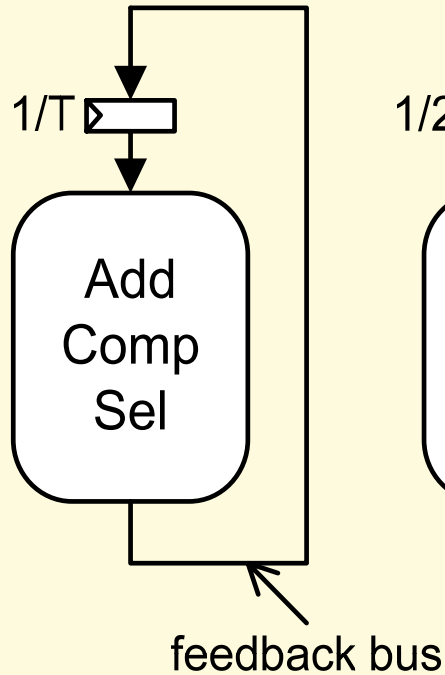
# Potential Saving of System Energy



- Significant energy saving with correct supply or delay selection
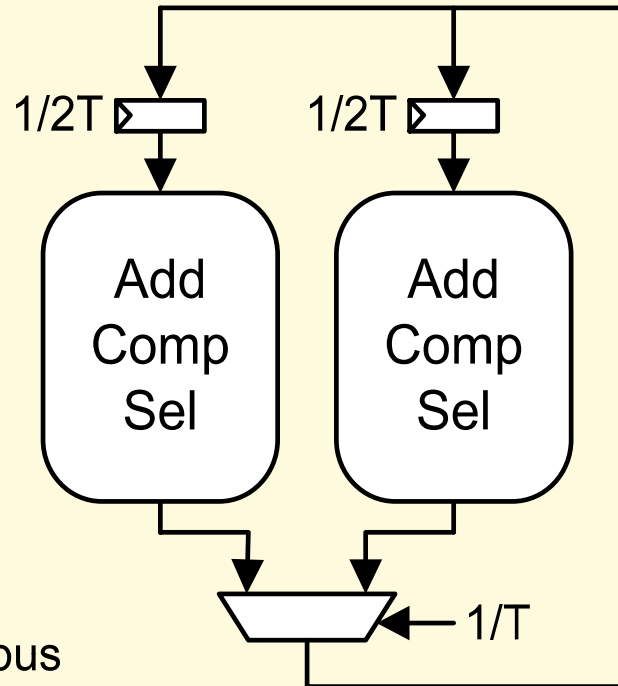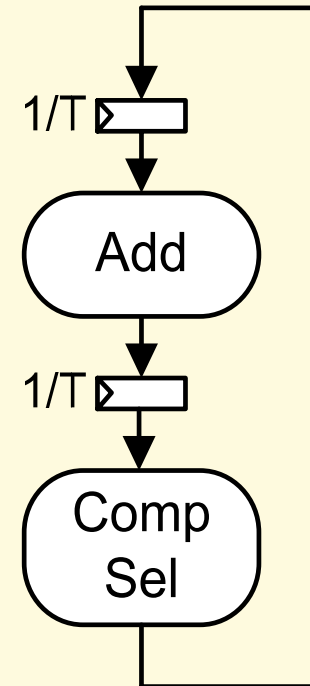
# Energy-Delay Improvement of Pipelined Stages

