

International Conference on Supercomputing 2002 Tutorial: Minimally Clocked Microprocessor Design

**Diana Marculescu (CMU), Dave Albonesi
(Rochester), Pradip Bose (IBM)**

Schedule

- **8:30-9:30 (Diana Marculescu)**
 - Trends and issues in clock distribution
 - How much asynchrony do we want?
 - Motivation for GAL S, synchronization issues, deadlock prevention, possible inter-clocking domain communication schemes
 - GAL S processors: power/performance evaluation
- **9:30-9:45 Break**
- **9:45-10:45 (Dave Albonesi)**
 - GAL S processors: power/performance evaluation (cont'd)
 - Workload characterization and impact on the use of fine-grained speed/voltage scaling
- **10:45-11:00 Break**
- **11:00-12:00 (Pradip Bose)**
 - Case study - LPX, an IP-CMOS based processor
- **12:00-12:20 (Diana Marculescu)**
 - Looking in the crystal ball: Where will GAL S be used?
 - Concluding remarks
- **12:20-12:30 Q&A**

But first... What is “minimally clocked”?

What is minimal clocking?...

- **Asynchrony in (almost) any shape or form...**
- **... But mostly:**
 - Globally Asynchronous, Locally Synchronous (GALS) designs
 - Self-timed logic within fully synchronous systems (Locally Asynchronous, Globally Synchronous)
 - Not fully asynchronous
- **Thus, “minimally clocked” encompasses design styles that**
 - Do not have a global timing reference
 - Rely on synchronization mechanisms
 - Are not fully asynchronous

Why minimally clocked?

- **Current design roadblocks:**
 - Increasing power density
 - Increasing design complexity

- **Shortening time-to-market require:**
 - Design and verification tools and methodology
 - Solutions for design testability and reliability

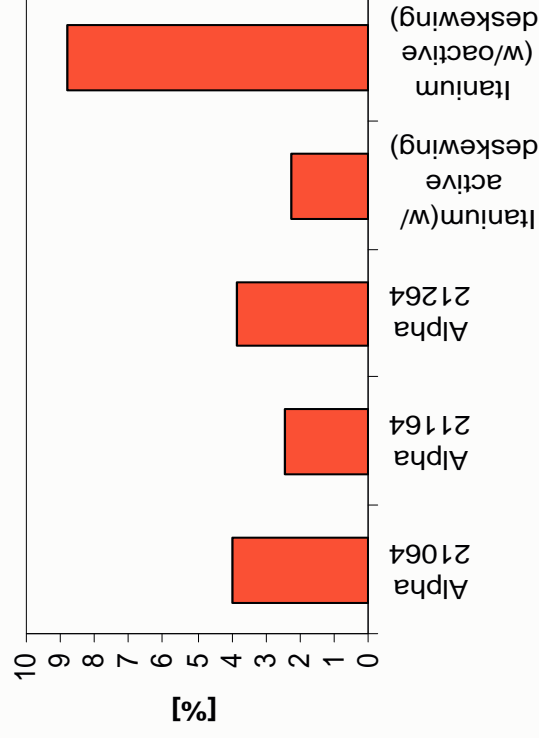
- **Drivers for minimally clocking:**
 - Design reuse (IPs independently designed)
 - Minimized design effort for global clock distribution
 - Allow fine-grain application - driven adaptability (voltage and speed scaling)

Trends and issues in clock distribution

Global clock distribution issues

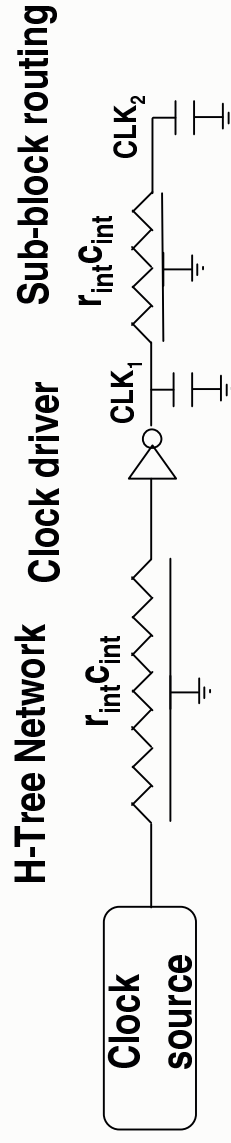
- **Global clock distribution is getting more difficult**
 - Larger die size and device count
 - Higher clock speed
 - Increasing clock power consumption (30 - 40% of total power)

