# ISCAS

## Vancouver • 2004

## Tutorial Notes

# Tutorial 6:
# Clocking and Synchronization issues in sub-100nm System on Chip (SoC) Designs

**Presented by:**
**Ramalingam Sridhar, State University of New York at Buffalo**
**Ram Krishnamurthy, Intel Corporation**
**Sanu K. Mathew, Intel Corporation**

**Sunday Afternoon, May 23, 13:15 - 16:15**

# ISCAS
## Vancouver • 2004

## Tutorial Notes

## Tutorial 6:
## Clocking and Synchronization issues in sub-100nm System on Chip (SoC) Designs

**Presented by:**
**Ramalingam Sridhar, State University of New York at Buffalo**
**Ram Krishnamurthy, Intel Corporation**
**Sanu K. Mathew, Intel Corporation**

**Sunday Afternoon, May 23, 13:15 - 16:15**
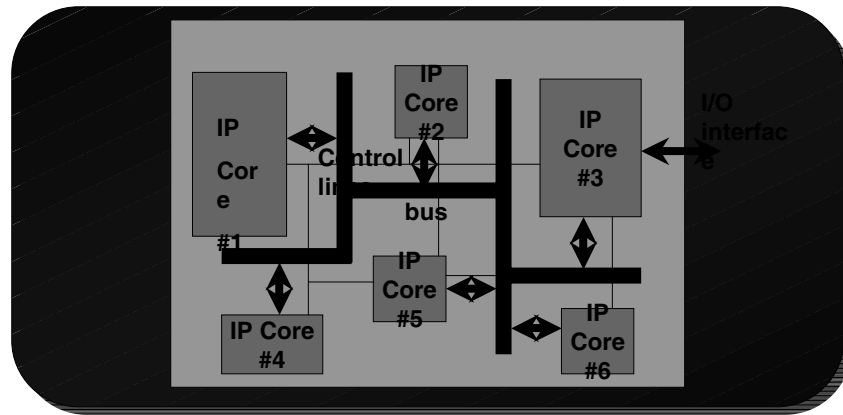
# Clocking and Synchronization in UDSM SOCs

**Professor Ramalingam Sridhar**
**University at Buffalo, SUNY**
**Buffalo, NY, USA**

---

# Outline

- **UDSM SoC – An Overview**
- **Design Challenges for SoCs**
- **Clocking and Synchronization Problems**
- **Potential Solutions**
  - Asynchronous Design
  - GALS
  - Current-mode Interconnects
  - Other Approaches
- **Delay Variations and Wave-pipelining**
- **Summary**

2

# What is SoC?



- Heterogeneous integration of components makes SoCs powerful, flexible and versatile

# UDSM SoCs Characteristics

- **Faster gates allow for systems operating at very high frequencies (multi-GHz)**
- **Smaller transistors result in**
  - **Higher Performance**
  - **Higher Density**
  - **More functionality**
- **Supply voltage and voltage swings reduced to keep power under control**
- **Threshold voltages scaled for high drive current and performance maintenance**

# SoC Requirements

- **To ensure seamless integration and operation of SoCs, different IPs have to be properly clocked and synchronized with one another**

- **Design of efficient clock distribution network and synchronization circuitry dictates proper functioning of the SoC**
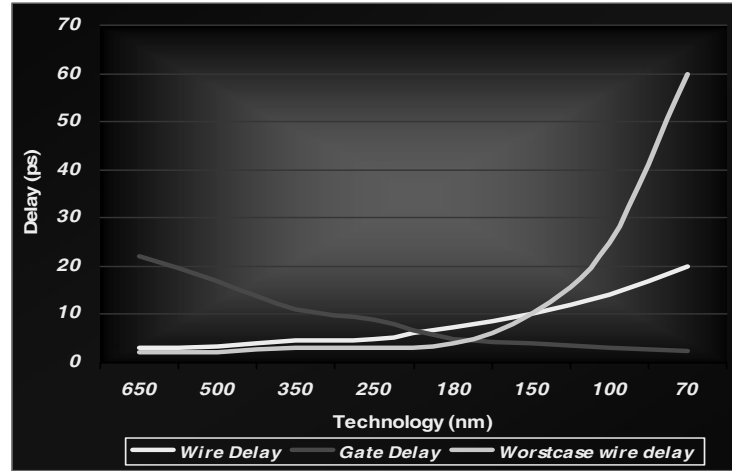
5

# Challenges of SoC design in the UDSM domain

- **Higher power consumption due to increasing sub-threshold and gate-oxide leakage currents**
- **Increasing interconnect delays limiting system throughput**
- **Synchronization problems due to unpredictable delay variations**
- **Environmental and process variations-induced timing uncertainties**
- **Signal Integrity problems due to closer integration and higher frequencies**
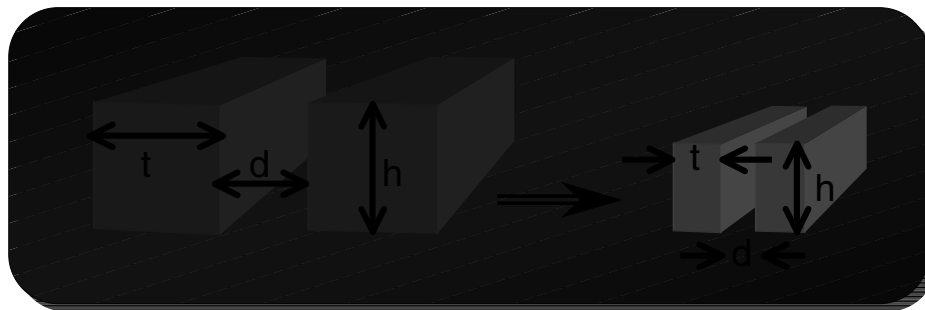- **Susceptibility to soft-error (single-event upsets)**

6

# Scaling of gate and wire delays



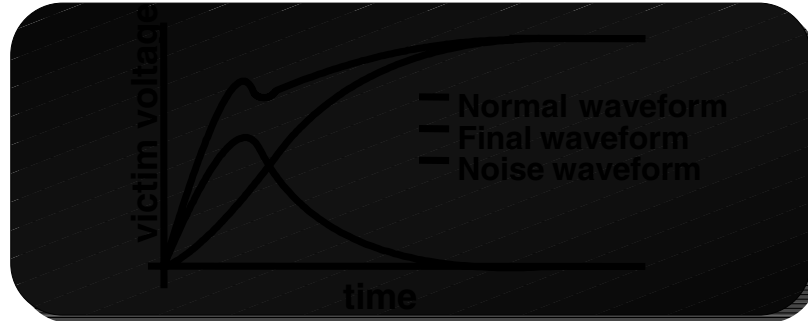- **Interconnect delays dominate system performance in sub-90nm designs**

7

# Interconnect Scaling



- **Wire resistance $R = \rho l/A$**
- **With scaling, $t$ and $h \downarrow \Rightarrow A \downarrow$ and $R \uparrow \Rightarrow$ delay$\uparrow$**
- **Hence aspect ratio ($h$)$\uparrow$ to reduce delay**
- **However, high aspect ratios (~3) and reduced $d$ lead to perfect formation of a parallel plate capacitor, inducing crosstalk between wires**

8

4

# Delay Variations in Interconnects



- **Crosstalk induces noise voltage in neighboring wires that can aid/degrade signal integrity**
- **Thus wire delay can vary widely and unpredictably in UDSM technologies**
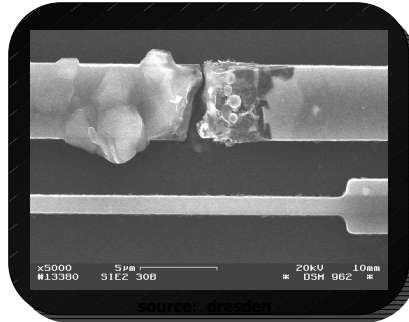- **Techniques that reduce delay and its variations necessary**

9

# Signal Integrity - Noise

- **In spite of Cu interconnects and low-K devices, higher layer interconnects are thin**
- **No low-K insulators between wires on same layers and cross-cap percentage is higher**
- **Hence, crosstalk induced delay variation is a serious challenge in UDSM**
- **Shielding, buffer insertions & net-reorder can still play crucial roles**

10

# Temperature Effects



- **Electromigration of the interconnects is accelerated due to increased resistance and hence heat generation in the wires**

- **Using low-k dielectrics for intra-level gap fill can cause significant increase in thermal effects owing to their lower thermal conductivity than $SiO_2$**
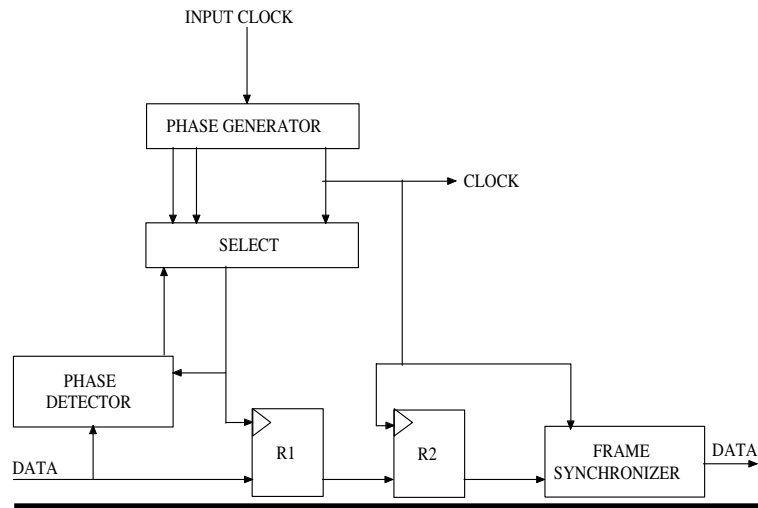
11

# Clocking and Synchronization of IP Cores

- **Popular synchronization techniques of the IPs**
  - **Mesochronous clocking**
  - **Plesiochronous clocking**

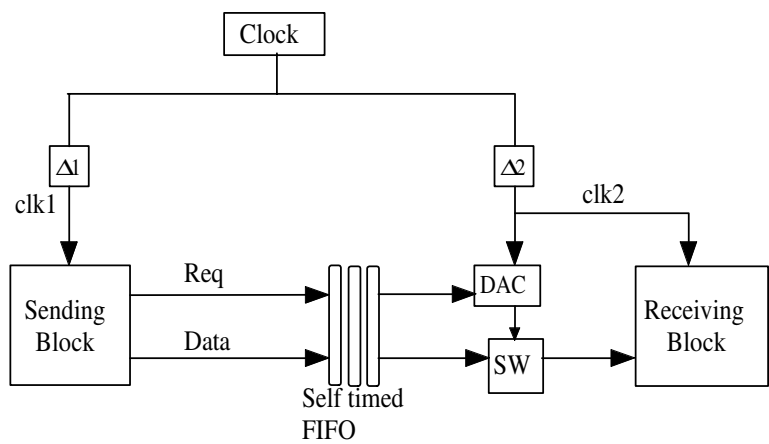- **Other schemes such as Pausable clocking are also existent**

12

# Block Diagram of a Re-synchronizer
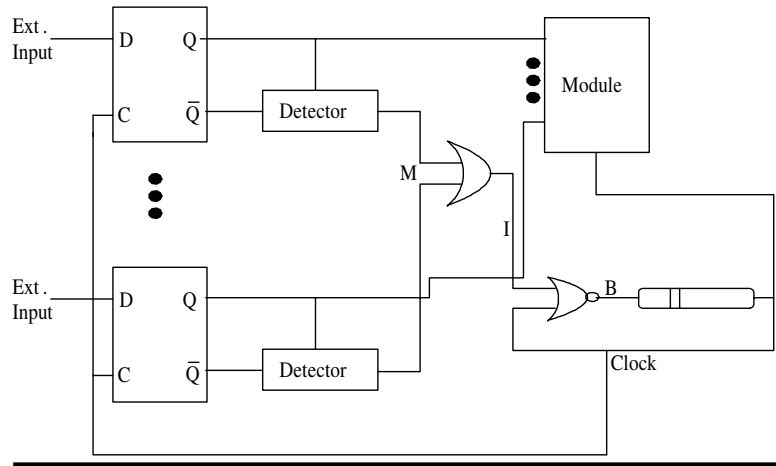
INPUT CLOCK

PHASE GENERATOR

CLOCK

SELECT

PHASE DETECTOR

R1

R2

FRAME SYNCHRONIZER

DATA

DATA

13

# Mesochronous Clocking

Clock

Δ1

clk1

Δ2

clk2

Sending Block

Req

Data

Self timed FIFO

DAC
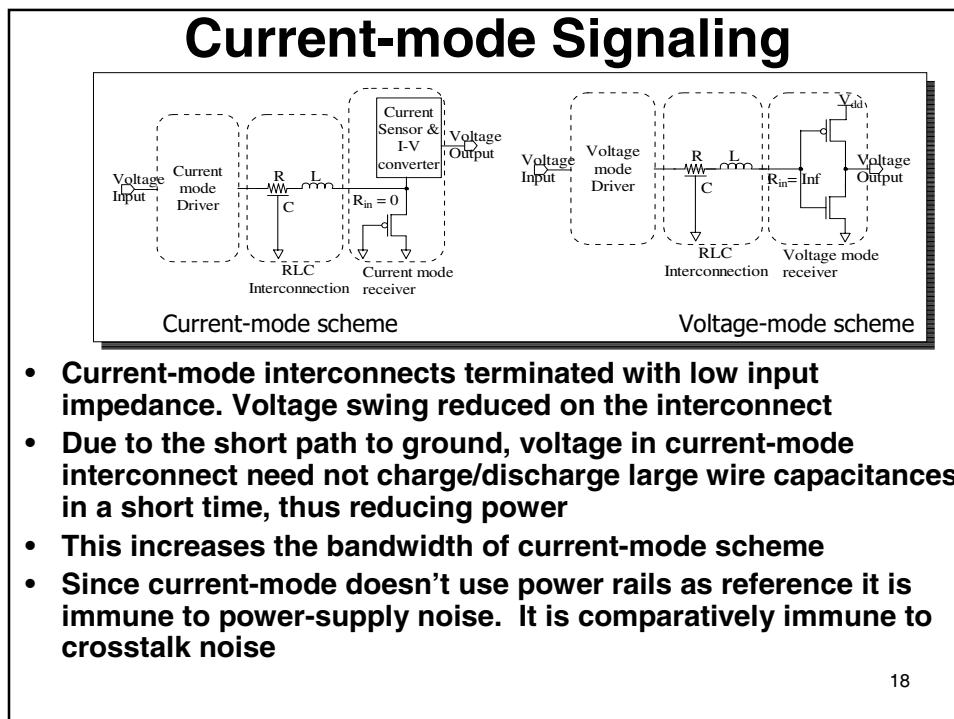
SW

Receiving Block

14

7

# Pausable Clocking



15

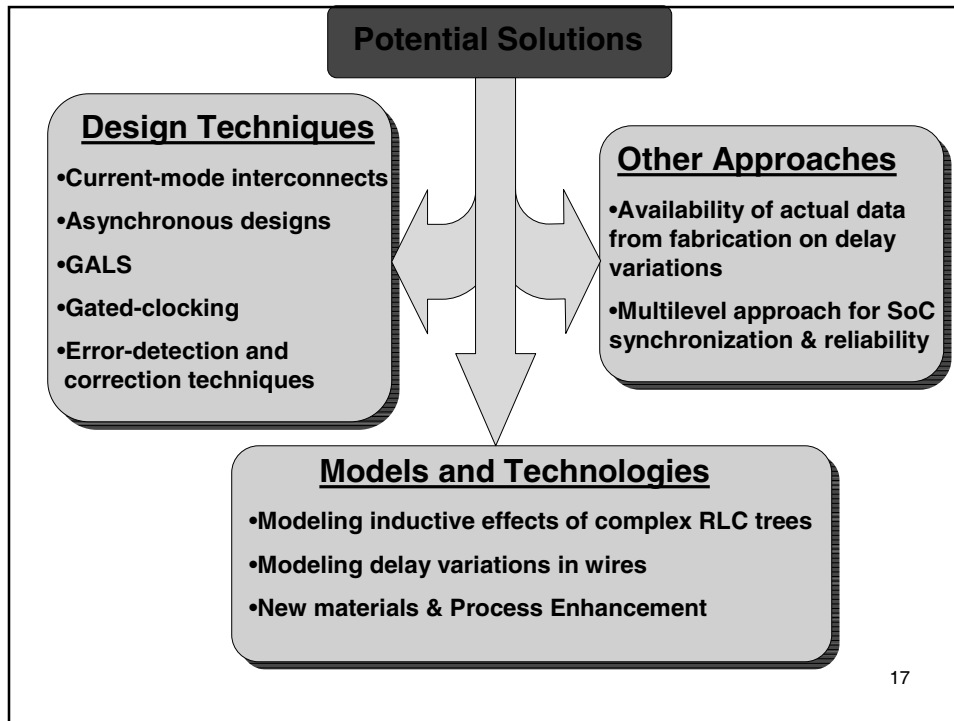# Clocking of different IPs

- **Clock signals should reach the respective components at appropriate times for successful operation of SoC**