# Architectural Advantages



(a) Reference

(b) Parallelism

(c) Pipelining

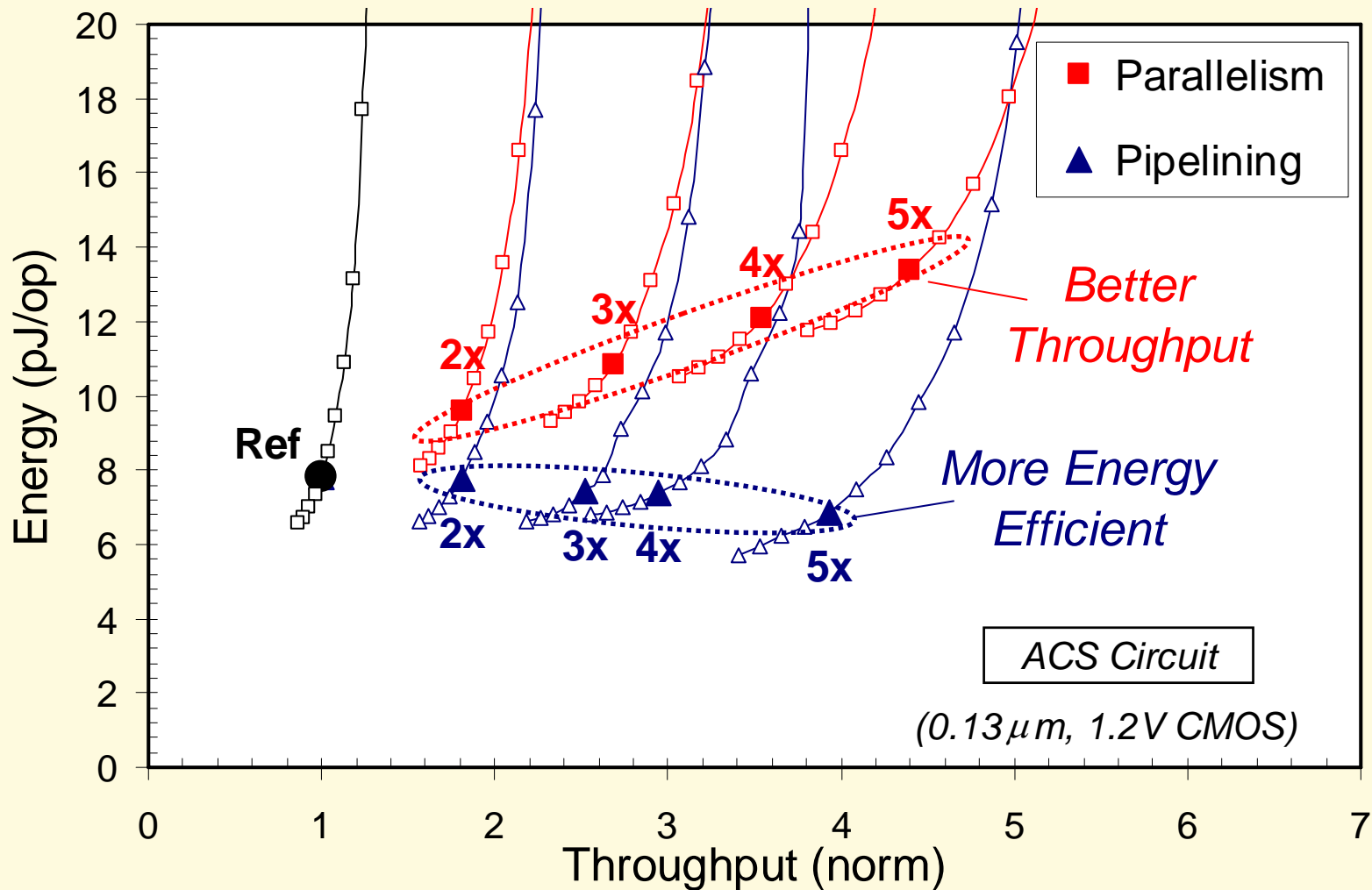- Feedback bus is typical
- Length = 128μm

- $\cong$ N • area
- $\cong$ N • bus length
- 1/N clock rate

- $\cong$ area
- $\cong$ bus length
- Same clock rate

# Energy-Throughput Comparison



- Pipelining is mostly more efficient in E-D domain!