- **Clock skew has become a large part of the cycle time**



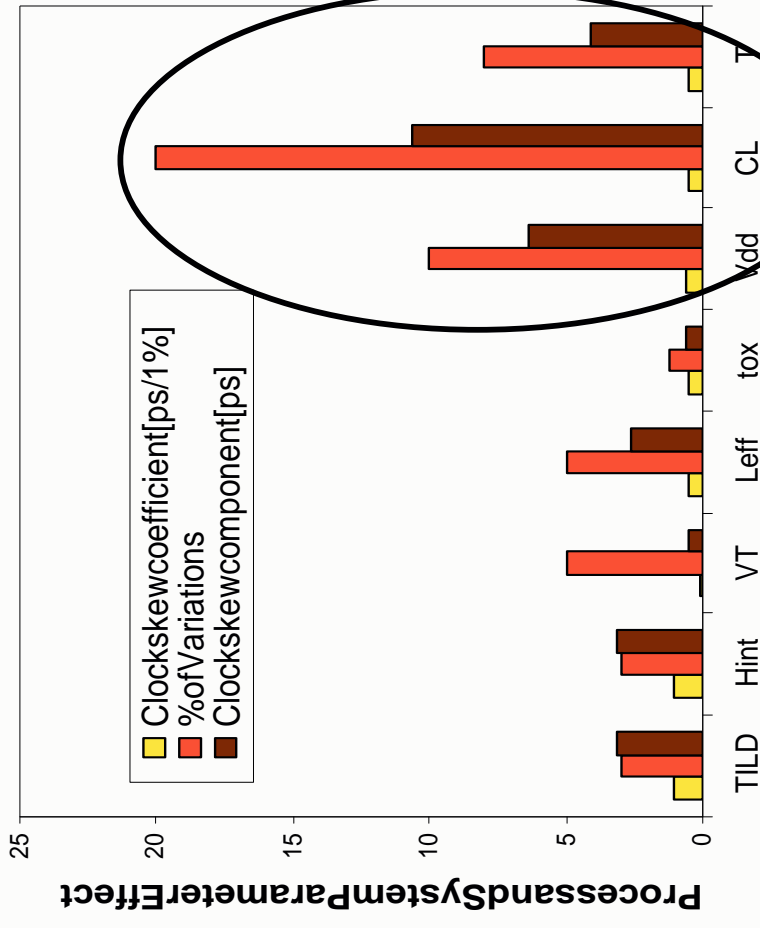
Clock skew

- **Case study: H-tree**
- **Skew components:**
 - Interconnect delay from clock source through the H-tree to the driver
 - Transistor delay of the sub-block driver
 - Internal wiring routing delay within the sub-block from clock driver to registers – also called *internal clock skew*



- **The first 2 components are most affected by physical parameter variation**

Clock skew components



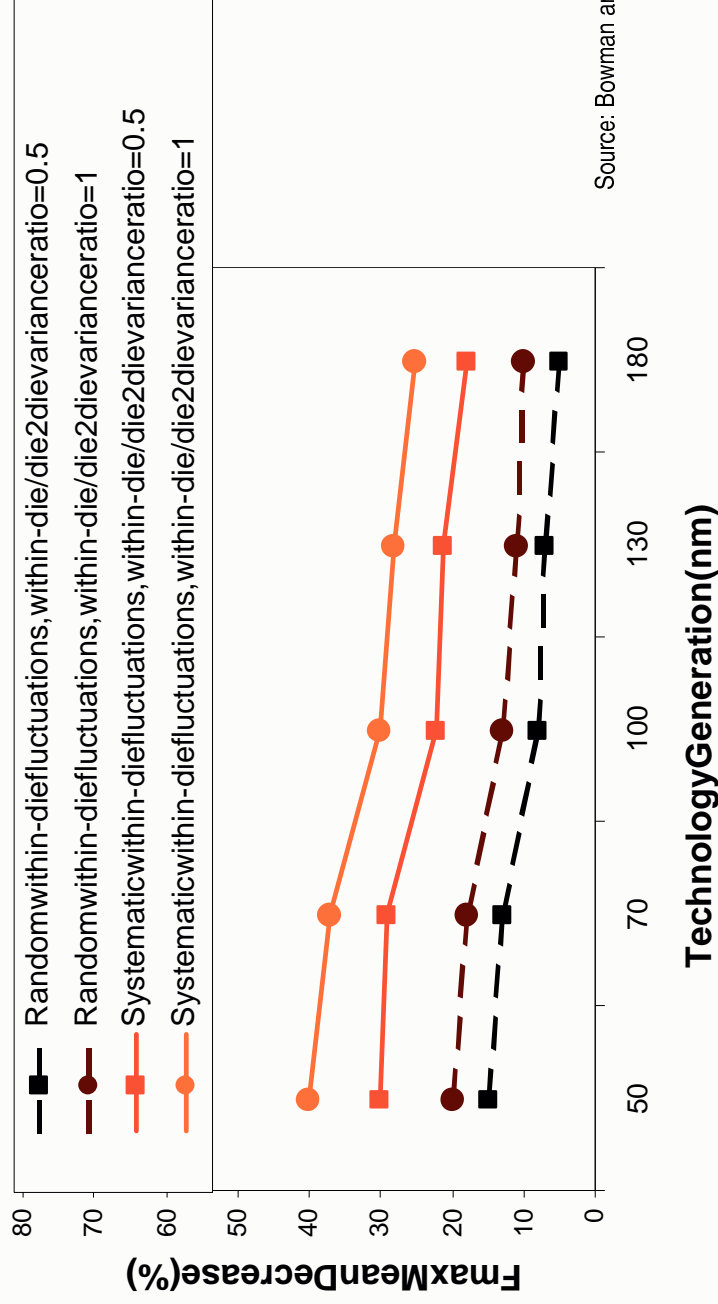
Source: Zarkesh-Ha, Mule and Meindl, 1999

- **Non-uniform distribution of clocked registers, Vdd and temperature may have the highest impact on clock skew**
 - Assumptions: 20,000 clocked registers and Alpha's temperature map

How bad can the variation be?

- **Clocked pipeline registers**
 - Worst case variations of up to a factor of 3
 - May translate into 20% variation among sub-capacitances - blockloading
- **Temperature variations**
 - Some structures are hotter than others...
 - ...And thus induce further increase in clock skew