- **Different Clock distribution schemes, such as, symmetric H tree, X tree and buffered clock tree schemes used to minimize skew**

16

## Potential Solutions

### Design Techniques

•Current-mode interconnects

•Asynchronous designs

•GALS

•Gated-clocking

•Error-detection and correction techniques

### Other Approaches

•Availability of actual data from fabrication on delay variations

•Multilevel approach for SoC synchronization & reliability

### Models and Technologies

•Modeling inductive effects of complex RLC trees

•Modeling delay variations in wires

•New materials & Process Enhancement

17

---

# Current-mode Signaling



Current-mode scheme          Voltage-mode scheme

- **Current-mode interconnects terminated with low input impedance. Voltage swing reduced on the interconnect**
- **Due to the short path to ground, voltage in current-mode interconnect need not charge/discharge large wire capacitances in a short time, thus reducing power**
- **This increases the bandwidth of current-mode scheme**
- **Since current-mode doesn't use power rails as reference it is immune to power-supply noise. It is comparatively immune to crosstalk noise**
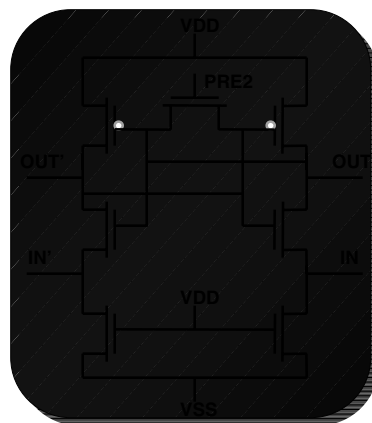
18

# Types of Current-mode Signaling

- **Differential and single-ended signaling**
  - **Single-ended signaling uses single wire per signal thus reducing the area and associated power overhead in differential signaling**
  - **Differential signaling transmits a differential current per signal (using 2 wires)**
- **Noise immunity of differential signaling significantly high compared to single-ended**
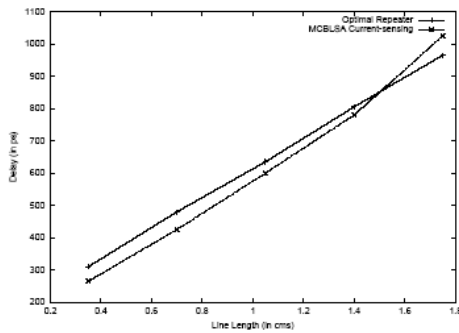
19

# Differential method implementation

- **MCBLSA interconnect delay is less compared to optimal repeater insertion (avg. 6%)**

- **Power consumed less than that in optimal repeater insertion**

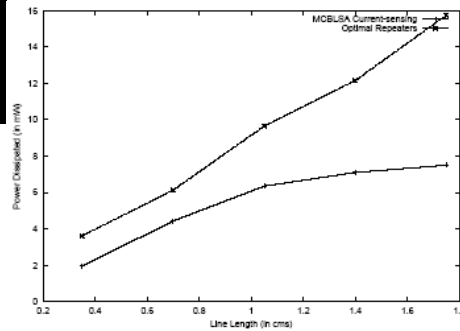- **Power consumption remains almost constant with increase in wire length**



**Modified Clamped Bit Line Sense Amplifier (MCBLSA) scheme, A. Maheswari *et.al.*, 2001**

20

# Differential method implementation



Delay of optimal repeater method vs. MCBLSA scheme

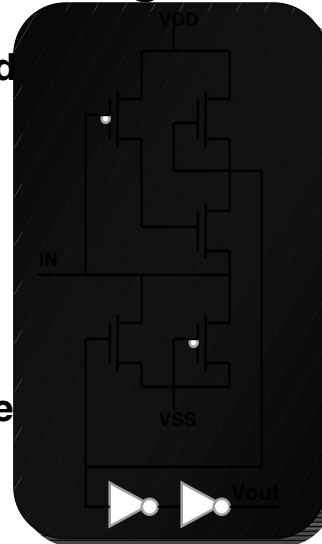**Power consumed in optimal repeater method vs. MCBLSA scheme**



---

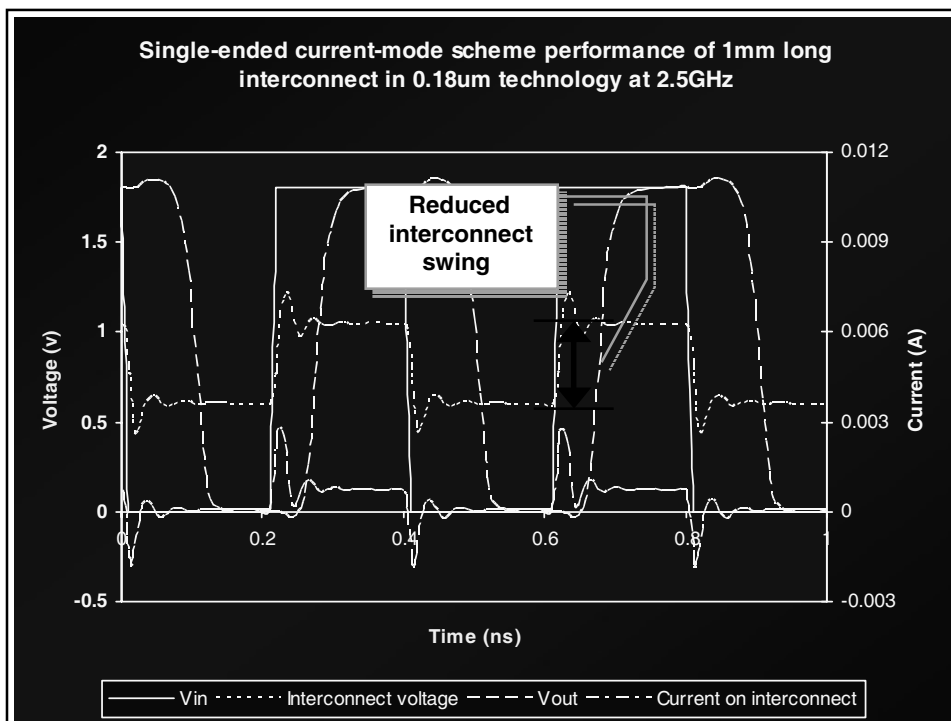# Differential method - disadvantages

- **Area overhead: Consumes almost double the routing area due to additional wire per signal**

- **Scope for reduction in power consumed if the additional wire is eliminated**

- **Single-ended design overcomes these disadvantages at the expense of noise immunity**

22

# Single-ended signaling

- **Less area and routing overhead compared to MCBLSA scheme**
- **Delay and power consumption improved by 11% over the MCBLSA scheme in 0.18µm technology**
- **Delay variations due to cross-talk increased from 28% in MCBLSA to 32% in this scheme**
- **This scheme could operate successfully up to 10GHz in 0.18µm technology**



**Single-ended current sensing circuit, R. Sridhar *et.al.*, 2004**



Single-ended current-mode scheme performance of 1mm long interconnect in 0.18um technology at 2.5GHz
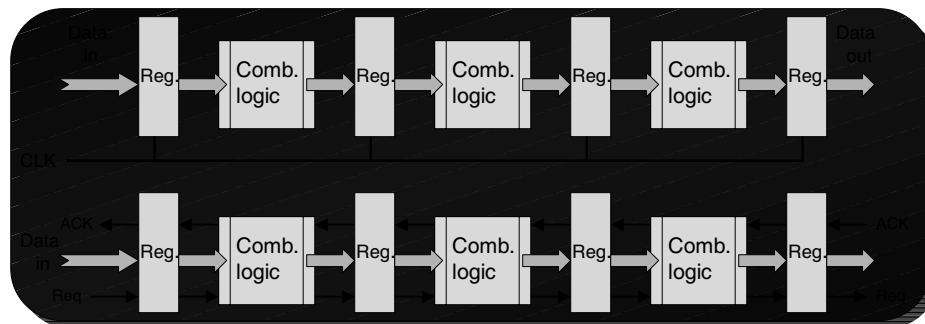
# Current-mode signaling - Summary

- **Current-mode mostly used in inter-chip signaling. On-chip use limited to CMOS SRAM circuits**

- **However its favorable characteristics in terms of bandwidth, noise immunity and power consumption makes it a good choice for on-chip global interconnects**

- **Voltage-mode circuits expected to saturate in the 5-10GHz domain. Current-mode circuits provide a promising alternative in this regard**

25

# Asynchronous Circuits



- **Explicit local synchronization between blocks**
- **Provides modularity for SoC design – plug and play compatibility**
- **Robust with regard to wire delay, temperature and process variations**
- **Helps reduce power supply noise by reducing current peaks around clock edges**

26

# Asynchronous Design – Hurdles

- **EDA tools**
  - **Capable of complex timing analysis**
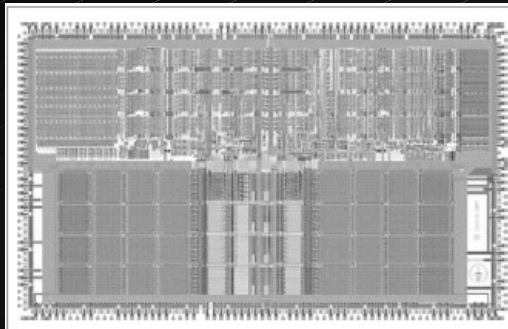
- **Difficult to test**

27

# Asynchronous Implementations

- **Amulet Microprocessors (asynchronous ARM, Univ. of Manchester, 2000)**
- **Caltech Asynchronous Microprocessors (asynchronous MIPS, 1998)**
- **Titac2 (Univ. of Tokyo, 1997)**
- **Intel's RAPPID Instruction Length Decoder (2001)**

28

# The miniMIPS processor
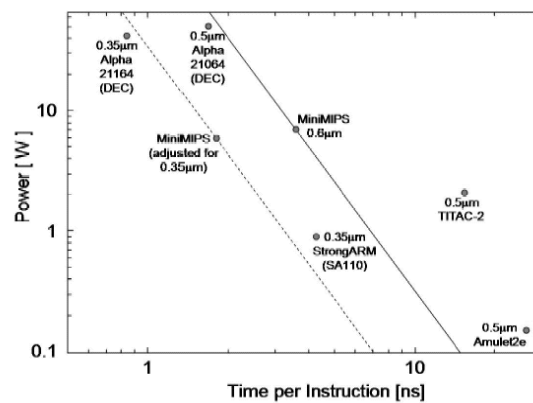


Asynchronous version of MIPS microprocessor

A. J. Martin *et.al.*
Caltech group, 1998

- **Fully asynchronous 32-bit RISC µP similar to MIPS R3000**
- **Implements most of MIPS-I ISA**
- **Has two 4-KB on-chip caches, an instruction cache and a direct-mapped write-through data cache**
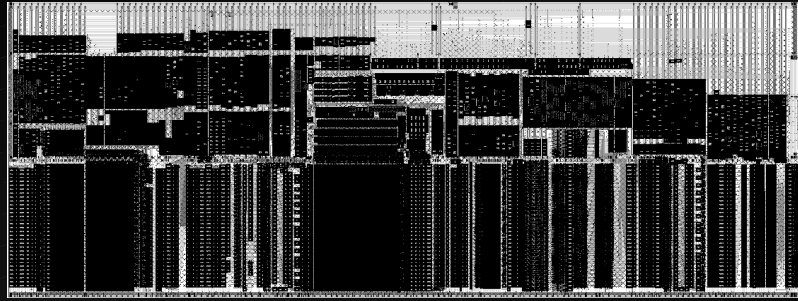
# miniMIPS – Results



- **180MIPS, 4W @ 3.3v; 100MIPS, 850mW @ 2v; 60MIPS, 220mW at 1.5v**
- **Fabricated chip reported to be up to 4 times faster than comparable synchronous commercial µPs**

# Amulet3 - Results



S. B. Furber *et.al.*, University Of Manchester, 2000

- **ARM9TDMI is the synchronous implementation closest to Amulet3**
- **ARM9 – operates up to 120MHz with 1.1MIPS/MHz, 1.8mW/MHz $\Rightarrow$ energy per instruction 610 MIPS/W**
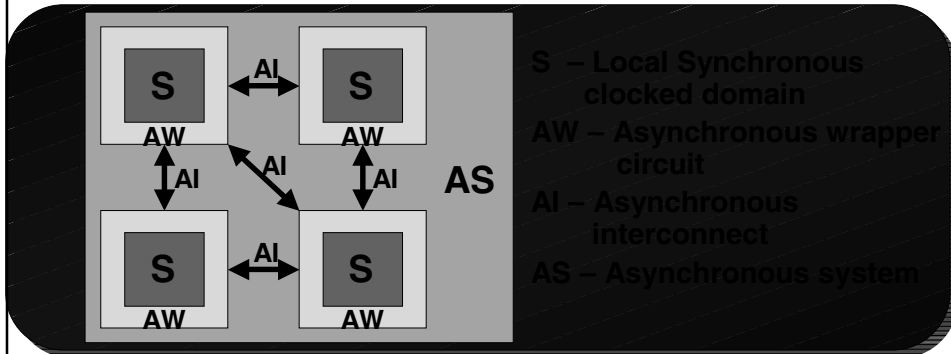- **Amulet3 – 220mW at 85 Dhrystone MIPS, energy per instruction 620 MIPS/W**

31

# Summary of Asynchronous Design

- **Asynchronous design provides a promising alternative to clocking and synchronization problems of synchronous designs**

- **Offers significantly improved performance in terms of throughput and performance while maintaining minimal power overhead**
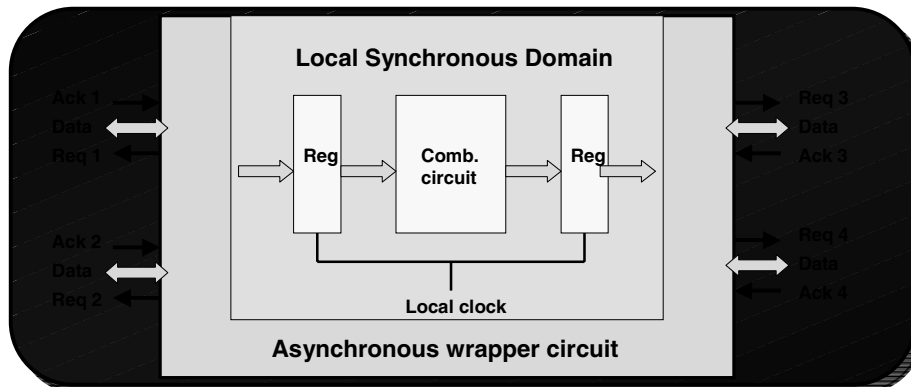
32