# ACS Unit Implementation

*Application:*
*64-state rate-$\frac{1}{2}$*
*Viterbi decoder*

| PM Mem. | PM Mem. |
|---------|---------|

| ACS0 | ACS4 |
| ACS1 | ACS5 |
| ACS2 | ACS6 |
| ACS3 | ACS7 |

| PM Mem. | PM Mem. |
|---------|---------|

**(a) Non-Pipelined**

- Not pipelined
- 8 ACS circuits
- Largest area
- *Least PM entries per ACS*

- Pipelined-2
- 4 ACS circuits
- Midway area
- *More PM entries / ACS*

| PM Mem. | PM Mem. |
|---------|---------|

| ACS0 | ACS2 |
| ACS1 | ACS3 |

| PM Mem. | PM Mem. |
|---------|---------|

**(b) Pipelined-2**

- Pipelined-4
- 2 ACS circuits
- Small area
- *Most PM entries / ACS*

| PM Mem. | PM Mem. |
|---------|---------|

| ACS0 |
| ACS1 |

| PM Mem. | PM Mem. |
|---------|---------|

**(c) Pipelined-4**

# Energy-Throughput Comparison



- Less energy for deeper pipeline at given throughput

# Designing a System
# for a Fixed Performance
# in the Energy-Delay Space

                        Energy-Efficient CMOS Circuit Design

# Pipelining and Parallelism for an ACS Circuit



(a) Single ACS

(c) Pipelined ACS

(b) Parallel ACS

Ideal Degree-N Parallelism
- N-time throughput
- (Area, routing) overhead

Ideal Depth-N Pipelining
- N-time throughput
- CSE overhead

# Energy-Delay Results for Parallelism



H. Q. Dao, B. R. Zeydel, V. G. Oklobdzija, "Energy-Efficient Optimization of the Viterbi ACS Unit Architecture", Proceedings of the Asian Solid-State Circuit Conference, A-SSCC 2005, Hsinchu, Taiwan, November 1-3, 2005.

# Energy-Delay Results for Pipelining



H. Q. Dao, B. R. Zeydel, V. G. Oklobdzija, "Energy-Efficient Optimization of the Viterbi ACS Unit Architecture", Proceedings of the Asian Solid-State Circuit Conference, A-SSCC 2005, Hsinchu, Taiwan, November 1-3, 2005.

# Energy-Delay Estimation Matches Complex Circuit Simulation

# Energy Efficiency of Architectural Choices
## (including supply scaling and circuit sizing)



*Slightly Better Delay*

5x 4x 3x 2x 1X

*More Energy Efficient*

H. Q. Dao, B. R. Zeydel, V. G. Oklobdzija, "Energy-Efficient Optimization of the Viterbi ACS Unit Architecture", Proceedings of the Asian Solid-State Circuit Conference, A-SSCC 2005, Hsinchu, Taiwan, November 1-3, 2005.

> 52%

▲ Parallelism
● Pipelining

—— *Circuit Sizing*
— · — *Supply Scaling*

ACS Circuit

**Energy [pJ/op]** (y-axis: 0 to 20)
**Delay [FO4]** (x-axis: 0 to 25)

# Eergy-Delay Results for Parallelism and Pipelining



ACS Energy [pJ/op] vs Delay [FO4]