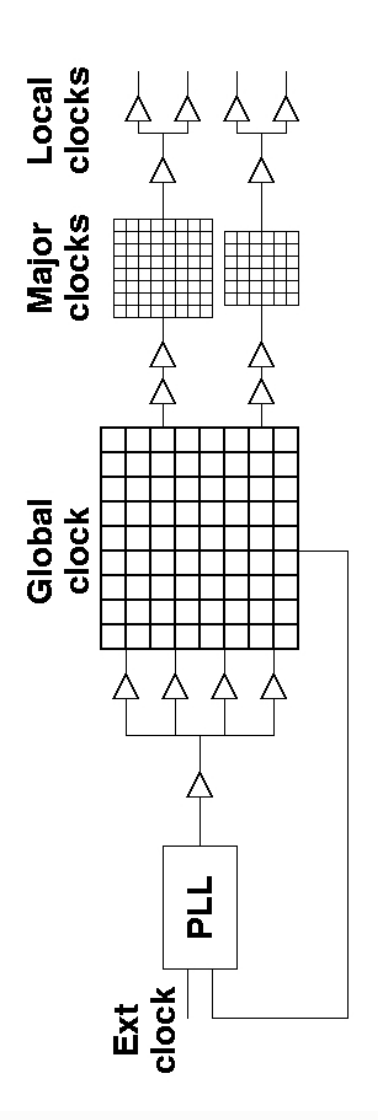
Adverse effect of process variation



- Significant effect on maximum clock speed
 - Upto 30% reduction in clock speed for 0.05um technology
- Need design styles that reduce these variations

Current solutions for clock distribution

- Traditionally done with PLL, metal grids and H-trees
 - Requires lot of die area, power, design effort
- Nowadays, clock distribution is done hierarchically
 - Global clock and major clocks
 - Deskewing circuits – may incur a lot of design effort and silicon



- ITRS 2001 (SIA Roadmap)
 - GAL S design style for alleviating clock distribution problems

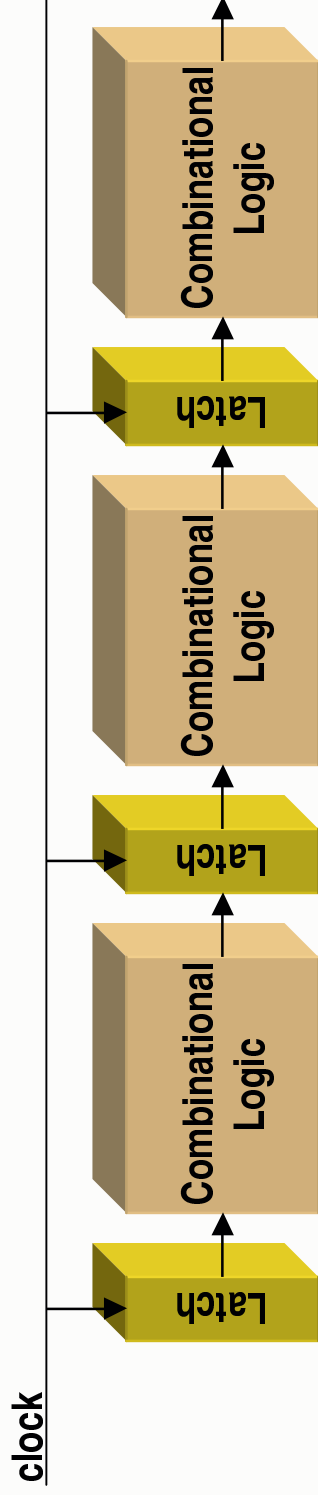
Schedule

- **8:30-9:30 (Diana Marculescu)**
 - Trends and issues in clock distribution
 - How much asynchrony do we want?
 - Motivation for GALs, synchronization issues, deadlock prevention, possible inter-clocking domain communication schemes
 - GAL processors: power/performance evaluation
- **9:30-9:45 Break**
- **9:45-10:45 (Dave Albonesi)**
 - GAL processors: power/performance evaluation (cont'd)
 - Workload characterization and impact on the use of fine-grained speed/voltage scaling
- **10:45-11:00 Break**
- **11:00-12:00 (Pradip Bose)**
 - Case study - LPX, an IPCMOS based processor
- **12:00-12:20 (Diana Marculescu)**
 - Looking in the crystal ball: Where will GALs be used?
 - Concluding remarks
- **12:20-12:30 Q&A**

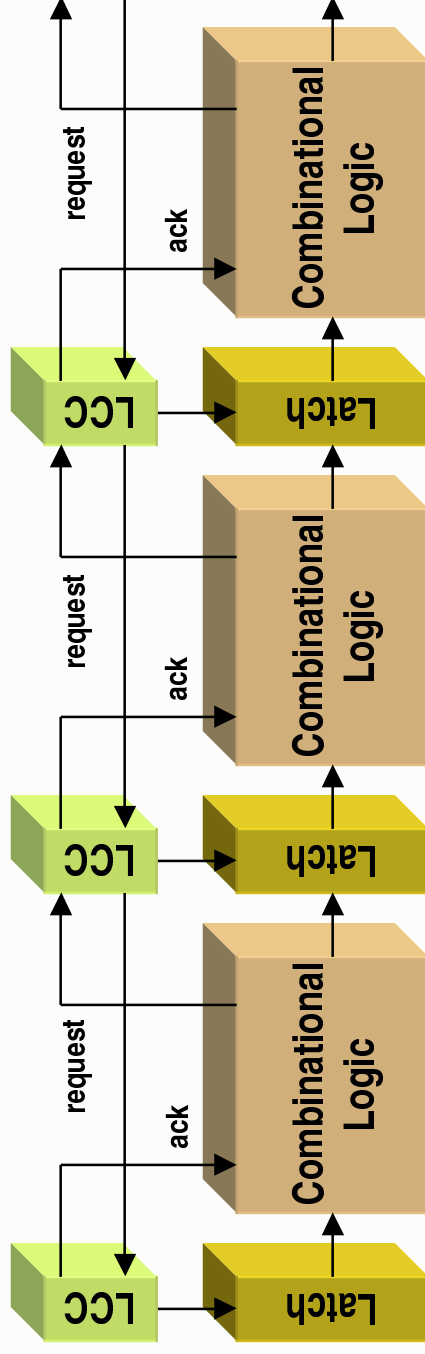
How much asynchrony do we want?