# Globally Asynchronous Locally Synchronous (GALS) Systems

S  – Local Synchronous
       clocked domain
AW – Asynchronous wrapper
       circuit
AI – Asynchronous
       interconnect
AS – Asynchronous system

- **Large asynchronous design difficult to realize**
- **GALS combines advantages of synchronous and asynchronous designs**

33

# Globally Asynchronous Locally Synchronous (GALS) Systems

Local Synchronous Domain

Reg | Comb. circuit | Reg

Local clock

Asynchronous wrapper circuit

Ack 1 / Data / Req 1
Ack 2 / Data / Req 2
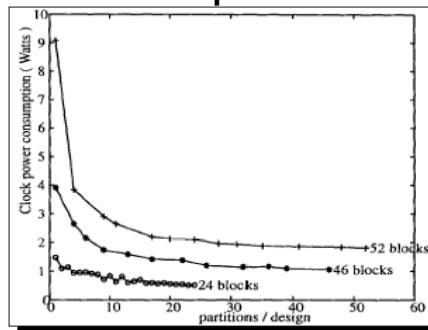
Req 3 / Data / Ack 3
Req 4 / Data / Ack 4

- **GALS immune to global clock skew which degrades performance significantly in synchronous systems**
- **Tolerant to global wire delays and its variations**

34

17

# GALS Implementation - Results

- **The partitioning of GALS is an important factor that affects final system peformance**
- **Some important observations of GALS implementation by J. Oberg *et.al.***
- **Partitioning the system into synchronous (SB) and asynchronous blocks plays a significant role in determining power and performance**
- **Allows the SB's to run at different clock speeds and still be synchronized**
- **Results show upto 70% reduction in power**
- **Overheads in GALS increases with number of partitions**
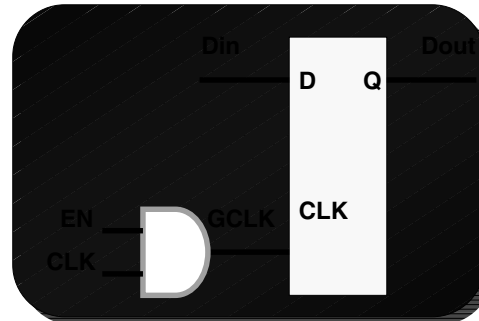


---

# GALS Implementation - Results

- **Some important observations of GALS implementation by D. Marculescu *et.al.***
- **A 5-clock domain GALS processor was implemented in 0.25μm technology**
- **Power consumption in GALS 10% lower compared to synchronous equivalent**
- **Drop in performance ranging between 5-15%**

36

# Gated Clocking

- **The clock input to a synchronous block is turned off when inactive**

- **This saves power dissipated in the block when inactive**



- **Clock combined with other signals using combinational logic to control clock input to a synchronous block**

- **This strategy fits well in particular for SoCs to reduce power consumption**

37

# Gated Clocking …

- **The effectiveness of this scheme depends on the frequency of operation of the various synchronous blocks**

- **The power saved in gated-clocked circuits is 30% lower than buffered clock circuits on an average**

- **The disadvantages of this scheme include**
    - **Area overhead**
    - **Requires complex timing analysis tools**
    - **Introduces additional clock skew due to the combinational elements**

38

# Modeling of wire delay variation

- **Delay uncertainty modeling is a novel idea and needs to be explored further for better system performance**
- **In clock distribution networks, this reduces timing violations and increases system reliability**
- **From the model, wire delay uncertainties can be estimated and circuits can be designed to tolerate the adverse effects of interconnects**

39

# Delay Uncertainty Tolerant Circuits

- **Techniques that design delay uncertainty tolerant circuits are crucial for future SoC designs**

- **Estimation of delay variations play an important role in the above designs**

40

# Delay Uncertainty Tolerant Circuits

- **Specialized circuits and techniques to counteract the UDSM impacts are necessary**
  - **Wave pipelining and other circuit level asynchronous and pipelining techniques provide an insight into this problem**
  - **Incorporating self-checks and some tolerance circuitry**
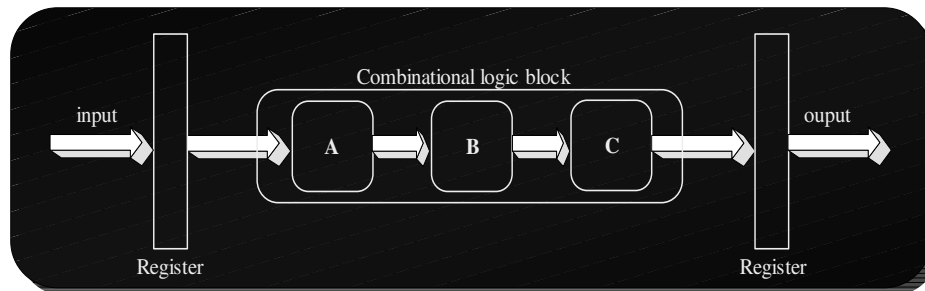- **Validation on adders and multipliers**

41

# Introduction to Wave Pipelining

- **WP is a design method that uses gate capacitance to hold information between successive stages**
- **Allows multiple sets of data to coexist, and thus increases the throughput**
- **Boosts the throughput of a system without additional registers**

42
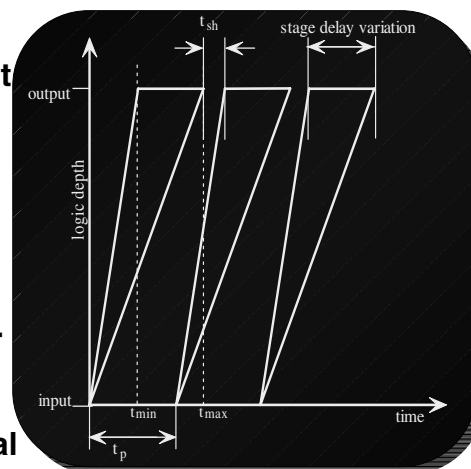
# Clocking Wave pipelined circuits



- **Input applied to *A* at t=0**
- **Output of *A* appears at t=$t_A$**
- **Similarly for *B* and *C* – $t_B$ & $t_C$ respectively**
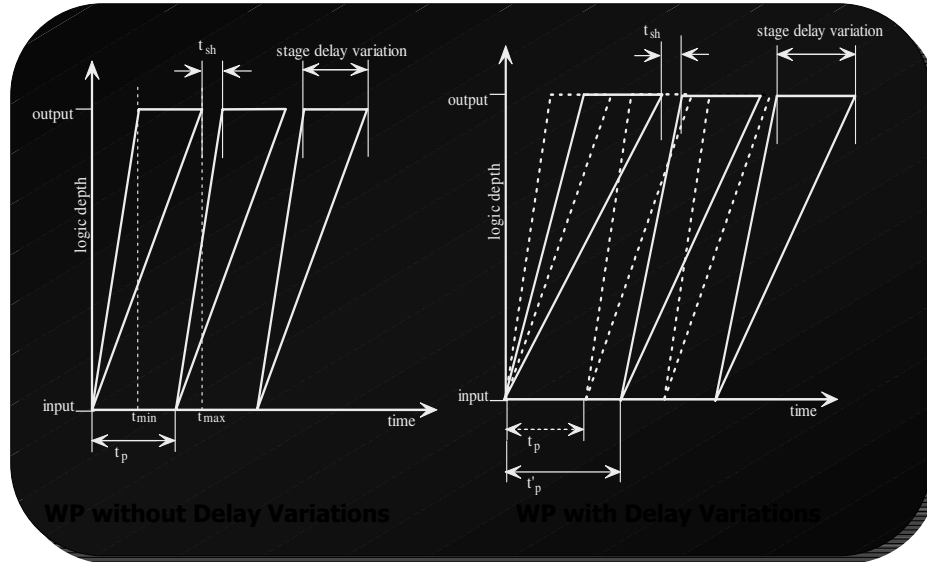
43

# Wave Pipelining – Timing Constraints

- **First input is applied to *A* @ time t=0 and held for at least $t_{SH}$**
- **Delays of *A* - $t_{min}$ & $t_{max}$**
- **New inputs applied at *A* such that new outputs appear not before $t_{max}$**
- **Hence new inputs at *A* should be applied only after time**
  **$t_p > (t_{max} - t_{mix}) + t_{SH} + \Delta$ where $\Delta$ is any unconditional clock skew**



44

# How Delay Variations affect WP



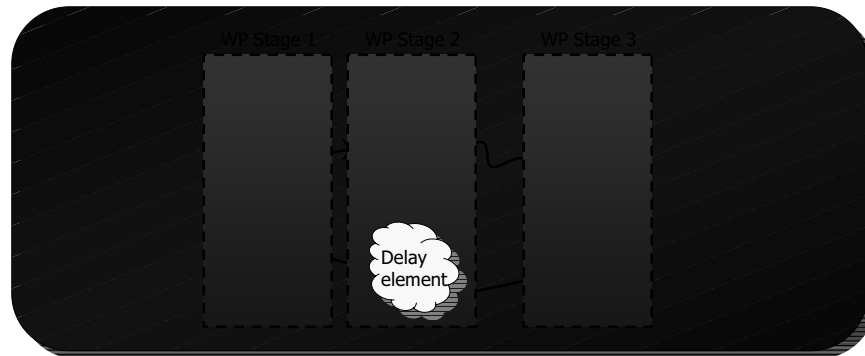WP without Delay Variations   WP with Delay Variations

45

# Wave Pipelining Design Principle

- **The basic design of WP involves equalizing delays between successive stages**

- **The cause of the delay variations could be:**
  - **Difference in propagation paths**
  - **Data dependent delays**
  - **Temperature and process variations**
  - **Power supply drift and noise**
  - **Signal noise such as crosstalk**

46

# Difference in Propagation Path



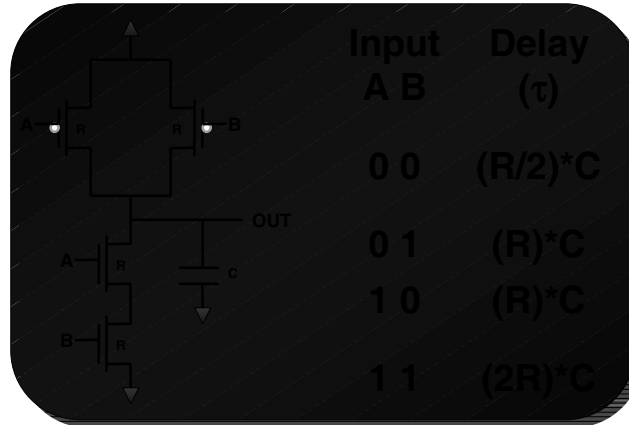- **Can be adjusted by inserting delay elements and gate sizing**

47

# Temperature and Process Variations

- **In WP, incorporating Nowka's delay variation model results in higher number of data waves as compared to the fixed frequency clocking**
- **This is achieved by generating clock signal whose frequency depends on the delay in the logic network**
- **Such a design takes into consideration the process, temperature and any local variations that affect delay**

48

# Data Dependent Delays



| Input A B | Delay ($\tau$) |
|-----------|----------------|
| 0 0 | (R/2)*C |
| 0 1 | (R)*C |
| 1 0 | (R)*C |
| 1 1 | (2R)*C |

- **Data dependent delays can be compensated by appropriate choice of design style**
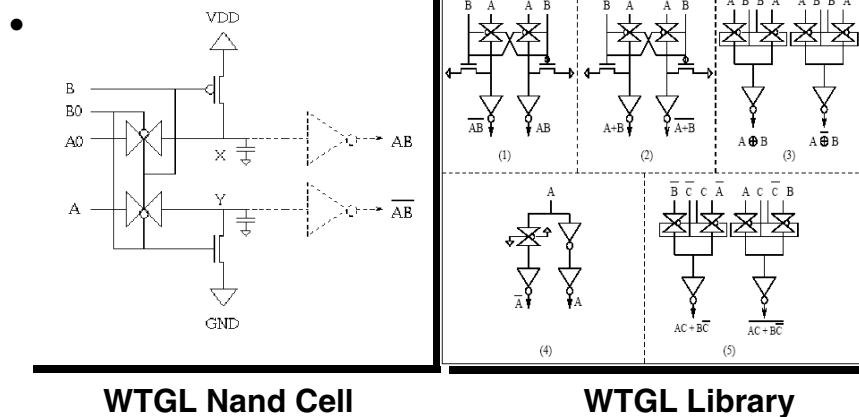
# Logic Design Styles

- **A suitable logic design style should have the following features:**
  - **Equal rise and fall times**
  - **Gate delay independent of current/previous input patterns**
  - **Gate speed adjustment with predictable effects**
  - **High noise immunity, low power, high speed**

# Wave-pipelined Transmission Gate Logic



**WTGL Nand Cell**          **WTGL Library**

51

# Advantages of WTGL Cells

- **Dual-rail logic – less prone to noise and error**
- **Has higher noise immunity**
- **Produces rail-to-rail full swing**
- **Delay easily adjusted by varying output inverter sizes**
- **Every WTGL has equal delay hence balancing easier**
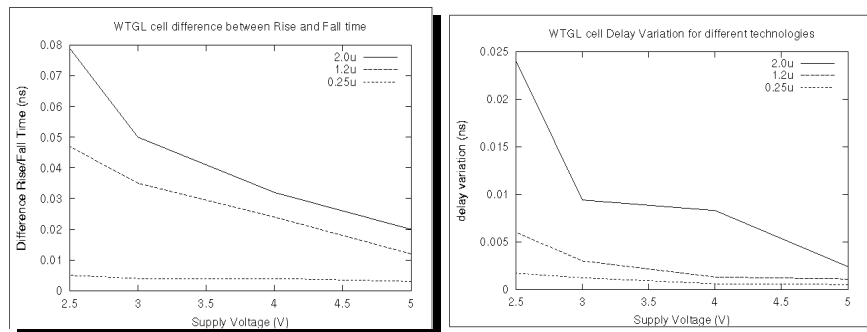- **High throughput and power efficient**

52

# WTGL in UDSM Technologies

- **WTGL has linear delay variation as a function of supply voltage**
  - Reducing supply voltage results in quadratic power saving
- **Rise/fall time also same in different supply voltage in deep sub-micron**
- **Output inverter in WTGL can be sized to have more current driving capability and do fine tuning**
  - In bigger scale, size of the inverter has large effect on power dissipation
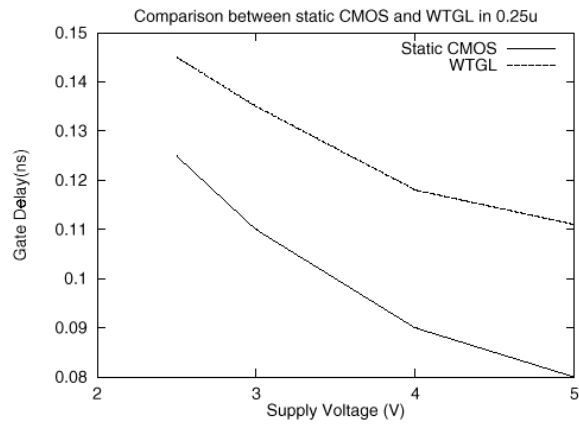
53

# WTGL in UDSM Technologies



•**Performance of WTGL cells is seen to improve as we scale down technology**

54

27

# WTGL in UDSM Technologies

Comparison between static CMOS and WTGL in 0.25u

Gate Delay(ns) vs Supply Voltage (V)

- Static CMOS ———
- WTGL – – – – – –

55

# Comparison of Logic Styles

| Logic Style | $T_{HL}$(ps) | $T_{LH}$(ps) | $T_{fall}$(ps) | $T_{rise}$(ps) | Power in 2um (uW) | Power in 0.25um (uW) |
|---|---|---|---|---|---|---|
| Balanced NAND | 100 | 126 | 230 | 256 | 89 | 7.5 |
| NPCPL | 78 | 99 | NA | NA | 266 | 33 |
| WTGL | 100 | 99 | 70 | 85 | 250 | 6.5 |
| Cross coupled | 64 | 64 | 95 | 85 | 780 | 140 |

56

# Summary of Different Logic Styles

- **Logic Styles compared.**

  – **Balanced NAND: Good power characteristics**

  – **WTGL: Faster with less delay variation.**

- **Logic styles not considered**

  – **NPCPL: Reduced output voltage swing**

  – **Cross Coupled: Large power dissipation**

57

# Implementation Results for a 4:2 Compressor

|  | Delay | Delay Variation | Power (mW) |
|---|---|---|---|
| **Balanced NAND** w/interconnect effect | 0.98ns | 0.032ns | 1.72 |
| **Balanced NAND** w/o interconnect effect | 0.75ns | 0.04ns | 1.33 |
| **WTGL** w/ interconnect effect | 0.19ns | 0.01ns | 1.27 |
| **WTGL** w/o interconnect effect | 0.16ns | 0.018ns | 1.15 |

58

# Results from WP Multiplier circuit

- **Multiplier has 7.8ns propagation delay and delay variation is 1.32ns.**
- **Best clock cycle is 2.3ns.**
- **Wave pipelined multiplier is 3.5 times faster than non pipelined multiplier.**
- **Conventional pipelined multiplier which has 5 stage has 6.4 ns propagation delay.**
  - **Each stage is 1.5n and clocking overhead is 0.9ns.**
  - **Clock period of conventional pipelined multiplier is 3.15ns.**
  - **Wave pipelined multiplier is 35% faster than conventional multiplier.**

59

# Results from WP Adder circuit

- **A 4-bit Brent-Kung WP adder was laid out in 180nm and simulated**

- **Simulations showed a latency of 1.9ns with delay variation of 12%**

- **The layout was enhanced using the above techniques and the new latency was 1.83ns with a delay variation of 6%**

60

# Wave Pipelining Summary

- **WP is a technique that can achieve maximum clock cycle without using additional registers.**
  - In deep sub-micron design, it helps to reduce clock skew and power dissipation by clock distribution.
- **For WP, WTGL cell is proper logic style in deep sub-micron design.**
  - Less input dependent delay variation.
  - Less power consumption.
  - Same rise and fall time.
  - Functionality.
  - Less affected by temperature variation.