Supply = 1.2V

Parallel-8

depth•degree = 8

Parallel-4
Pipeline-2

depth•degree = 4

Parallel-4

Parallel-2
Pipeline-2

Parallel-2

Reference
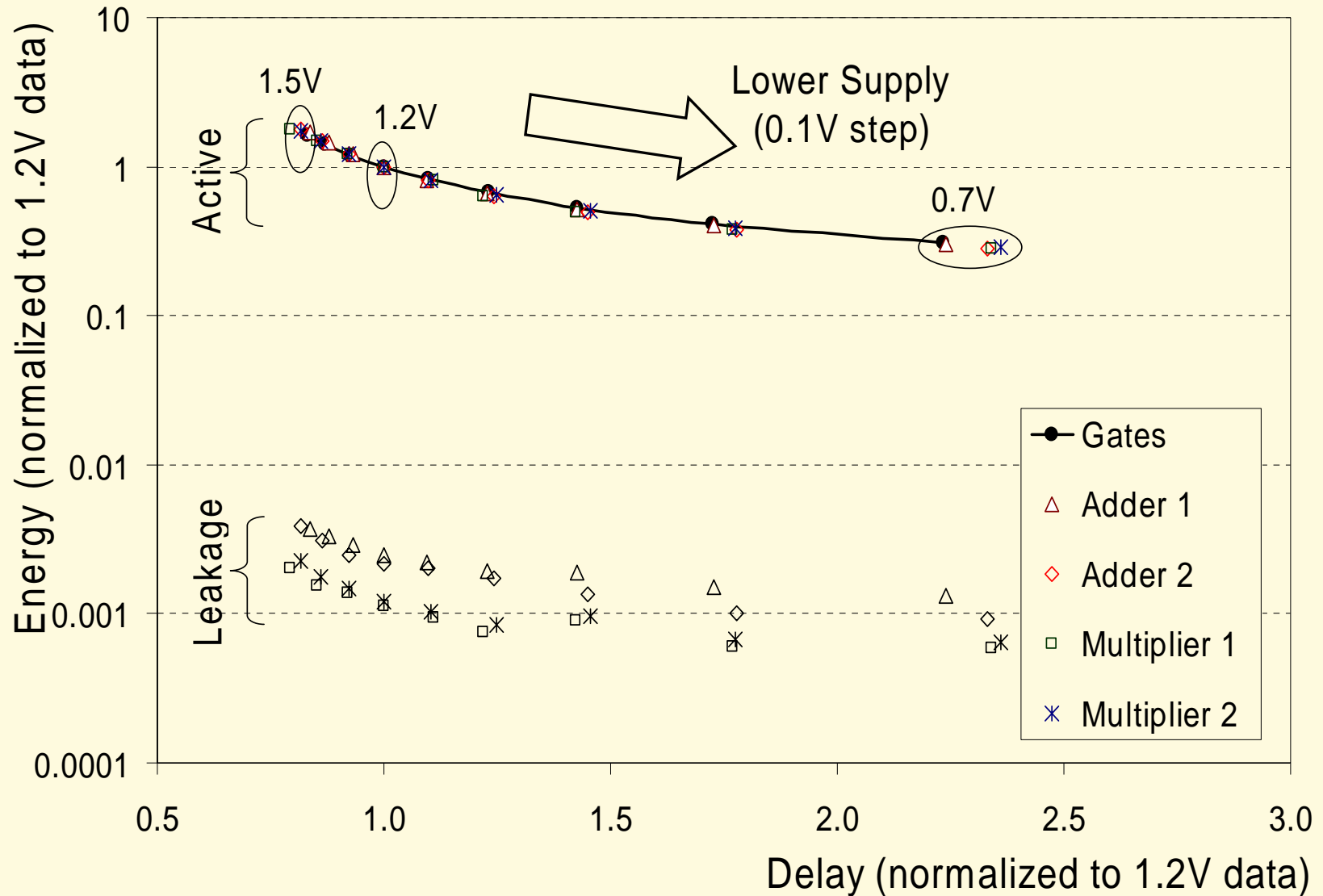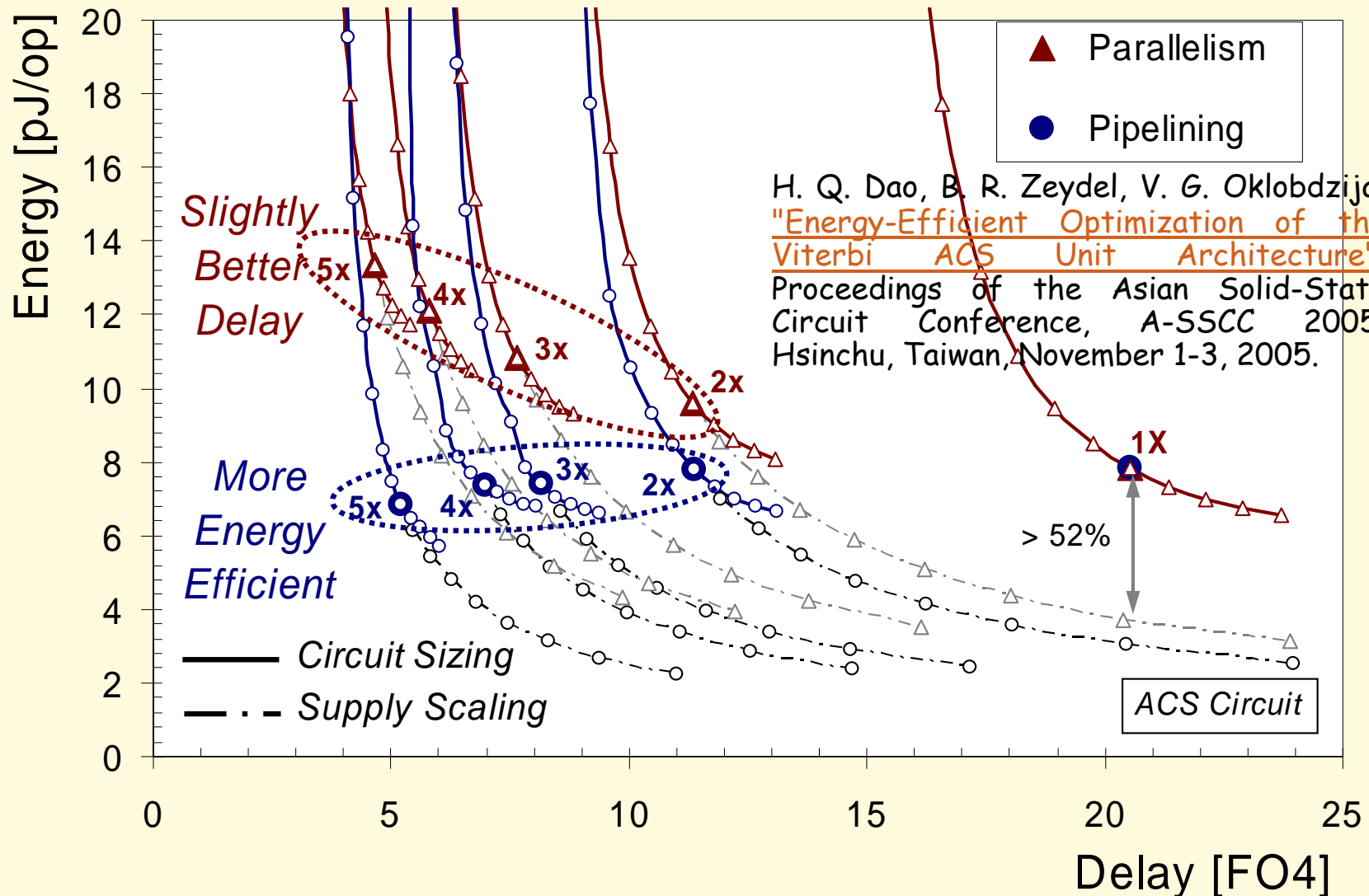
Parallel-2
Pipeline-4