Synchronous vs. asynchronous pipelines

■ Synchronous

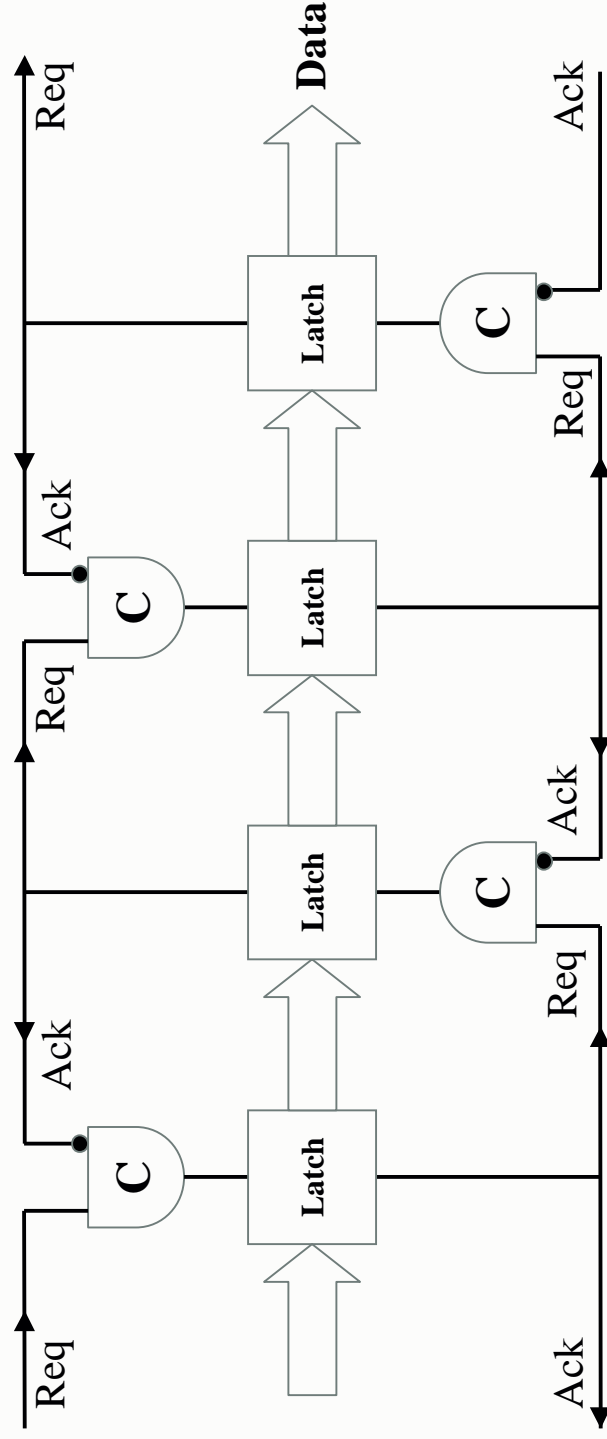


■ Asynchronous

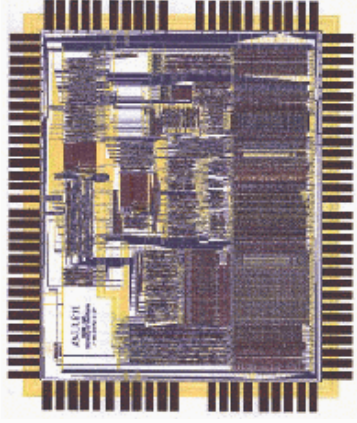


Micropipelines

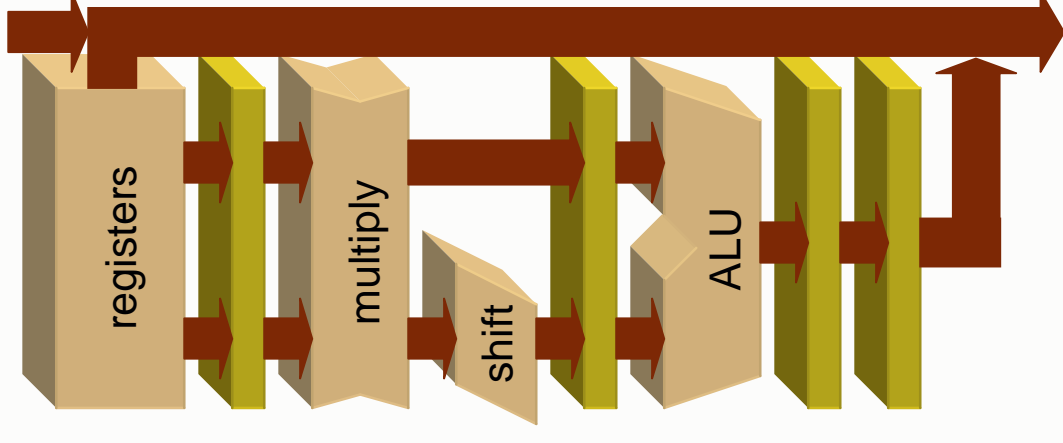
- Introduced in 1989 [Sutherland, CACM 1989, Turing Award Lecture]
- Handshaking based on transition signaling instead of standard logic levels