61

# WP and SoC design

- **The ideas learnt from WP technique in terms of minimizing delay variations could be applied to different IP cores of SoCs for minimum delay variation between the components**

- **Helps in improving reliability and performance of the overall system**

62

# Inductance Modeling & Estimation

- **Existing inductance modeling reflect only isolate wires**
- **Current techniques not adequate to analyze large complicated networks of RLC trees**
- **Techniques that model and estimate the inductive effects of such large RLC trees necessary to improve system reliability**

63

# Design Approaches

- **Designing circuits that detect and correct or tolerate errors due to wire delay uncertainties**

- **Availability of fabrication data for accurate modeling & estimation**

64

# Design Approaches…

- **Network-on-Chips (NoC) by Luca Benini and DeMichelli's group**
  - **Attempts to model the SoC into different layers similar to the network OSI model**
  - **The interconnects form the lowest layer – the physical layer**
  - **One of ideas behind this concept is routing data rather than routing wires**

65

# Development of New Materials

- **In addition to design approaches, new materials and processes that overcome the problems of inductance, noise and delay uncertainties need to be explored**

- **This helps in conforming to Moore's Law for a longer time**

66

# Using a Multi-level Approach

- **A consorted effort from all design perspectives should be adopted to achieve reliable performance in UDSM SoC designs**
- **A multi level approach that applies the model developed to check the inductive and delay unpredictability effects should be followed at each level of design flow**

67

# SoC Design Flow

```
High level                    Layout
description                   Synthesis

Apply model &                 Apply model &
check timing closure          check timing closure

RTL                           Place &
Synthesis                     Route

Apply model &                 Apply model &
check timing closure          check timing closure
```

68

# Summary

- **Sub-90nm designs create many challenging problems for VLSI designers**
- **A Key challenge is the unpredictable behavior of the interconnect characteristics resulting in delay variations.**
- **New techniques such as current-mode interconnection scheme and results from other circuit domain could be helpful in dealing with this problem.**
- **Also, prevention and correction both should be considered in achieving signal and function integrity**
- **An approach that spans all levels of design should be developed.**

69

# References

- ITRS 2003 Edition, http://public.itrs.net/Files/2003ITRS/Home2003.htm
- http://www.ifw-dresden.de/ifs/31/gfa/em-life.pps
- S. Kim and R. Sridhar, "Hierarchical synchronization scheme using self-timed mesochronous interconnections," in *Proc. ISCAS*, vol. 3, pp. 1824–1827, June 1997
- K. Yun and R. Donohue, "Pausible clocking: A fi rst step toward heterogeneous systems," in *Proc. ICCD*, pp. 118–123, October 1996
- A. Maheswari and W. Burleson, .Current sensing techniques for global interconnects in very deep submicron (VDSM) CMOS,. in *Proceedings of IEEE Workshop on VLSI*, 2001, pp. 66.70
- S. B. Furber, D. A. Edwards, and J. D. Garside, "AMULET3: A 100 MIPS Asynchronous Embedded Processor," in Proc. Intl. Conference on Computer Design (ICCD), 2000.
- T. Nanya, A. Takamura, M. Kuwako, M. Imai, T. Fujii, M. Ozawa, I. Fukasaku, Y. Ueno, F. Okamoto, H. Fujimoto, O. Fujita, M. Yamashina, M. Fukuma :*"TITAC-2: A 32-bit Scalable-Delay-Insensitive Microprocessor",* HOT Chips IX,Stanford, pp.19-32 (Aug. 1997)
- Marly Roncken, Ken Stevens, Rajesh Pendurkar, Shai Rotem, and Parimal Pal Chaudhuri. "CA-BIST for Asynchronous Circuits: A Case Study on the RAPPID Asynchronous Instruction Length Decoder." The Fifth International Symposium on Advanced Research in Asynchronous Circuits and Systems, pp. 62--72, April 2000.
- Hemani, T. Meincke, S. Kumar, A. Postula, T. Olsson, P. Nilsson, J. Oberg, P. Ellervee, and D. Lundqvist. "Lowering Power Consumption in Clock by Using Globally Asynchronous Locally Synchronous Design Style". Proceedings of the 1999 Design Automation Conference, pp. 873-878.
- F. Klass and J. Mulder, .Use of CMOS Technology in Wave Pipelining,. In *Proceedigs of International Symposium on VLSI Design (Bangalore)*, pp. 303.308, January 1992.
- X. Zhang and R. Sridhar, .CMOS Wave Pipelining Using Transmission Gate Logic,. in *Proceedings of IEEE International ASIC Conference and Exhibit*, pp. 92.95, May 1994.
- H. S. Park, .Impact of Deep Sub-Micron Technology on Wave Pipelining,.Master's thesis, The State University of New York at Buffalo, September 2001.
- A. Narasimhan, Crosstalk Tolerant Wave Pipelined Systems, Master's Thesis, The State University of New York at Buffalo, August, 2003

70

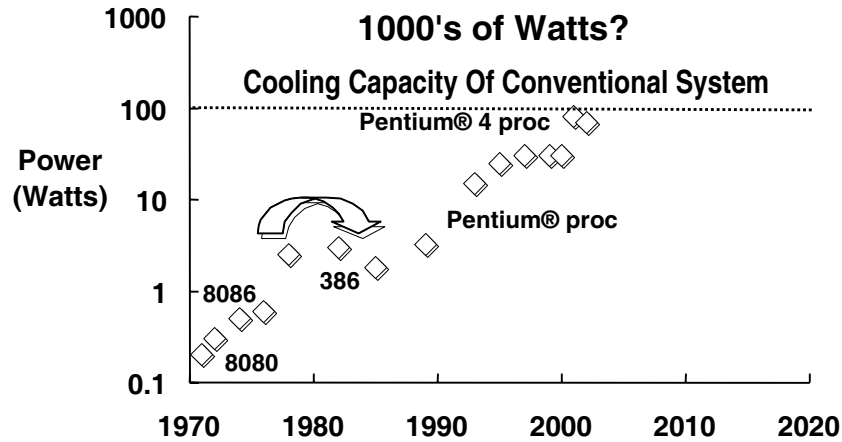# High-Performance CMOS Circuits for Sub-90nm SOC Technologies

**Sanu Mathew and Ram Krishnamurthy**

**Circuit Research, Intel Labs**

**Intel Corporation, Hillsboro, OR**

**Intel** **Labs**

---

# Outline

- **Challenges & Circuit Solutions:**
  - **Leakage power reduction: Dual-Vt and Forward Body Bias**
  - **Stand-by leakage reduction: Sleep transistor design**
  - **Dual-Vcc switching + leakage power reduction**
  - **Dual-Vcc interface: split-output level converters and write-port latches**
  - **Dynamic leakage-tolerant Conditional/burn-in keeper**
  - **Process parameter variation tolerant dynamic circuits**
  - **Pseudo-static & Self Reverse Biased bitlines**
  - **Static split-decoder register file technologies**
  - **Source follower and transition encoded interconnects**

2

# Microprocessor Power Trend



- **C scales by 30% per generation…**
- **…but Vcc scales by 10-15% only**
- **Must maintain or reduce power in future**

3

# Leakage vs. Switching Power



- **$I_{off}$ increase 3-5X per generation**
- **Active leakage power > 50% of total power**
- **Aggressive active leakage control required**

4

# Functionality with High Leakage

**Clock**

**Pk_0**

**Inv_out**

$I_{Leak}$

**Dyn_out**

$M_{11}$  $M_{1j}$  $M_{1K}$

$M_{21}$  $M_{2j}$  $M_{2K}$

Keeper / pulldown ratio

1.6
1.2
0.8
0.4
0

**Sub-70nm**

1X   3X  5X  10X  20X

**Subthreshold + gate leakage**

M. Anders et al, 2001 Symp. VLSI Circuits

**Sub-70nm Dynamic Circuit Active Leakage Tolerance:**

- **Cache, RF, Arrays, Bitlines most affected**
- **Keeper sizes > 50% of pulldown strength**
- **High contention $\Rightarrow$ degraded performance**
- **Slow keeper shutoff $\Rightarrow$ high short-circuit power**

5

---

# Bitline Leakage Tolerance

Current ($\mu A/\mu m$)

2510
2010
1510
1010
510
10

$L_{poly}$ ($\mu m$)  0.18   0.13   0.1   0.07   0.05

Robustness (Noise Margin / $V_{cc}$)

1.2
1
0.8
0.6
0.4
0.2
0

110°C   **Bitline Robustness**

**Bitline $I_{on}$**

**Bitline $I_{off}$ (16 cells)**

- **Bitline $I_{on}/I_{off}$: 60% $\downarrow$ per generation**
- **Leakage tolerant bitline techniques required**

6

3

# On-chip Interconnect RC

**Interconnect Delay**

**30%↑ per generation**

**30%↓ per generation**

**Typical Gate Delay**

Delay (ns): 1.0, 0.1, 0.01, 0.001

Technology Node (nm): 250, 200, 150, 100, 50

- **RC/μm increases 40-60% per generation**
- **Local inter-gate wires dominate critical-path delays**
- **Global wire lengths not scaling by 0.7x**
- **Copper, low-K ILD: modest benefit**

7

# Increasing Number of Repeaters

**# full-chip repeaters (normalized to 0.18um)**

□ **without interconnect lengths scaling by 0.7x**

7, 6, 5, 4, 3, 2, 1, 0

Technology Node (μm): 0.18, 0.13, 0.1, 0.07, 0.05, 0.035

- **RC/μm is only one side of the story…**
- **Finer pipeline global buses ⇒ more flop-repeaters**
- **Exponential increase in bus repeaters aggravates power problem**

8

4

# Process Parameter Variation Tolerance



- **Significant variation in IOFF (hence $F_{max}$ spread)**
- **Worsening with process scaling**
- **Excess leakage dies: lack in robustness**
- **Low leakage dies: over-designed for robustness**

> **Process parameter variation tolerant circuit techniques**

9

# Outline

- **Challenges & Circuit Solutions:**
  - **Leakage power reduction: Dual-Vt and Forward Body Bias**
  - **Stand-by leakage reduction: Sleep transistor design**
  - **Dual-Vcc switching + leakage power reduction**
  - **Dual-Vcc interface: split-output level converters and write-port latches**
  - **Dynamic leakage-tolerant Conditional/burn-in keeper**
  - **Process parameter variation tolerant dynamic circuits**
  - **Pseudo-static & Self Reverse Biased bitlines**
  - **Static split-decoder register file technologies**
  - **Source follower and transition encoded interconnects**

10

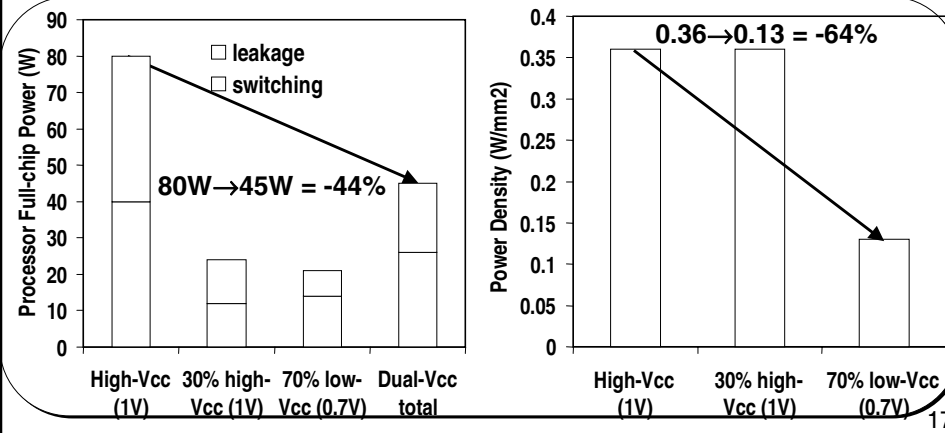# Active Leakage Reduction: Dual $V_t$ design

- **Motivation: Exploit two $V_t$'s provided by process**
  - **High-$V_t$ with nominal $I_{off}$**
  - **Low-$V_t$ with 10X higher $I_{off}$ (15% better delay)**

**Measured leakage (1.2V) 130nm dual-$V_t$ process**



Normalized active leakage vs DIBL (mV/V), showing Low-$V_t$ and High-$V_t$ regions with 10X difference.