Pipeline-4

Pipeline-2

H. Q. Dao, B. R. Zeydel, V. G. Oklobdzija, "Energy-Efficient Optimization of the Viterbi ACS Unit Architecture", Proceedings of the Asian Solid-State Circuit Conference, A-SSCC 2005, Hsinchu, Taiwan, November 1-3, 2005.
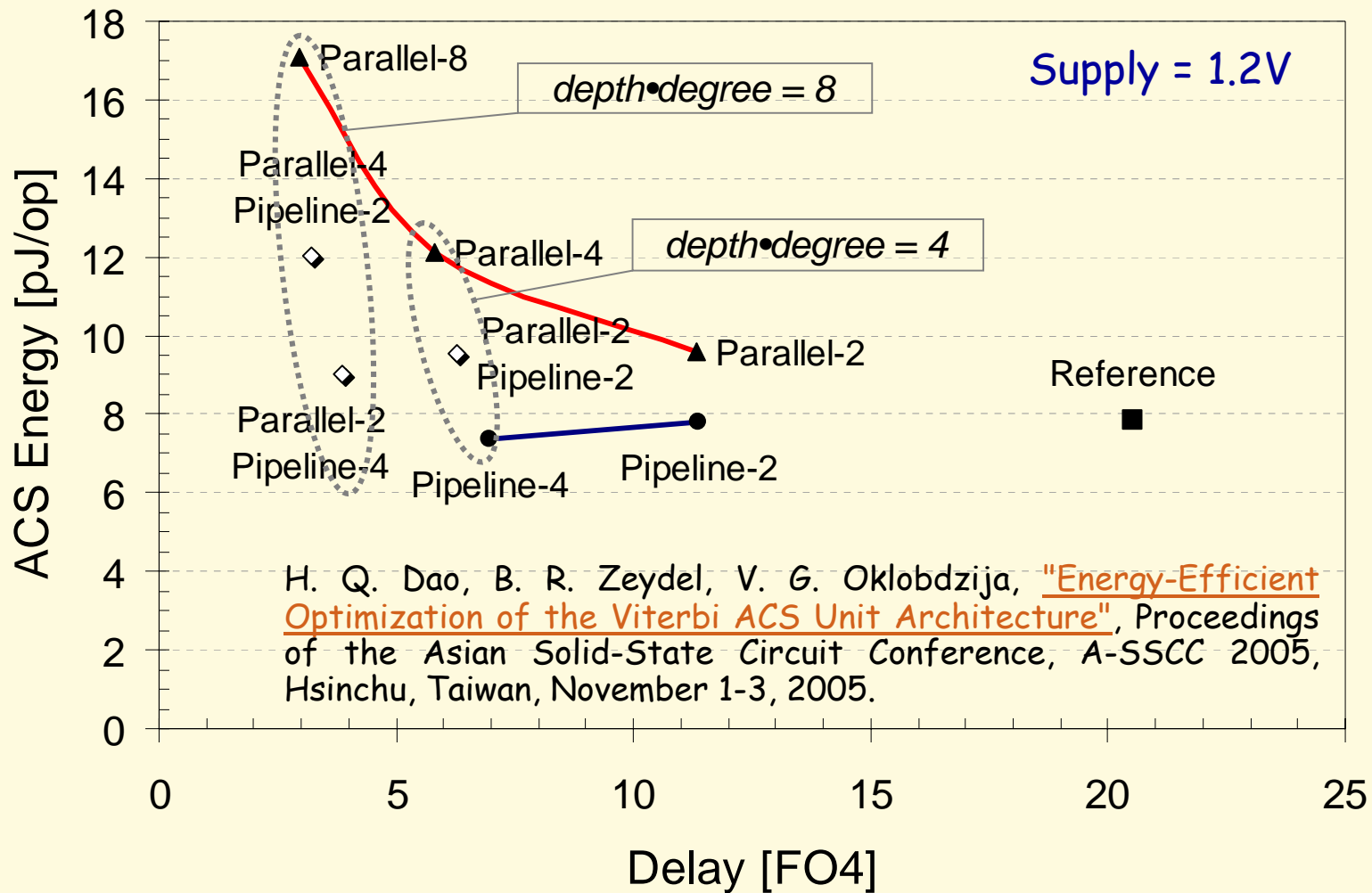
## Pipelining has the Best Energy-Delay Tradeoff

# Is it possible to lower the Energy ?



- **Reduce Energy for same Delay!**
- **Improve Delay for same Energy!**

# Achieved Energy Savings in KS and HC Adders



- Simulation of 64-bit static adders confirms saving!