AMULET1

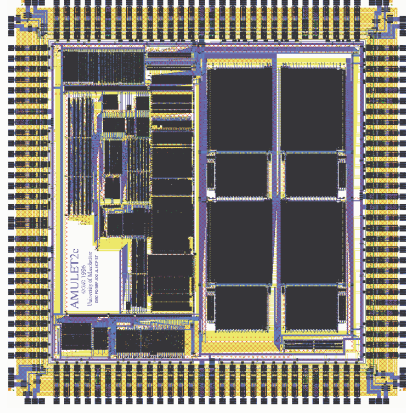


Source: <http://www.cs.man.ac.uk/amulet/index.html>

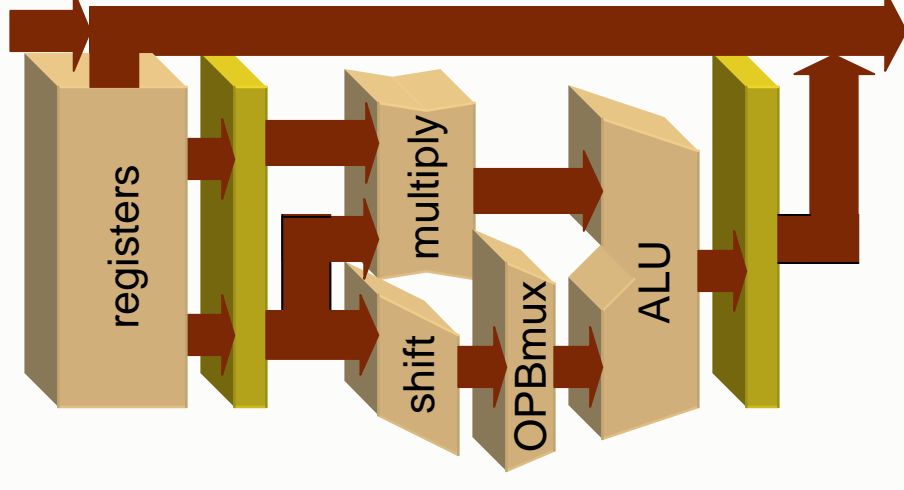


- Designed between 1990-1993
- No bypassing, relies on register locking mechanism
- Each functional unit is micropipelined
- AMULET1 Lessons
 - Micropipelines good on chip, hard to debug at board level
 - Execution pipeline deeper than needed

AMULET2e

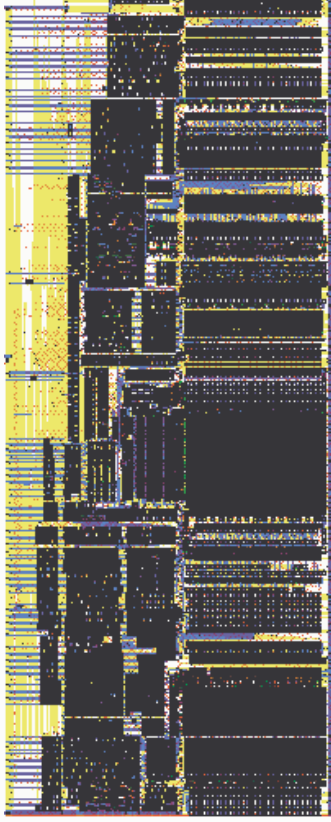


Source: <http://www.cs.man.ac.uk/amulet/index.html>



- Released 1996
- New execution pipeline
- Register forwarding
- Branch Target Cache
- 4K 64-way set associative cache

AMULET3



Source: <http://www.cs.man.ac.uk/amulet/index.html>

- Released in 1998
- Developed to prove commercial viability of asynchronous design
- Performance:
 - Maximum throughput: 120 MIPS
 - Power consumption: 155 mW @ 120 MIPS
 - 113000 transistors in processor core occupy 3sqmm on .35 micron process

- Third register port added making nearly all instructions single cycle
- Reorder buffer allows out-of-order execution
- Branch predictor improved to require fewer memory cycles
- Interrupts routed to prefetch stage to handle power-saving “Halt” operations more cleanly

Comparative performance

	AMULET1	AMULET2	AMULET3	ARM710	ARM810
Process	1 μm	0.5 μm	0.35 μm	0.6 μm	0.5 μm
Core Area	25 mm^2	25 mm^2	21 mm^2	32 mm^2	76 mm^2
Transistors	58,374	93,000	113,000	570,295	836,022
Power	152 mW	140 mW	155 mW	120 mW	500 mW
MIPS/W	77	285	780	192	172

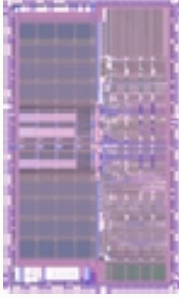
CAM – Caltech Asynchronous Microprocessor



Source: <http://async.cs.caltech.edu>

- **The world's first asynchronous microprocessor**
- **Fabricated in 1989**
- **16-bit RISC machine with 16 general-purpose registers**
- **Peak performance in HP 1.6 μ m CMOS**
 - 5MIPS at 2V drawing 5.2mA of current
 - 18MIPS at 5V drawing 45mA
 - 26MIPS at 10V drawing 105mA

miniMIPS



Source: <http://async.cs.caltech.edu>

- **Asynchronous MIPS R3000**
- **Designed between 1995 and 1998**
- **Peak performance in HP 0.6 μm CMOS**
 - 180MIPS and 4W at 3.3V and 250
 - 100MIPS and 850mW at 2.0V and 250
 - 60MIPS and 220mW at 1.5V and 250
- **Expected 280 MHz, got about 180 MHz due to:**
 - HP process
 - Layout errors

Comparative performance

	MiniMIPS @3.3V	MiniMIPS @2V	AMULET2e	StrongARM SA110
Process	0.6 μm	0.6 μm	0.5 μm	0.35 μm
Power	7 W	1 W	150 mW	120 mW
MIPS ² /E $10^{21} \text{ J}^{-1}\text{S}^{-2}$	3100	3375	365	14420

- Better than Amulet2e in terms of performance per unit of energy
- Difficult to compare – different process technology, ISAs, etc.

So... Why not fully asynchronous?