**Logic path between latch boundary**

**Goal: selective high-$V_t$ usage for leakage power reduction**

11

---

# Active Leakage Reduction: Dual $V_t$ design



High Vt — Number of paths vs Delay

Low Vt — Number of paths vs Delay

Dual Vt — Number of paths vs Delay

**J. Tschanz et al, 2002 Symp. VLSI Circuits**

- **Methodology:**
  - **Low-$V_t$ on critical paths (best delay)**
  - **High-$V_t$ on non-critical paths: 10X lower leakage**
- **Selective low-$V_t$ and high-$V_t$ insertion enables >1X and <<10X active leakage power**
- **Challenges: Tool-flow, accurate slack estimation**

12

# Standby Leakage Reduction: Sleep Transistor design

- **Motivation: Cut off power supply in sleep-mode**
    - **Insert "sleep" transistor between main supply and functional unit's supply rails**
    - **Latches tied to main supply rails: retain state**

**sleep transistor**

**Virtual V$_{cc}$**

**Functional Unit**

**Virtual V$_{ss}$**

**sleep transistor**

**Standby leakage benefit for 5% delay penalty**

13

---

# Switching + Leakage Reduction: Forward Body Bias

**Vbp**

**Vdd**

**+Ve**

**-Ve**   **Vbn**

**Vcc: 1, 1.05, 1.1 ... 1.5V**

Normalized total power

**110°C**
**α = 0.1**
**ZBB**
**FBB**
**1.2V**   **500mV**
**1.1V**

4
3
2
1
0

0.6  0.8  1  1.2  1.4

**Frequency (GHz)**

FBB/ZBB leakage ratio

30
20
10
0

**27°C**

0.6  0.8  1  1.2  1.4

**Frequency (GHz)**

**A. Keshavarzi et al, 2002 Symp. VLSI Circuits**

**20% power reduction at 1GHz**

**8% ↑ frequency at iso-power**

**20X ↓ idle-mode leakage**

# Outline

- **Challenges & Circuit Solutions:**
  - **Leakage power reduction: Dual-Vt and Forward Body Bias**
  - **Stand-by leakage reduction: Sleep transistor design**
  - **Dual-Vcc switching + leakage power reduction**
  - **Dual-Vcc interface: split-output level converters and write-port latches**
  - **Dynamic leakage-tolerant Conditional/burn-in keeper**
  - **Process parameter variation tolerant dynamic circuits**
  - **Pseudo-static & Self Reverse Biased bitlines**
  - **Static split-decoder register file technologies**
  - **Source follower and transition encoded interconnects**

15

# Switching + Leakage Reduction: Dual Supply Design

- **Motivation: power-optimized performance**
- **High (regular) supply for critical units**
- **Lower supply for non-critical units**
- **Lower $V_{cc}$ generated off-chip or regulated on-die**

| Performance-critical units (high Vcc) | Level Converters | Non-critical units (low Vcc) |
| --- | --- | --- |
| | | Low Vcc on-die regulator(optional) |

16

# Switching + Leakage Reduction: Dual Supply Design

- **Cumulative processor-level power and power density benefit**

- **Challenges: level converters, low-$V_{cc}$ distribution**



Left chart — Processor Full-chip Power (W), with legend "leakage" and "switching", annotation "80W→45W = -44%". Categories: High-Vcc (1V), 30% high-Vcc (1V), 70% low-Vcc (0.7V), Dual-Vcc total.

Right chart — Power Density (W/mm2), annotation "0.36→0.13 = -64%". Categories: High-Vcc (1V), 30% high-Vcc (1V), 70% low-Vcc (0.7V).

17

# Switching + Leakage Reduction: Dual Supply Design

- **Active leakage benefit with lower supply voltage**

- **Exponential subthreshold and gate leakage reduction**



Left chart — **Measured Leakage in 1.2V, 130nm process**, Normalized Leakage vs Voltage (V), with labels "Subthreshold lkg" and "Gate lkg".

Right chart — **130nm L1 cache leakage**, Leakage Energy (Normalized) vs VCC (V), with labels "w.c. corner", "Nominal corner", and "79%".

18

9

# Register File Energy Breakup

**90nm CMOS, 1.2V, 110°C simulation**



Pie chart:
- Dec 8%
- GBL 9%
- LBL 7%
- Clk 4%
- Read/Write Select Drivers 35%
- Bitcells 37%

- **Active leakage = 83% of total energy**
- **Address decoder and read/write select drivers contribution = 43% of total energy**
- **Goal: Lower supply voltage on decoder and read/write select drivers to reduce total energy**

19

---

# Dual Supply Register File

S. Hsu et al, 2003 Symp. VLSI Circuits



Decoder — RS0 · · · RS7 — D0 ... D7 — LBL0 — OUT — GBL0 — GBL3 — GBL4 — GBL7 — 4:1 — 2:1 — 4:1 — D

**Address decoder and Read/write select drivers (Low-Vcc = 0.9V)**

**Bitcells, local and global bitlines (High-Vcc = 1.2V)**

- **Local bitline PMOS keeper: self level converting**
- **Bitcells stability unaffected: operating on high-Vcc**

20

**10**

# Dual-Vcc RF Energy-Delay

**90nm CMOS, 110°C simulation (high-Vcc = 1.2V)**

Worst Case Read Cycle Time (ps) — left axis: 135, 140, 145, 150, 155, 160

Total Energy (pJ) — right axis: 80, 85, 90, 95, 100, 105, 110, 115

6.5GHz

14%    23%

7.3GHz

Decoder and R/W Select Drivers Supply (V): 0.8  0.9  1  1.1  1.2  1.3

- 6.5GHz dual supply operation (7.3GHz at 1.2V)
- 23% total energy reduction with dual supply design

21

---

# Dual Vcc Clocking

**130nm CMOS, 30°C measurements**

Energy (pJ): 0, 100, 200, 300, 400, 500

Global Grid

LCBs

-71%    -49%    -21%

Regular | Global Grid & LCBs low-V$_{cc}$ | Global Grid only low-V$_{cc}$ | LCBs only low-V$_{cc}$

Scan ctl

Output FIFO

Misc

Sched

Clock

Input FIFO

RF

FIFO

ALU

BB ctl

**5GHz 130nm Integer Execution Core**

R. Krishnamurthy et al, 2002 Symp. VLSI Circuits

- Goal: Combat increasing clock power in μP's
- High-V$_{cc}$ (1.2V) on datapath, Low-V$_{cc}$ (0.8V) on clock
- 21-71% clock energy reduction

22

# Dual-Vcc DC Power Impact



**32-bit processor pass-gate latch:**
**DC power = 290μW**

**LCB DC power = 350μW**

- DC power free latches and LCBs required to enable practical dual-Vcc core/cache interface or low-Vcc clocking

23

# Split-output Level Converter LCB



**Conventional CVSL LCB**          **Split-output LCB**

- **Contention in CVSL LCB degrades delay**
- **Split-output LCB decouples CVSL stage from output driver stage**
  - **Fast level conversion due to low contention**
  - **Reduced fanin load on clock grid**

24

**12**

# LCB Energy-Delay Comparisons



High-$V_{cc}$=1.2V, Low-$V_{cc}$=0.8V
130nm CMOS, 30°C simulation

Conventional CVSL

Split-output

-16%

47%

LCB Energy (pJ) vs LCB Delay (ps)

R. Krishnamurthy et al, 2002 Symp. VLSI Circuits

| LCB Scheme | Fanin cap (fF) | Total area (mm²) | CVSL-stage contention energy (pJ) |
|---|---|---|---|
| Conventional CVSL | 8.2 | 15.5 | 0.085 |
| This work | 7.1 (-14%) | 13.8 (-11%) | 0.039 (-54%) |

- **Effective low-energy alternative to CVSL LCB**

25

# Write-port Pass-transistor Latch



low-$V_{cc}$

high-$V_{cc}$    Clk    high-$V_{cc}$

D    0    Q

- **Dense 9T design**
- **No local D# or Clk# inverters:**
  - **Low setup time**
  - **Low output glitch**
- **No DC power penalty**

CLK    D    OUT

Peak glitch: 20mV

Voltage (V) vs Time (ns)

26

13

# Write-port Latch Comparisons

**130nm CMOS, 30°C simulation**

| Latch Scheme | Number of Transistors | D→Q (ps) | O/P glitch (mV) | Energy (pJ) |
|---|---|---|---|---|
| Pass-gate high-$V_{cc}$ | 11 | 37 | 18mV | 0.76 |
| Sense-amp dual-$V_{cc}$* | 11 | 44 | 60mV | 0.59 |
| Write-port dual-$V_{cc}$ | 9 | 38 | 20mV | 0.66 |

**\* H. Kawaguchi et al, JSSC, May 1998**

- **Data activity = 0.1, Clock activity = 1.0**
- **Optimized for constant fanin and fanout load**
- **12% energy reduction vs. pass-gate latch**
- **14% delay reduction vs. sense-amp latch**

27

# Process Skew Sensitivity



**High-$V_{cc}$ = 1.2V, Low-$V_{cc}$ = 0.8V**

- **Comparable delay spread across fast/slow corners**

28

# Supply Variation Sensitivity

**High-$V_{cc}$ = 1.2V $\pm$15%, Low-$V_{cc}$ = 0.8V $\pm$15%**



D-Out Delay Spread (ps)

- Pass-gate high-$V_{cc}$: 45.6 / 37 / 32.8 — **35%**
- Sense-amp Dual-$V_{cc}$: 52.5 / 44 / 39.5 — **30%**
- Write-port Dual-$V_{cc}$: 48.4 / 38 / 33.2 — **39%**

- **Comparable delay spread for supply variations**

29

---

# Adaptive Vcc: Variation-tolerant Circuits

- **Motive: change Vcc adaptively to reduce impact of parameter variations**
  - **Large Fmax vs. Isb spread (worsening with scaling)**
  - **Lower Vcc on leakage-limited circuits (subject to stability limits)**
  - **Higher Vcc on speed-limited circuits (subject to reliability limits)**



5.3 mm

4.5 mm

21 sub-sites within 1 die

**J. Tschanz et al, 2002 Symp. VLSI Circuits**

Fixed Vdd: 1.05V
Adaptive Vdd: 20mV resolution

Adaptive Vdd + body bias
Adaptive Vdd + WID body bias

30

15

# Outline

- **Challenges & Circuit Solutions:**
  - **Leakage power reduction: Dual-Vt and Forward Body Bias**
  - **Stand-by leakage reduction: Sleep transistor design**
  - **Dual-Vcc switching + leakage power reduction**
  - **Dual-Vcc interface: split-output level converters and write-port latches**
  - **Dynamic leakage-tolerant Conditional/burn-in keeper**
  - **Process parameter variation tolerant dynamic circuits**
  - **Source follower and transition encoded interconnects**

31

# Improving Dynamic Leakage Tolerance: Keeper Upsizing



A. Alvandpour et al, 2001 Symp. VLSI Circuits

- **Robustness = DC Noise Margin / Vcc**
- **Traditional noise engineering $\Rightarrow$ diminishing ROI**

32

# Dual Vt Scaling Trends



R. Krishnamurthy et al, 2001 Great Lakes VLSI Symp.

- **Replace NMOS pulldowns with high-$V_t$**
- **Good one-time solution for 130nm node**
- **15-30% degradation for both high- and low-$V_t$ in sub-130nm**
- **Dual-Vt bitlines don't scale well beyond 130nm**

33

# Leakage-tolerant Conditional Keeper Domino Technology



A. Alvandpour et al, 2001 Symp. VLSI Circuits

34

# Leakage-tolerant Conditional Keeper Domino Technology



**130nm 16-wide domino bitline**



- **Motivation:**
- **Weak keeper (low contention) during evaluation window**
- **Strong keeper activated only if dynamic node "high"**
- **20% delay reduction at same robustness**
- **High-performance "dual-$V_t$ enabler"**

35

---

# Burn-in Tolerant Dynamic Circuits



**A. Alvandpour et al, 2002 CICC**

- Leakage sensitive circuits not functional at burn-in
  - Elevated supply and temperature
- Larger keepers increase delay at "normal" condition
- Conditional keeper enables functional burn-in testing
- 2X lower noise during burn-in
- 50% better delay than upsizing BKM keeper

36

# Leakage Variations Impact



**Dynamic 8-way Bitline**



**DC noise robustness**

- **Dynamic circuit NMOS pulldown leakage variation:**
    - **Keeper size determined for target robustness at worst-case leakage corner**
    - **Excess leakage dies: fail to meet target robustness**
    - **Lower leakage dies: over-designed for robustness** 37

# Delay and Robustness Spread

**90nm CMOS, 1.2V, 110°C simulations**



- **Fast corner keeper sizing is sub-optimal for delay** 38

## Variable Strength Keeper Size

**90nm CMOS, 1.2V, 110°C simulations**

Target noise robustness

**23% speedup**

Nominal corner

Worst-case corner

11%
10%
9%
8%

2%

3

2%

4%

1%

5%

2%

Normalized delay

2.5
2.1
1.7
1.3
1.0

DC robustness (normalized to Vcc)

0.13   0.18   0.23   0.28   0.33   0.38

- **Goal: downsize keeper on nominal leakage dies**   39

## Process Compensating Dynamic Circuit Technology

**3-bit programmable conditional keeper**

b[2:0]

clk

W    2W    4W

RS0    RS1    RS7

D0    D1    D7

LBL0

N0

LBL1

C. Kim et al, 2003 Symp. VLSI Circuits

i1
i2
=
i2    i1
s

- **Shared-NAND: 2 less NMOS devices, dense layout**   40

20

# Robustness Squeeze



- **5X reduction in robustness failing dies**

41

# Delay Squeeze



- **10% opportunistic speedup**

42

# Keeper Ratio Distribution



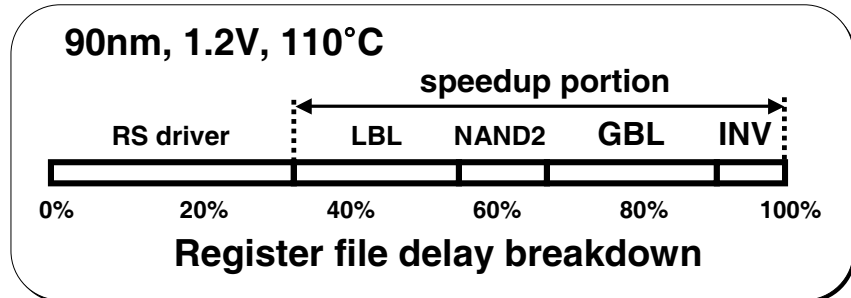- **Keeper downsized in 92% of dies**

# 128x32b 2R2W PCD Register File



- **Single-ended read, 8 bitcells/LBL, 8-way GBL**
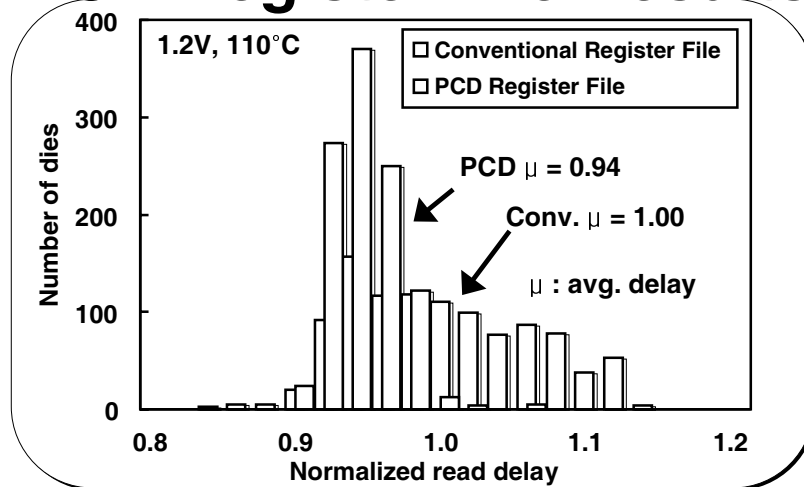- **Keeper folded into existing layout templates**

# PCD Register File Delay, Energy

**90nm, 1.2V, 110°C**

speedup portion

RS driver      LBL    NAND2    **GBL**    INV

0%      20%      40%      60%      80%      100%

**Register file delay breakdown**

- **Speeds up 67% of RF critical path delay**
- **2% worst-case total energy overhead**

45

# PCD Register File Results

**1.2V, 110°C**

☐ **Conventional Register File**
☐ **PCD Register File**

**PCD** $\mu$ = 0.94

**Conv.** $\mu$ = 1.00

$\mu$ : avg. delay

Number of dies: 400, 300, 200, 100, 0

Normalized read delay: 0.8, 0.9, 1.0, 1.1, 1.2

| Read Delay Benefit | 5.5% |
|---|---|
| Robustness Failing Dies | 0.2% (5X ▼) |
| Read Delay Variation: $\sigma/\mu$ | 6.1%→2.3% (2.7X ▼) |

46

# Outline

- **Challenges & Circuit Solutions:**
  - **Leakage power reduction: Dual-Vt and Forward Body Bias**
  - **Stand-by leakage reduction: Sleep transistor design**
  - **Dual-Vcc switching + leakage power reduction**
  - **Dual-Vcc interface: split-output level converters and write-port latches**
  - **Dynamic leakage-tolerant Conditional/burn-in keeper**
  - **Process parameter variation tolerant dynamic circuits**
  - **Pseudo-static & Self Reverse Biased bitlines**
  - **Static split-decoder register file technologies**
  - **Source follower and transition encoded interconnects**