# Achieved Energy Savings in Representative Adders



- Comparison of high-performance 64-bit adders
  - Delay and Energy Optimized

# Reduction of Hot-Spots



Delay Optimized

Energy Optimized

Reduces Hotspots

- Energy minimization improves hotspots!

# Accomplishments
## (June 30, 2003 to June 30, 2004)



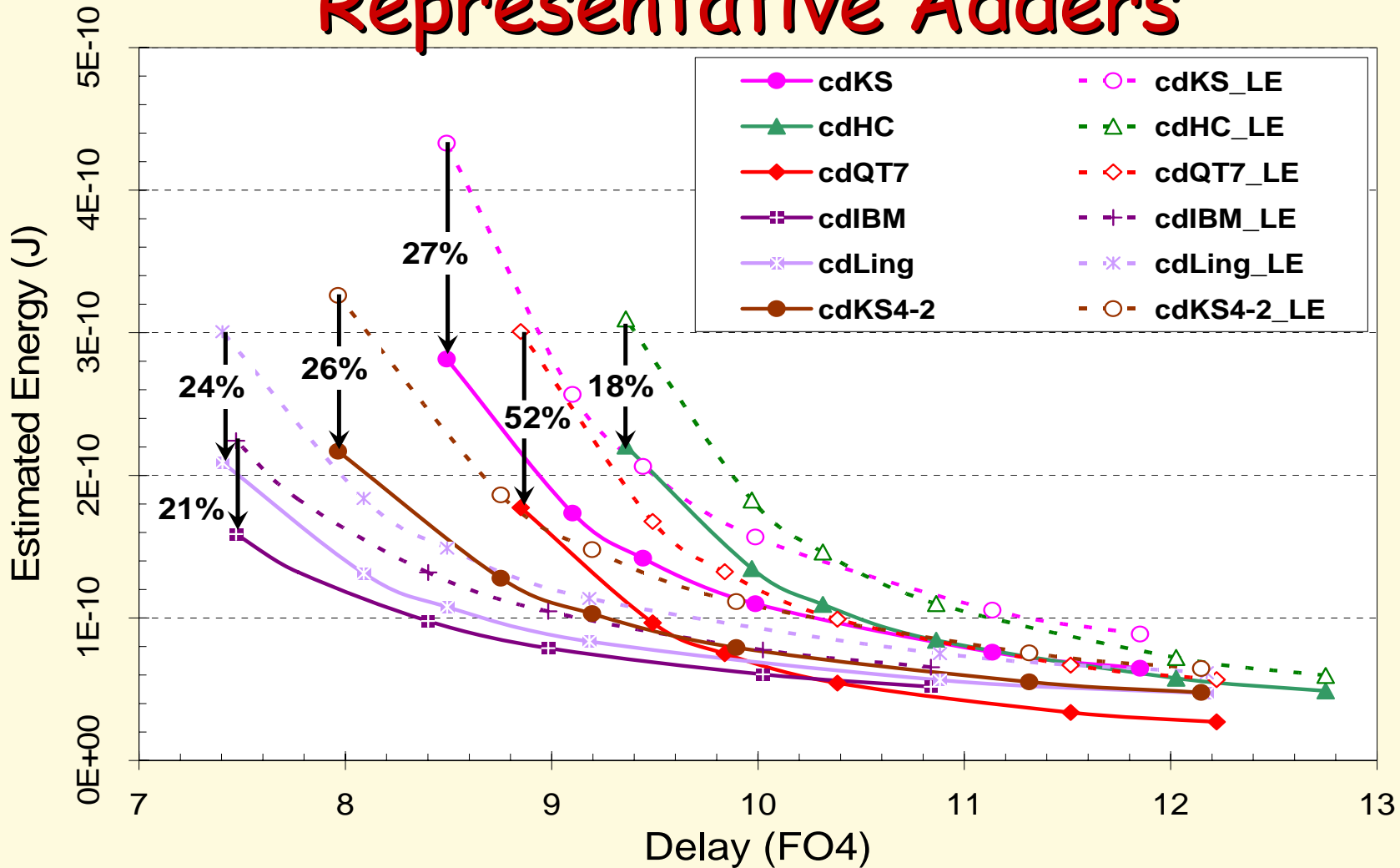90nm technology

Energy [pJ] vs Delay [ns] graph with "Optimized Design", "Initial Design", "Delay Saving", and "Energy Saving" curves and annotations.

Worst Case Energy Vector With 100% Input Activity

Collaboration with Intel AMR

Energy-Efficient CMOS Circuit Design

# Summary

- Energy-Efficient Design requires:
  - Early structure comparison in energy-delay space
  - Early Layout/Floorplanning
  - Optimization using energy minimization objective function

- LE does not guarantee a good design

- Our method of **energy-minimization** focuses on reducing power
- The same principles hold for other logic functions

# Future Work

- Methodology improvement
  - General design rules
  - Algorithms with fast convergence
  - Guidelines for close-to-optimal solutions

- CAD tool
  - Custom vs. standard cell library

- Improve gate modeling & characterization
  - Worst-case vs. single-switching,... or between?
  - Process variation

April 19, 2010

# Publications on Energy-Delay:

- Vojin G. Oklobdzija, Bart R. Zeydel, Hoang Dao, Sanu Mathew, Ram Krishnamurthy, "*Energy-Delay Estimation Technique for High-Performance Microprocessor VLSI Adders*", Proceedings of the International Symposium on Computer Arithmetic, ARITH-16, Santiago de Compostela, SPAIN, June 15-18, 2003.