- **Fully asynchronous benefits**
 - No clock, very low power
 - Average case performance may be better

- **However...**
 - Complete migration to asynchronous is not feasible immediately
 - Lack of mature CAD tools for fully asynchronous design...
 - ...and inertia

- **GALS provides an intermediate solution**
 - Assumes locally clocked, fully synchronous regions
 - Globally asynchronous communication among regions

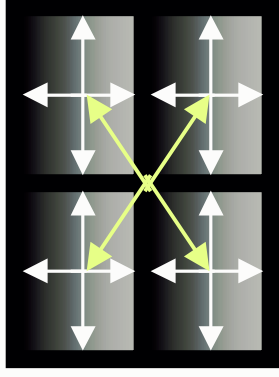
Schedule

- **8:30-9:30 (Diana Marculescu)**
 - Trends and issues in clock distribution
 - How much asynchrony do we want?
 - Motivation for GALs, synchronization issues, deadlock prevention, possible inter-clocking domain communication schemes
 - GALs processors: power/performance evaluation
- **9:30-9:45 Break**
- **9:45-10:45 (Dave Albonesi)**
 - GALs processors: power/performance evaluation (cont'd)
 - Workload characterization and impact on the use of fine-grained speed/voltage scaling
- **10:45-11:00 Break**
- **11:00-12:00 (Pradip Bose)**
 - Case study - LPX, an IPCMOS based processor
- **12:00-12:20 (Diana Marculescu)**
 - Looking in the crystal ball: Where will GALs be used?
 - Concluding remarks
- **12:20-12:30 Q&A**

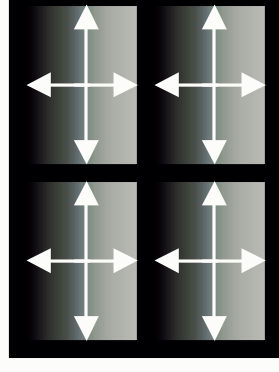
GALS systems: Why and how ?

GALS design benefits

Synchronous



GALS



- **Lack of a global clock**
 - Less design effort and cost for dealing with clock skew
 - Less power consumption

- **Potential for fine-grain adaptability**
 - Different speeds among synchronous blocks
 - Different voltages, and hence potential for lower power consumption
 - Can be used with on-the-fly application-driven adaptation

GALS design issues

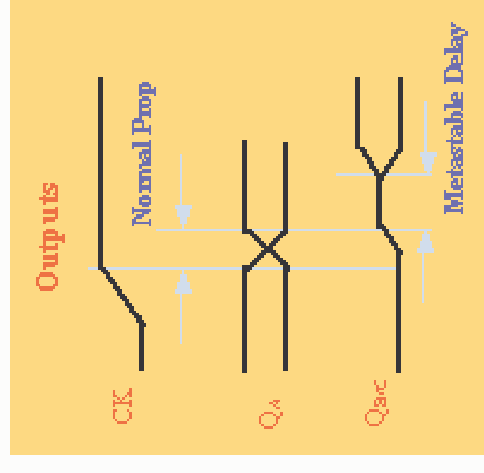
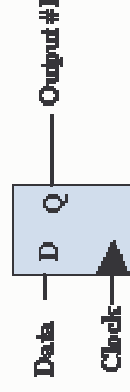
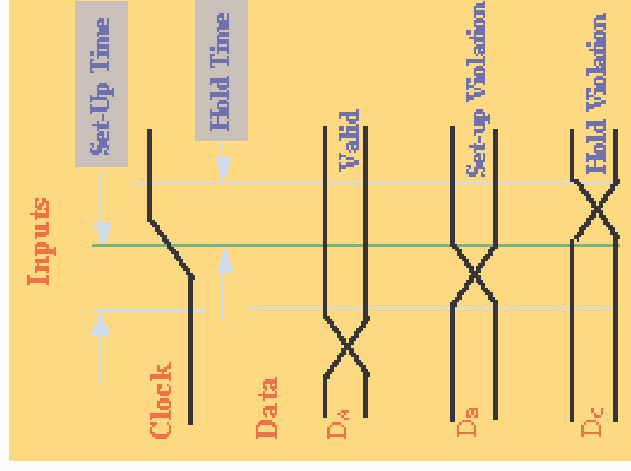
- **Metastability resolution**
 - Always a problem when interfacing clock domains
 - Synchronizers and arbiters
 - Possible solution: mixed -timing FIFOs [Chelcea and Nowick, DAC2001, WCED2002]

- **Local clock generation**
 - Ring oscillators [Muttersbach et al., ASYNC2000]

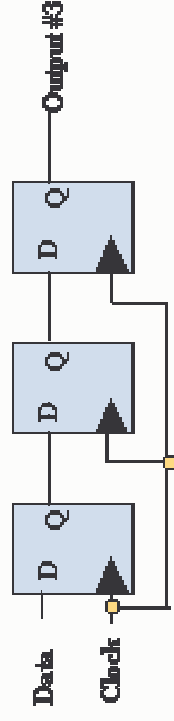
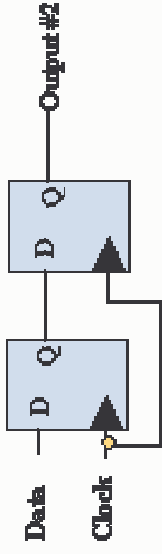
- **Failure modeling**
 - Synchronization failures can happen
 - The goal is to maximize Mean Time Between Failure (MTBF)