47

# Pseudo-static Bitline Technique



R. Krishnamurthy et al, 2001 Symp. VLSI Circuits

- **Goal: $V_{GS} = -V_{cc}$ and $V_{DS} = 0$ on deselected bitline's access transistors**
- **No oxide stress or additional bias voltages**

48

# Pseudo-static Bitline Technique

**Measured 130nm Leakage**

$V_G=V_B=0V$, $V_D=1.2V$

conventional

703X

Pseudo-static

**Normalized leakage**

1.E+4
1.E+3
1.E+2
1.E+1
1.E+0
1.E-1

-1.2    -0.9    -0.6    -0.3    0

$V_{GS}$ (V)

This work

Dual-Vt

Low-Vt

- **130nm measurement: 703X reduction in bitline leakage (~4 process generations)**
- **Scalable replacement for dual-Vt bitlines**
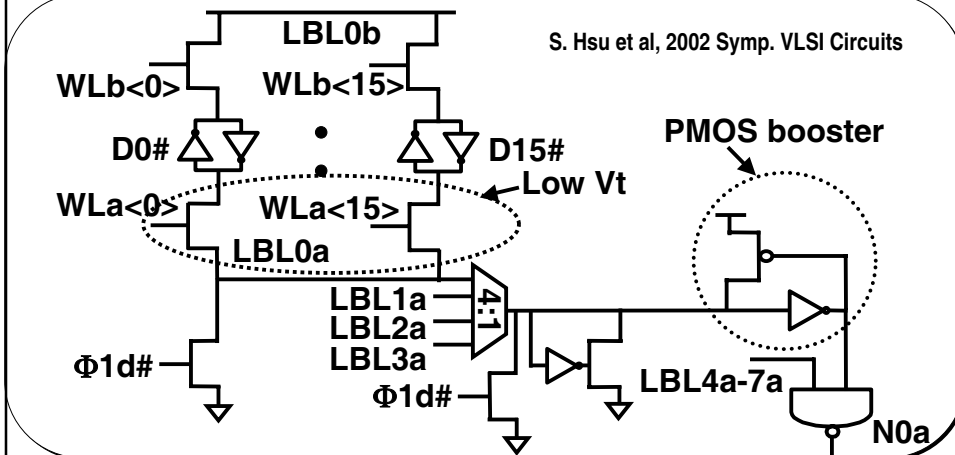
49

# 6GHz 130nm Pseudo-static RF

| LBL Scheme | Read Delay | DC robustness (DC noise margin/Vcc) | Energy/transition (normalized) |
|---|---|---|---|
| | | 130nm, 1.2V, 110C | |
| Low-Vt | 158ps | 0.072 | 1.0 |
| Dual-Vt | 178ps | 0.157 | 0.95 |
| This work | 165ps | 0.214 | 1.02 |

- **256-entryx32-bit 4-read, 4-write ported register file**
- **Single-cycle latency & throughput: performance critical**
- **6GHz operation (8% read delay improvement) with simultaneous 36% robustness benefit over dual-$V_t$**
- **Scalable to sub-130nm technologies**

50

25

## Source Follower Self Reverse-Bias Bitline

**S. Hsu et al, 2002 Symp. VLSI Circuits**

LBL0b

WLb<0>   WLb<15>

**PMOS booster**

D0#   D15#

Low Vt

WLa<0>   WLa<15>

LBL0a

LBL1a
LBL2a
LBL3a

4:1

$\Phi$1d#

$\Phi$1d#

LBL4a-7a

N0a

- **16-way pre-discharged local bitline**
- **NMOS source follower pull-up**
- **PMOS booster ensures full-swing transition**
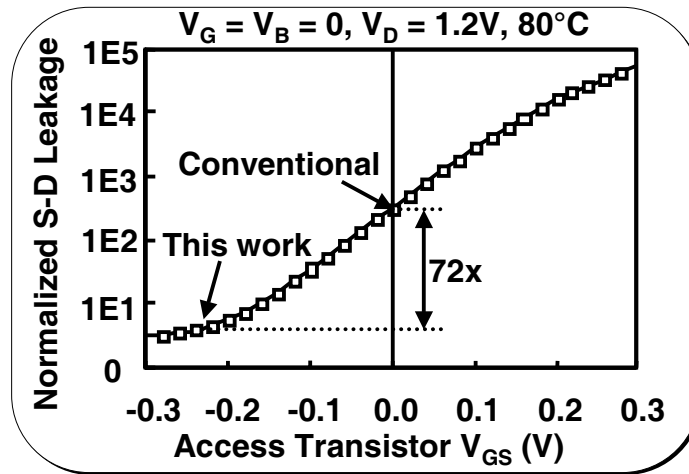- **Bitcell re-optimized for read/write stability**

51

## Self Reverse-Bias Waveform

$V_{cc}$ = 1.2V, 110°C

Voltage (V)

1.2

0.9

0.6

0.3

0

$\Phi$1d#

$T_{settle}$ = 40ps

LBL0a

$V_{GS}$ = -220mV

0   30   60   90   120   150   180

**Time (ps)**

- **$V_{GS}$ = -220mV on access transistors**
- **Bitline settles to its "natural" state**
- **Self-limiting $V_{GS}$:**
  - **Enables low-$V_t$ access transistors**

52

26

# Leakage Reduction Measurement

$V_G = V_B = 0$, $V_D = 1.2V$, $80°C$



- **130nm low-Vt NMOS leakage measurement**
- **72X bitline leakage reduction**
- **7X leakage reduction over dual-V$_t$**

53

# 4.5GHz Single-Cycle L0 Cache using SRB
## 130nm CMOS, 1.2V, 110°C simulation

| LBL Scheme | Read Delay | Energy/ Transition |
|---|---|---|
| Conventional | 247ps | 26.88pJ |
| This work | 220ps (-11%) | 28.14pJ (+5%) |

| LBL Scheme | DC robustness (DC noise margin/Vcc) | LBL droop/rise |
|---|---|---|
| Conventional | 0.114 | 350mV |
| This work | 0.233 | 220mV |

- **11% total read path delay improvement**
- **5% energy overhead due to PMOS boosters**
- **2X simultaneous DC robustness increase**
- **37% reduction in bitline droop / rise**
- **Scalable replacement for dual-Vt bitlines**

54

# Split Decoder Conditional Precharge Scheme

**AD<7:5>**

**3:8**

**BE<7:0>**

**AD** **8**

**BE<0>**

**5:32**

**To Bank0**
**RS/WS<31:0>**
**Column Sel <3:0>**

**AD<4:0>**

**BE<7>**

**5:32**

**To Bank7**
**RS/WS<255:224>**
**Column Sel <31:28>**

S. Hsu et al, 2003 Symp. VLSI Circuits

- **Register file organized into 8 banks x 32 entries**
- **1 first-level 3:8 address decoder:**
  - **Produces 8 Bank Enable signals BE<7:0>**
- **8 second-level 5:32 address decoders (1 per bank)** 55

---

# Conditional Precharge Design

**BE<7:0> (from split decoder)**

**BE<0>**

**AD<3>**
**AD<4>** **LCP0**

**AD<3>**
**AD<4>** **LCP1**

**AD<3>**
**AD<4>** **LCP2**

**AD<3>**
**AD<4>** **LCP3**

**x8 Banks**

**RS0**

**D0**

**LCP0**

**LBL0**

**RS7**

**D7**

**OUT**

- **LCP statically anchors deselected LBL**

56

28

# Conditional Precharge Design

**BE<7:0> (from split decoder)**

BE<0>

AD<3> AD<4> → LCP0
AD<3> AD<4> → LCP1
AD<3> AD<4> → LCP2
AD<3> AD<4> → LCP3

**x8 Banks**

RS0
D0
RS7
D7

LCP0
LBL0
OUT

LBL0 OUT
Column Sel<0>#
Column Sel<0>
BE<0>

GBL0

LBL3 OUT
Column Sel<3>#
Column Sel<3>

GBL3
GBL4
GBL7

4:1
4:1
2:1 — D

- **LCP statically anchors deselected LBL**
- **BE statically anchors deselected GBL**
- **Deselected LBLs and GBLs strongly held at "1"**  57

# Worst-case Leakage Vector

0
0
$I_{leak}$
LCP0 = 0
0
0
1
$I_{leak}$
OUT

1
1
0
$I_{leak}$
BE<0> = 0
1
1
1
0
$I_{leak}$
GBL3
GBL4
GBL7

4:1
4:1
2:1 — D

- **Conditional precharge layout area overhead <1%**
  - **Extra device folded into keeper layout template**
  - **Precharge control signals routed on upper layer** 58

**29**

# 6.5GHz Single-cycle Split Decoder RF
## 90nm CMOS, 1.2V, 110°C simulation

| LBL/GBL Scheme | DC robustness (Deselected BL) | DC robustness (One entry enabled) | Noise recovery time |
|---|---|---|---|
| Conventional static | 220mV | 310mV | 23ps |
| Conventional dynamic | 192mV | 196mV | Non-recoverable |
| This work | 365mV | 312mV | 23ps |

- **DC robustness for deselected bitline:**
  - **65% (90%) higher than static (dynamic) bitlines**
- **Same DC robustness as static with 1 entry enabled**
- **Full recovery from AC noise (same as static)**
- **1% read delay penalty due to excess diffusion cap**

59

# Outline

- **Challenges & Circuit Solutions:**
  - **Leakage power reduction: Dual-Vt and Forward Body Bias**
  - **Stand-by leakage reduction: Sleep transistor design**
  - **Dual-Vcc switching + leakage power reduction**
  - **Dual-Vcc interface: split-output level converters and write-port latches**
  - **Dynamic leakage-tolerant Conditional/burn-in keeper**
  - **Process parameter variation tolerant dynamic circuits**
  - **Pseudo-static & Self Reverse Biased bitlines**
  - **Static split-decoder register file technologies**
  - **Source follower and transition encoded interconnects**

60

# Interconnect Reality



$$T_{clk} = T_{clk-Q} + N*T_{bitslice} + T_{setup} + T_{skew}$$

Unit "B"

Unit "A"

- Microprocessor global bus limitations:
- Ideal world: $T_{repeater} = T_{RCsegment} = T_{bitslice}/2$
- Real world:
- Don't get repeaters where you want
- Floorplan decides repeater locations
- Performance << ideal $T_{clk}$

61

# PMOS-Boosted Source Follower



IN

OUT

$C_L$

Output Voltage (V)

Input

PSF Out

CMOS Out

Time (ns)

Early strike effect

R. Krishnamurthy et al, 2001 VLSI Circuits Symp.

- **Source follower NMOS to accelerate driver strength**
  - **NMOS begins fast pullup $\Rightarrow$ "early" strike effect**
  - **PMOS follow-through completes full-swing transition**
- **Full-swing CMOS driver robustness**

62

# L1 Cache Bus Results



- **8% driver delay reduction for same fanin and area**
- **10% simultaneous peak current reduction (7% lower decoupling capacitance)**
- **Effective alternative to upsizing CMOS drivers**

63

---

# Background: Conventional Dynamic Bus



- Domino timing applied to interconnect
- Monotonic transitions
  - **Reduced collinear capacitance**
    - **Static (worst case) = 2X**
    - **Dynamic (worst case) = 1X**
  - **Φ2 repeater required – susceptible to noise**
- Higher transition activity when input = 1
- Static CMOS inverters drive all segments

64

# Dynamic Bus Advantages

| | Worst Capacitance | Worst Inductance | |
|---|---|---|---|
| **Static** | $2C_L+C_O$ | | **Adds** |
| **Dynamic** | $C_L+C_O$ | | **Subtracts** |

- Capacitance effects reduced
  - **Collinear capacitance reduced 2X**
  - **Orthogonal capacitance unchanged**
- Inductance effects reduced
  - **Can oppose transition for static bus**
  - **Can reduce capacitive effects for dynamic bus**

65

# Transition-Encoded Bus

**M. Anders et al, 2002 VLSI Circuits Symp.**

D1 — encode — decode — FF

$\Phi1$  $\Phi2$  $\Phi1$

- Encoder circuit
  - **XOR of previous and current input**
  - **Domino compatible output**
- Decoder circuit
  - **XOR of previous output and bus state**

66

# Transition-Encoded Bus



40
Delay Improvement %
30
20
10
0
−10

0.18μm Pentium® 4 Simulations

1    6    11    16

Length (mm)    M. Anders et al, 2002 VLSI Circuits Symp.

- Transition only when current input != previous input
- Dynamic bus performance but energy profile of static bus
- Energy scales linearly with input switching activity
- 79% of full-chip buses: 10%-35% delay improvement

67

---

# TED Bus Advantages



D1  encode  ⊳∘〰 • • • D1  ⊳∘〰 • • •  decode  FF
Φ1          Φ2                        Φ1

- Dynamic bus performance improvement
  - **Collinear capacitance reduction**
- Static bus energy
  - **Transition dependent switching activity**
- Noise-insensitive Φ2 repeater required
  - **Regains noise immunity of CMOS inverter**

68

# Encoder Circuit



- Domino-compatible XOR of current and previous inputs
- Small FF holds previous state

69

# Decoder Circuit



- Original data recovered
  - **When Bus = 1, Out transitions**
  - **When Bus = 0, Out does not change**
- Flip-flop must be initially reset
- No clock required for decode

70

# Noise Insensitive Φ2 Repeater



- Sized for $V_{CC}/2$ trip point
  - **Same as static CMOS inverter**
- Latching behavior of domino repeater
- Minimal delay penalty

71

# Energy Comparison



Enables delay, energy improvement for performance-critical buses

72

# Peak Current Comparison



Peak Current (normalized) vs Delay (normalized), showing "Static" and "Transition-encoded" curves. Labeled: 9mm metal3, 130nm process, 1.2V, 30ºC

Enables peak current reductions
across all delays due to smaller transistors

73

# Advantages Summary

|  | Equal delay | Equal transistor width | Equal driver size |
|---|---|---|---|
| Delay reduction | 0% | 19% | 22% |
| Total transistor width reduction | 32% | 0% | -20% |
| Peak current reduction | 49% | 30% | 17% |
| Energy increase | 9% | 16% | 19% |

● Averaged over 3-9mm buses

● Metal3 in 130nm technology, 1.2V, 30ºC

74

# Summary

- Performance demand continues, barriers: Power, Leakage, Interconnect, Variations
- To get E.D = $(0.7)^4$, $V_t$ has to scale aggressively
- Active & standby leakage reduction strategies
- Leakage-tolerant bitline technologies
- Process variation tolerant circuit technologies
- High-speed SF & transition encoded interconnects

75

# Outline

- **Motivation**
- **High-performance Adders**
    - **6.5GHz 32-bit Han-Carlson ALU**
    - **4GHz 32-bit Address Generation Unit**
- **Low-power Multiplier**
    - **1 GHz 16x16 multiplier**

76

# High-performance trends

100000

10000

1000

MHz 100

10

1

0.1

Pentium® 4 proc

-30 GHz

Pentium® proc

386

8086

8080

1970    1980    1990    2000    2010    2020

- ● **Frequency doubles every generation**
- ● **Performance-critical units** ⎫ **Single-cycle**
  - ● **ALUs & AGUs**           ⎬ **throughput**
  - ● **Register files, L0 Caches** ⎭ **& latency**

77

# Motivation

Cache

Processor
thermal
map

Temp
(ºC)

108
105.2
102.4
99.53
96.7
93.86
91.03
88.2
85.36
82.53
79.7
76.87
74.03
71.2
68.37
65.53

Execution
core

ALUs &
AGUs

- ● **ALUs: performance and peak-current limiters**
- ● **High activity ⇒ thermal hotspot**
- ● **Goal: high-performance energy-efficient design**