- Hoang Q. Dao, Bart R. Zeydel, Vojin G. Oklobdzija, "*Energy Minimization Method for Optimal Energy-Delay Extraction*", Proceedings of the European Solid-State Circuits Conference, ESSCIRC 2003, Estoril, PORTUGAL, September 16-18, 2003.

- V. G. Oklobdzija, B. R. Zeydel, H. Q. Dao, S. Mathew, R. Krishnamurthy, "Comparison of High-Performance VLSI Adders in Energy-Delay Space", *IEEE Transaction on VLSI Systems*, Volume 13, Issue 6, pp. 754-758, June 2005.

- Hoang Q. Dao, Bart R. Zeydel, Victor Zyuban, and Vojin G. Oklobdzija, "On Energy Optimization of Digital Systems", The fourth annual IBM Austin Conference on Energy-Efficient Design, ACEED 2005, Austin, Texas, March 1-3, 2005.

- H. Q. Dao, B. R. Zeydel, V. G. Oklobdzija, "Energy-Efficient Optimization of the Viterbi ACS Unit Architecture", *Proceedings of the Asian Solid-State Circuit Conference, A-SSCC 2005, Hsinchu, Taiwan, November 1-3, 2005.*

- H. Q. Dao, B. R. Zeydel, V. G. Oklobdzija, "Energy Optimization of Pipelined Digital Systems Using Circuit Sizing and Supply Scaling", *IEEE Transaction on VLSI Systems,*Vol. 14, Issue 2, Feb. 2006 pp. 122-134.

- S. K. Hsu, S. K. Mathew, M. A. Anders, B. R. Zeydel, V. G. Oklobdzija, R. K. Krishnamurthy, S. Y. Borkar, "*A 110 GOPS/W 16-bit Multiplier and Reconfigurable PLA Loop in 90-nm CMOS*", IEEE Journal of Solid-State Circuits, Vol.41, No.1, January 2006.

- B. R. Zeydel, D. Baran, V. G. Oklobdzija, "*Energy Efficient Design of High-Performance VLSI Adders*", IEEE Journal of Solid-State Circuits, June 2010.