Metastability

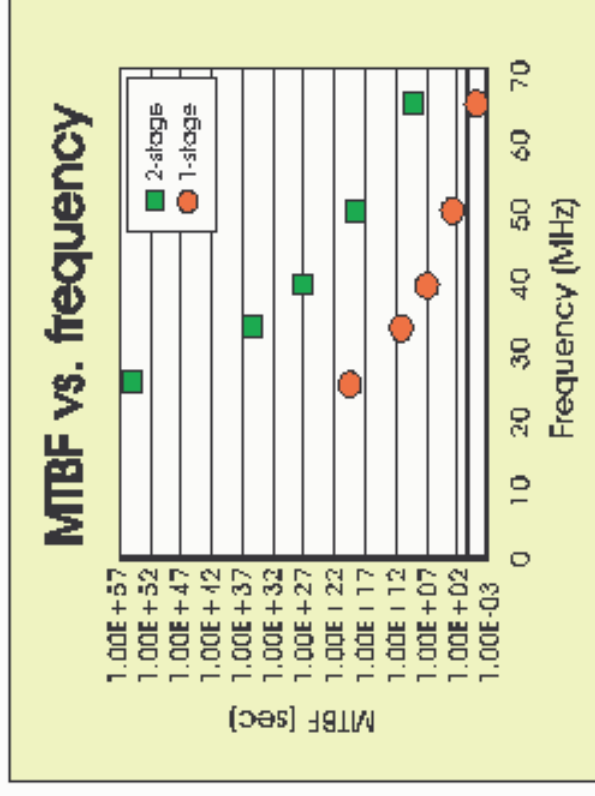
- Asynchronous signals may violate
 - Setuptime
 - Holdtime
 - Both
- Hence the possibility of a metastable state



Using synchronizers



- Adding a second FF will reduce the chance of the output going metastable.
- The output from the first flip flop may go valid, before the second flip flop is clocked.
- Adding yet another FF will reduce the probability that its output will be unstable even more.

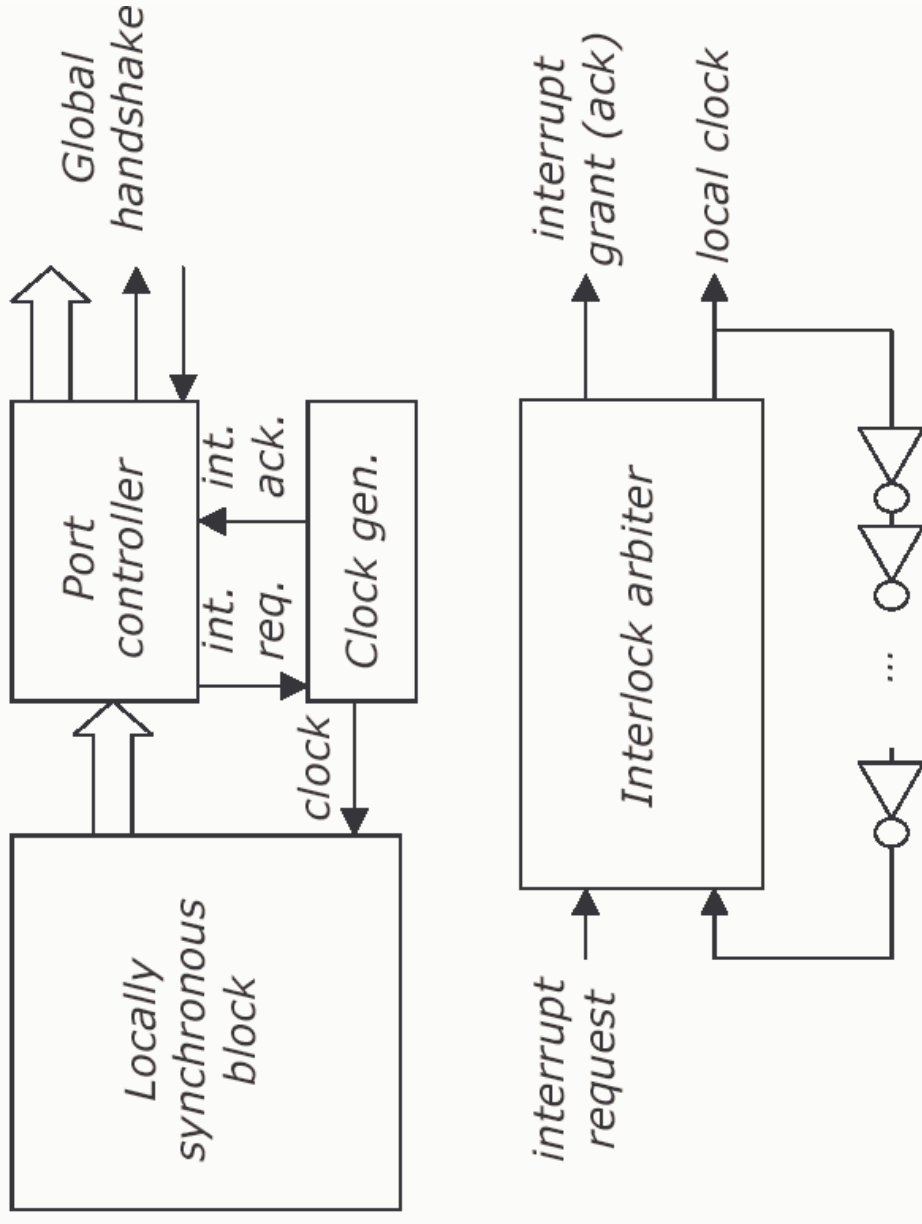


Source: Texas Instruments

Pausable clocks

- **Some applications cannot live with a significant failure rate:**
 - Multiprocessing
 - Cryptography and security applications
 - Mission critical applications
- **Idea for accommodating metastability [Yun et al., ICCD 1996]:**
 - Pause the passive phase of a local clock, long enough to allow metastability to resolve. Then the clock can be generated by an interlock arbiter.