78

# 32-bit ALU architecture



Mux control   Shift control

External operands → 6:1 Mux → 5:1 Mux → Adder core → O/p Mux → Sum

External operands → 6:1 Mux → 2:1 Mux → Adder core

Mux control   Sign control   Loopback bus

**Multiple ALUs clustered together in the execution core ⇨ High power density**

---

# A 6.5GHz, 130nm Single-ended Dynamic ALU

**[M. Anders et al, ISSCC 2002]**

Intel Labs

$$\text{Sum}_i = A_i \oplus B_i \oplus \text{Carry}_{i-1}$$

$$\text{Carry}_i = A_i \cdot B_i + (A_i + B_i)\text{Carry}_{i-1}$$

Intel Labs

---

**Partial Sum**

$$\text{Sum}_i = \overbrace{A_i \oplus B_i} \oplus \text{Carry}_{i-1}$$

$$\text{Carry}_i = A_i \cdot B_i + (A_i + B_i)\text{Carry}_{i-1}$$

Intel Labs

**Partial Sum**

$$\text{Sum}_i = A_i \oplus B_i \oplus \text{Carry}_{i-1}$$

$$\text{Carry}_i = A_i \cdot B_i + (A_i + B_i)\text{Carry}_{i-1}$$

**Generate**    **Propagate**

---

**Partial Sum**

$$\text{Sum}_i = A_i \oplus B_i \oplus \text{Carry}_{i-1}$$

$$\text{Carry}_i = A_i \cdot B_i + (A_i + B_i)\text{Carry}_{i-1}$$

**Generate**    **Propagate**

$$\text{Carry}_i = G_i + P_i \cdot \text{Carry}_{i-1}$$

# 32-bit Kogge-Stone Adder

**PG**

| 31 | 30 | 29 | 28 | 27 | 26 | 25 | 24 | 23 | 22 | 21 | 20 | 19 | 18 | 17 | 16 | 15 | 14 | 13 | 12 | 11 | 10 | 9 | 8 | 7 | 6 | 5 | 4 | 3 | 2 | 1 | 0 |

**Carry-merge gates**

**XOR**

- **Critical path = PG+5+XOR = 7 gate stages**
- **Generate,Propagate fanout of 2,3 ⎫ Energy**
- **Maximum interconnect spans 16b ⎰ inefficient**

85

---

# 32-bit Han-Carlson adder core

| b63 | b62 | b61 | b60 | b59 | PG generator | b3 | b2 | b1 | b0 | 2N |

Carry-merge0 — 2P

Carry-merge1 — 2N

Carry-merge4 — 2N

Odd carry generator — 2P

Sum XOR

| Odd bit | Even bit |
|---------|----------|
| ▽ | $CM_0$ |
| ▽ | $CM_1$ |

- **Carry-merge done on even bitslices**
- **50% fewer carry-merge gates vs Kogge-Stone**
- **Extra logic stage generates odd carries**

86

# Han Carlson carry-merge tree



- **Single rail adder core**
- **CSG circuit generates dual-rail carry**

# Complementary signal gen.



- **Domino-compatible Carry/$\overline{\text{Carry}}$**
- **Permits a single-rail carry-merge tree design**
- **Not time-borrowable – Penalty absorbed by placing gate at $\Phi_2$ boundary**

# CSG: Timing Diagram

Evaluation Phase

Pre-charge Phase

$\Phi 2$

**Input is setup during precharge**

$Cin_i$

$\overline{Carry}_i$

$Carry_i$

# CSG: Simulation waveforms

Complementary Node

Clk    True Node

**(a) Cin=0**

True Node

Clk    Complementary Node

**(b) Cin=1**

# Partial sum generator



- **Generates domino-compatible partial sum**
- **Placing the gate at $\Phi_1$ boundary mitigates output noise-glitches**

# Dynamic XNOR: Simulation Waveforms

# Dyn. XNOR: Noise Sensitivity



**Fast process corner**

**Slow process corner**

- Mismatch in input evaluation times can cause output noise

# Han-Carlson ALU Organization



- Single-rail dynamic 9-stage low-$V_t$ design

# Odd-bits CSG Sum Generation



- **Final carry-merge CSG(dual-rail carry output)**
  - → **pass-transistor sum XOR**

95

# Even-bits CSG Sum Generation



- **Domino-compatible sum**
- **Dual-rail sum from single-ended g inputs**

96

# Die Micro-photograph

- **130nm 6-metal dual-Vt CMOS**

- **Scheduler:**

  - **210µm x 210µm**

- **ALU:**

  - **84µm x 336µm**



Scheduler

ALU

---

# Delay and Power Measurements



25ºC

$F_{max}$ (GHz)

Supply Voltage (V)

25ºC

Power (mW)

Leakage Power (mW)

Design target

Supply Voltage (V)

- **6.5GHz at 1.1V, 25ºC**
- **Power: 120mW total, 15mW leakage**
- **Scalable to 10GHz at 1.7V, 25ºC**

## Improvements Over Dual-rail Domino

| Area | 50% |
|---|---|
| Performance (Delay) | 10% |
| Active Leakage | 40% |
| Robustness | equal |

- Leakage reduced by eliminating dual-rail logic
- Robustness not compromised
- CSG improves both area and performance

# A 4GHz 130nm Address Generation Unit with 32-bit Sparse-tree Adder Core

**[S. Mathew et al, VLSI Symp. 2002]**

Intel Labs

# Outline

- **Address Generation Unit (AGU) organization**
- **Sparse-tree adder core**
- **Dual-$V_t$ semi-dynamic design**
- **Sub-130nm scaling trends**
- **Summary**

101

# AGU Architecture



- **Single-cycle latency and throughput**
- **Effective Address = Base + Index*Scale +**
  **(Segment +Displacement)**
- **2-phase address computation**

102

51

# AGU Operation: Phase 1



- **Index pre-scaled via 3-bit barrel shifter**
- **3:2 compressor renders partial address:**
  - **Carry-save format**
- **Adder in pre-charge state**

103

# AGU Operation: Phase 2



- **Carry-save to binary format conversion:**
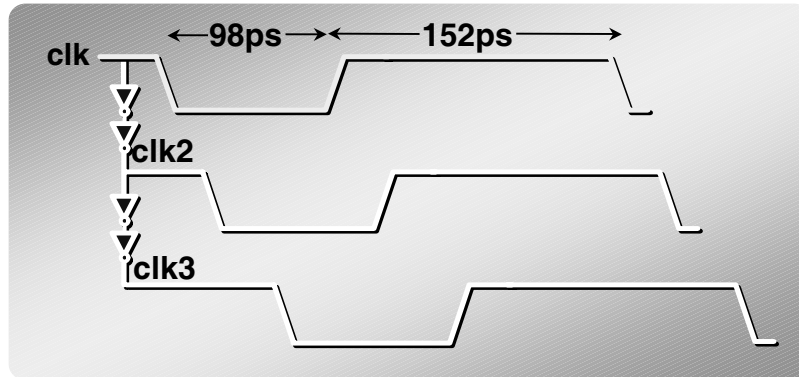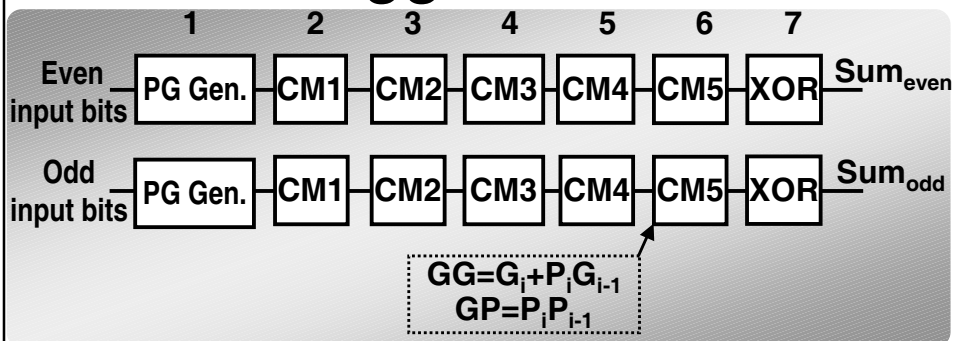  - **2's complement parallel 32-bit adder**

104

# Timing Diagram

**clk** ←98ps→

**Latches open**

**3:2 Compression done**

# Timing Diagram

**clk** ←98ps→ ←152ps→

**Adder inputs setup here**

**Addition done**

# Timing Diagram



- **Seamless time-borrowable clock boundaries**
- **152ps (6.6GHz) 32-bit adder core required**

# High-performance Adders: Kogge Stone



| | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|---|---|---|---|---|---|---|---|
| **Even input bits** | PG Gen. | CM1 | CM2 | CM3 | CM4 | CM5 | XOR | $Sum_{even}$ |
| **Odd input bits** | PG Gen. | CM1 | CM2 | CM3 | CM4 | CM5 | XOR | $Sum_{odd}$ |

$$GG = G_i + P_i G_{i-1}$$
$$GP = P_i P_{i-1}$$

- **Generate all 32 carries:**
  - **Full-blown binary tree $\Rightarrow$ energy-inefficient**
- **# Carry-merge stages = $\log_2(32) \Rightarrow$ 5 stages**

# Kogge-Stone Adder



- **Critical path = PG+5+XOR = 7 gate stages**
- **Generate,Propagate fanout of 2,3** ⎫ **Energy**
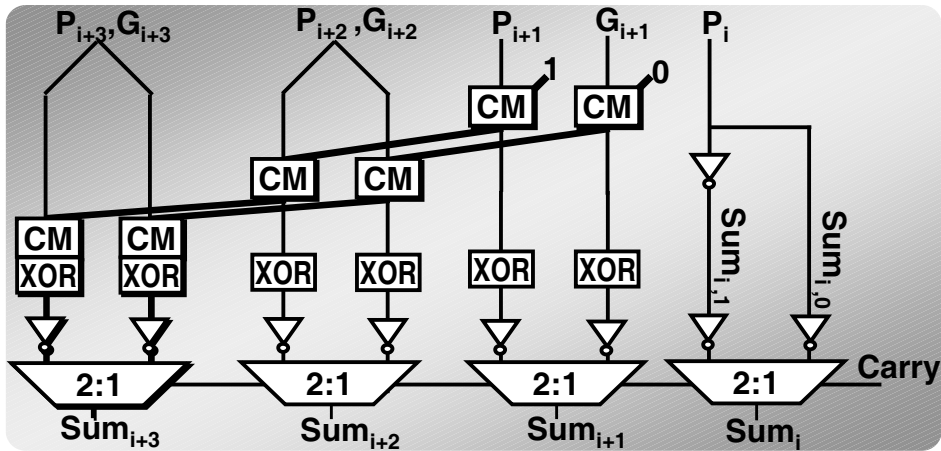- **Maximum interconnect spans 16b** ⎰ **inefficient**

109

# Sparse-tree Adder Architecture



$C_{23}$    $C_{19}$    $C_{15}$    $C_{11}$    $C_7$    $C_3$

- **Generate every 4th carry in parallel**
- **Side-path: 4-bit conditional sum generator**
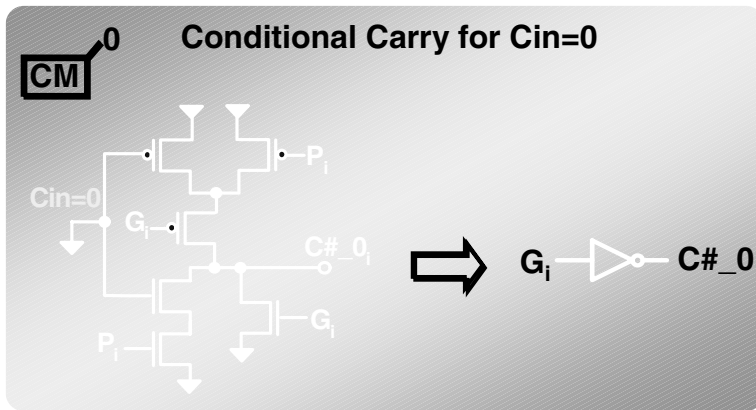- **73% fewer carry-merge gates⇒energy-efficient**

110

55

# Non-critical Sum Generator



- **Non-critical path: ripple carry chain**
- **Reduced area, energy consumption, leakage**
- **Generate conditional sums for each bit**
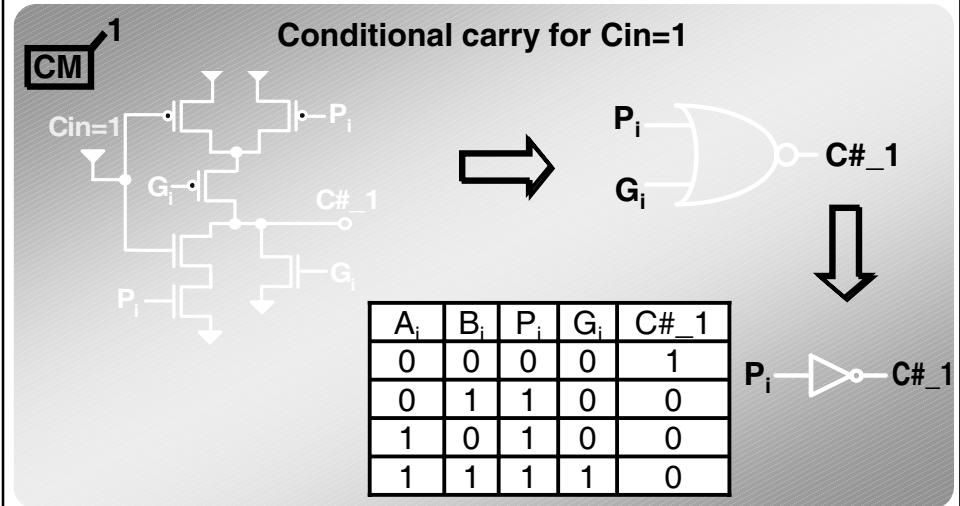- **Sparse-tree carry selects appropriate sum** 111

# Optimized First-level Carry-merge



- **Carry-merge stage reduces to inverter**
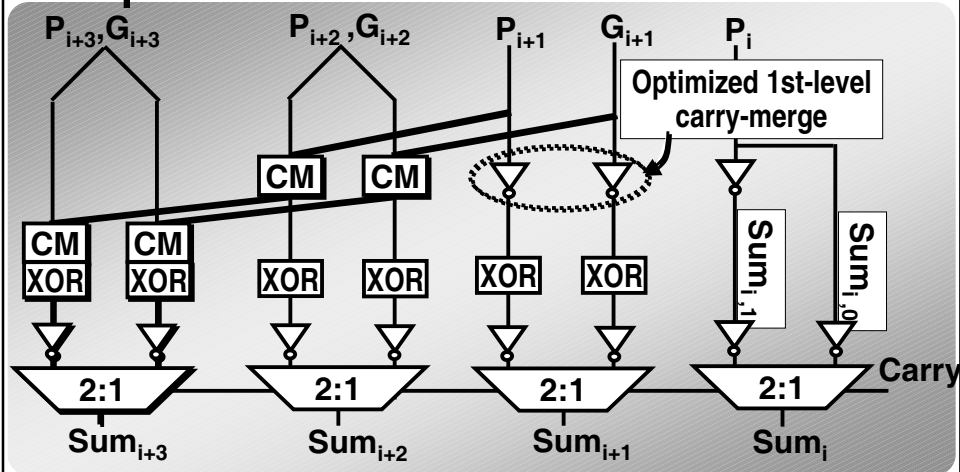- **Conditional carry_0 = $G_i$#**

112

# Optimized First-level Carry-merge

**CM** **1**

**Conditional carry for Cin=1**

Cin=1

$P_i$

$G_i$

C#_1

$P_i$

$P_i$

$G_i$

$G_i$

C#_1

$P_i$ — C#_1

| $A_i$ | $B_i$ | $P_i$ | $G_i$ | C#_1 |
|-------|-------|-------|-------|------|
| 0 | 0 | 0 | 0 | 1 |
| 0 | 1 | 1 | 0 | 0 |
| 1 | 0 | 1 | 0 | 0 |
| 1 | 1 | 1 | 1 | 0 |

- **$P_i$ & $G_i$ correlated**
- **Conditional carry_1 = $P_i$#**

113

---

# Optimized Sum Generator

$P_{i+3}, G_{i+3}$   $P_{i+2}, G_{i+2}$   $P_{i+1}$   $G_{i+1}$   $P_i$

**Optimized 1st-level carry-merge**

CM    CM

CM    CM    XOR    XOR    XOR    XOR    $Sum_{i,1}$    $Sum_{i,0}$

XOR   XOR

2:1    2:1    2:1    2:1    **Carry**

$Sum_{i+3}$    $Sum_{i+2}$    $Sum_{i+1}$    $Sum_i$

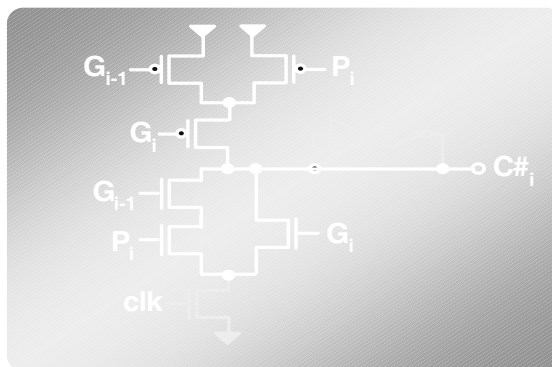- **Optimized non-critical path: 4 stages**

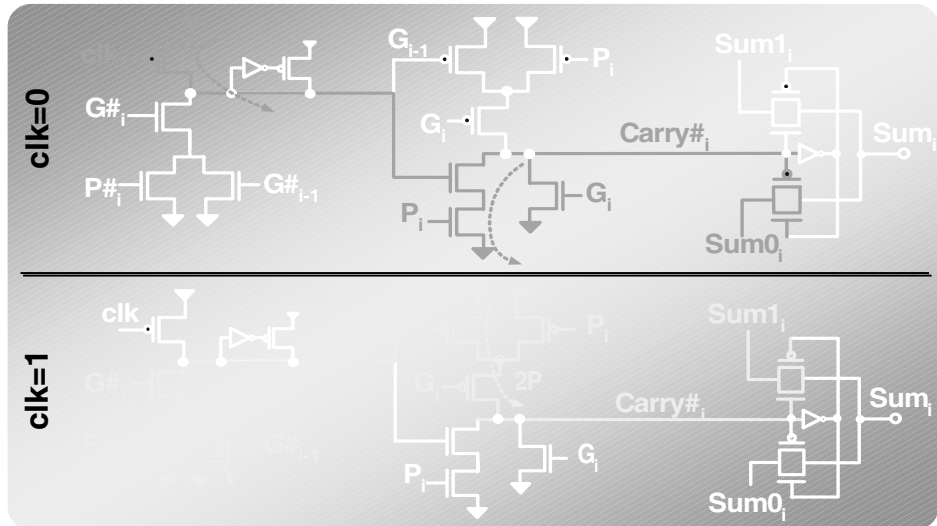114

# Adder Core Critical Path



- •Critical path: 7 gate stages $\Rightarrow$ same as KS
- •Sparse-tree: single-rail dynamic
- •Exploit non-criticality of sum generator
- •Convert to static logic$\Rightarrow$Semi-dynamic design

# 1ˢᵗ-level Carry-merge: Static Latch



- • Holds state in pre-charge phase
- • Prevents pre-charging of static stages

# Domino-Static Interface



- **Sum=Sum0 during pre-charge**
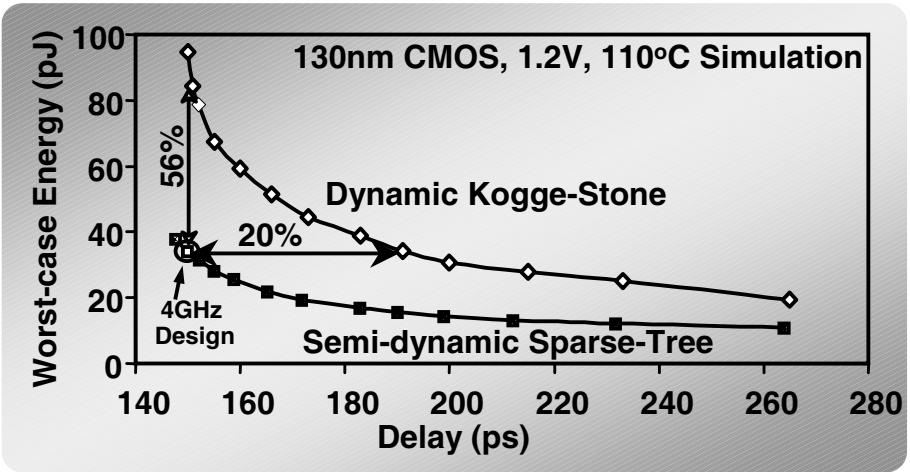- **Mux output resolves during evaluation**

# Sparse-tree Architecture

- **Performance impact: (20% speedup)**
  - **33-50% reduced G/P fanouts**
  - **80% reduced wiring complexity**
  - **30% reduction in maximum interconnect**

- **Power impact: (56% reduction)**
  - **73% fewer carry-merge gates**
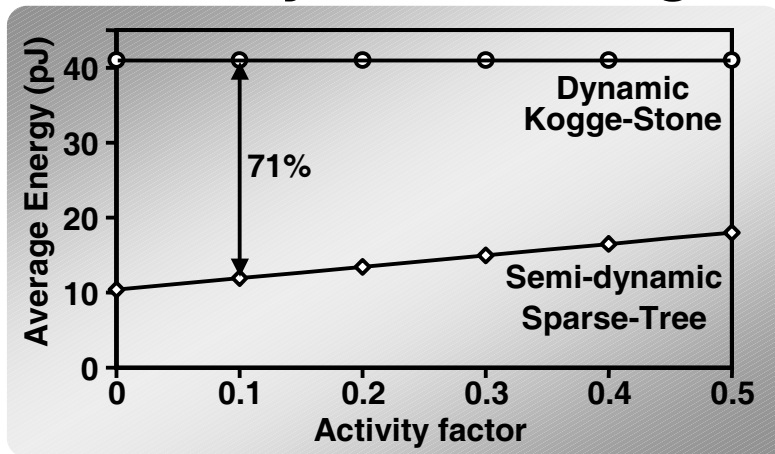  - **50% reduction in average transistor size**

# Energy-delay Space



Chart: Worst-case Energy (pJ) vs Delay (ps)

**130nm CMOS, 1.2V, 110ºC Simulation**

- Dynamic Kogge-Stone
- Semi-dynamic Sparse-Tree
- 56%
- 20%
- 4GHz Design

Y-axis: Worst-case Energy (pJ): 0, 20, 40, 60, 80, 100
X-axis: Delay (ps): 140, 160, 180, 200, 220, 240, 260, 280

- **20% speedup over Kogge-Stone**
- **56% worst-case energy reduction**
  - **Scales with activity factor**

119

# Semi-dynamic Design



Chart: Average Energy (pJ) vs Activity factor

- Dynamic Kogge-Stone
- Semi-dynamic Sparse-Tree
- 71%

Y-axis: Average Energy (pJ): 0, 10, 20, 30, 40
X-axis: Activity factor: 0, 0.1, 0.2, 0.3, 0.4, 0.5

- **Static sum generators : low switching activity**
- **71% lower average energy at 10% activity**

120

# Dual-V$_t$ Allocation

| 130nm CMOS, 1.2V, 110$^o$C Simulation | | |
|---|---|---|
| | Low-V$_t$ | Dual-V$_t$ |
| Delay | 152ps | 152ps |
| Switching Energy | 36pJ | 34pJ (-6%) |
| Leakage Energy | 0.9pJ | 0.4pJ (-56%) |

- **Exploit non-criticality of sidepaths**
  - **Use high-V$_t$ devices**
- **0% performance penalty**
- **56% reduction in active leakage energy**

121

# Scaling Performance

| | 130nm | 100nm |
|---|---|---|
| Delay | 152ps | 102ps (-33%) |
| Switching Energy | 36pJ | 18pJ (-50%) |
| Leakage Energy | 0.9pJ | 0.7pJ (-23%) |

- **Average transistor size = 3.5$\mu$m**
  - **Reduces impact of increasing leakage**
- **80% reduction in wiring complexity**
  - **Reduces impact of wire resistance**
- **33% delay scaling, 50% energy reduction**

122

# Summary

- **4GHz AGU in 1.2V, 130nm technology**

- Sparse-tree adder architecture described

- 20% speedup and 56% energy reduction

- Semi-dynamic design:
  - **Energy scales with switching activity**

- Dual-$V_t$ non-critical paths:
  - **Low active leakage energy**

- **6.5GHz ALU and scheduler at 1.1V, 25ºC**
  - Scalable to 10GHz at 1.7V, 25ºC

123

# A 90nm 1GHz 22mW 16x16-bit 2's Complement Multiplier for Wireless Baseband
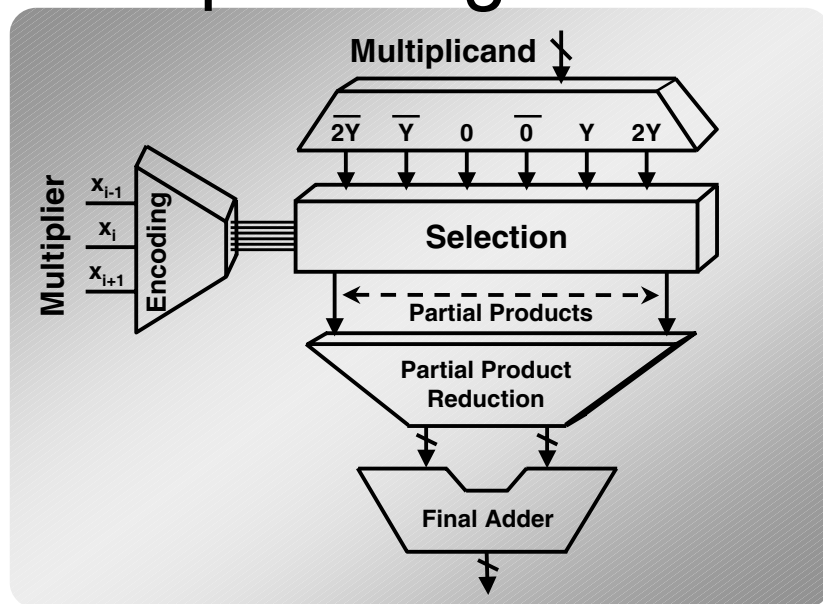
Zeydel et al. VLSI Symp. 2003

**Intel** **Labs**

# Outline

- **Multiplier Block Diagram**

- **Booth Encoding and Select**

- **Optimized Partial Product Reduction**

- **Signal Arrival Optimized Final Adder**

- **90nm Energy-Delay Results**

- **Conclusions**

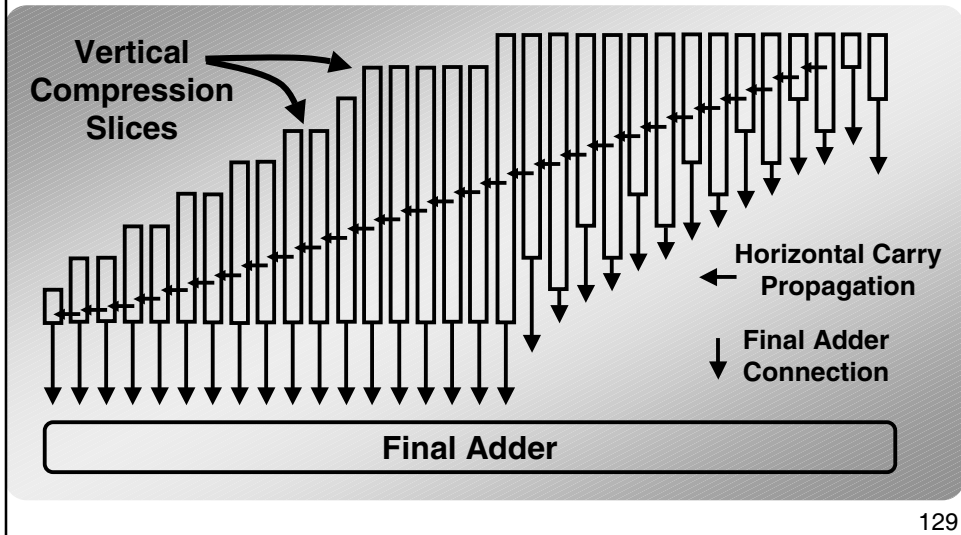# Multiplier Organization

# Booth Encoding and Selection



127

# Booth Encoding Circuits



128

# Vertical Partial Product Compression Tree

**Vertical Compression Slices**

**Horizontal Carry Propagation**

**Final Adder Connection**

**Final Adder**

129

# Partial Product Reduction Tree

**Vertical Compression Slices**

0 0 1 1 0 1 0

FA    FA

FA

FA

FA

**Time**

130

**65**

# PPRT 3:2 Compressor



Isolates 3-stack
From output load

**0.8V 110° C**

| Input | Worst-case Delay |
|-------|------------------|
| Sum   | 240ps (1x)       |
| Carry | 165ps (0.69x)    |

131

# Final Adder



8-bit
Conditional Sum

16-bit
Variable Block
Adder

8-bit
Ripple Carry
Adder

Input Arrival Delay
(Normalized)

Bit Position

132

**66**

# Fast Ripple Full Adders

# Carry Skip Adder Organization

# Energy Savings



Worst Case Total Energy (Normalized)

Conventional:
- Final Adder
- PPRT
- Booth Encoder
- Clock

This Work:
- Final Adder
- PPRT
- Booth Encoder
- Clock

20%
same
15%
14%

Total Savings = 12%

135

# Energy Based Delay Optimization



Initial Sizing Energy Distribution:
- Partial Product Reduction 23%
- Booth Select 18.5%
- Booth 6.7%
- Final Adder 8%
- Clk 43%

Energy Distribution For Final Sizing:
- Booth Select 18%
- Partial Product Reduction 23%
- Booth 12%
- Final Adder 11%
- Clk 35%
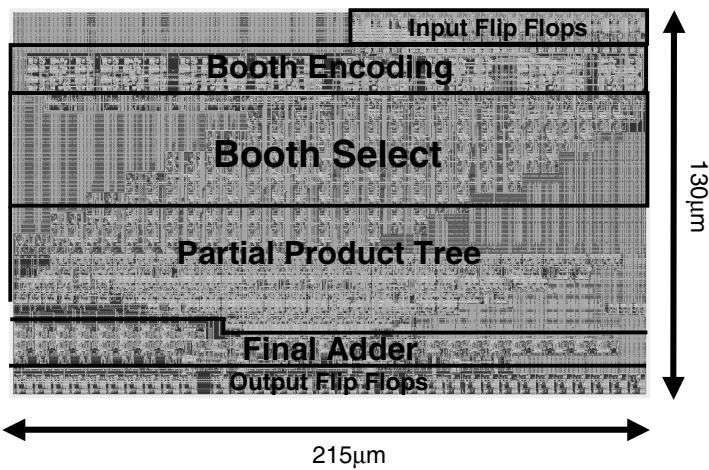
136

# Simulation Results

# Multiplier Layout

**0.8V: 500MHz, 3mW**

**1.2V: 1GHz, 22mW**

# Summary

- 16-bit multiplier features:

  - **Efficient Booth Encoding and Select**

  - **Delay and Area Optimized Partial Product Reduction Tree**

  - **Signal Profile Optimized Final Adder**

  - **Energy Optimized Sizing**

- Enables multi-mode operation

  - **1GHz at 1.2V 22mW**

  - **500MHz at 0.8V 3mW**

139