Pausable clocks

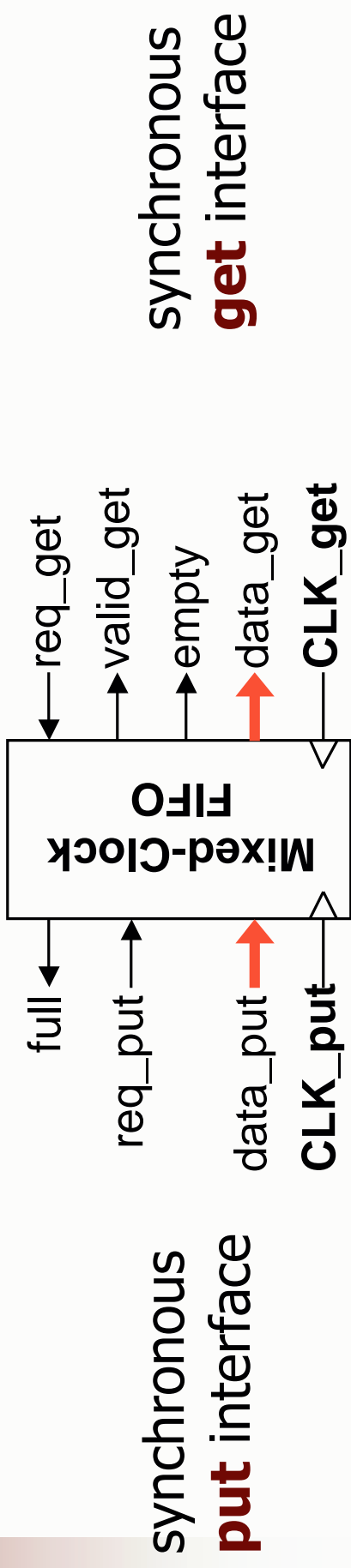


Source: Yun et al., 1996

Mixed timing FIFOs

- **Mixed timing or asynchronous FIFOs [Chelcea and Nowick, DAC 2001, WCED 2002]**
 - Low-latency
- **Modular and scalable:**
 - Define interfaces for each combination of Synchronous or Asynchronous domains
 - Combine interfaces to design new async/sync FIFO's
- **High throughput:**
 - In steady state: no synchronization overhead, no failure probability
 - Enqueue/Dequeue data items: one/cycle
- **Low area overheads**
- **Also, solve issue of long interconnect delays between domains**

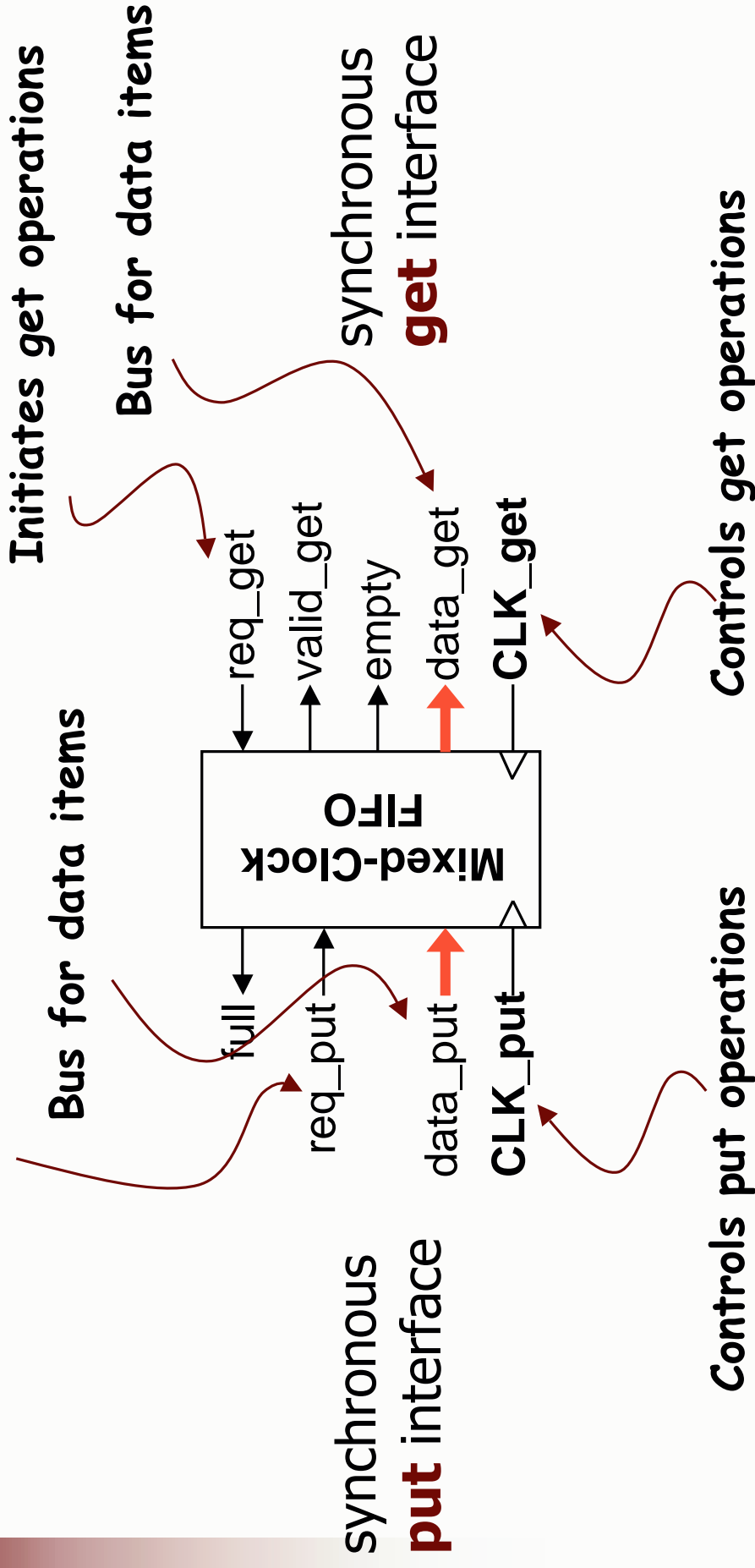
Mixed-timing FIFOs



Source: Chelcea and Nowick, 2001, 2002

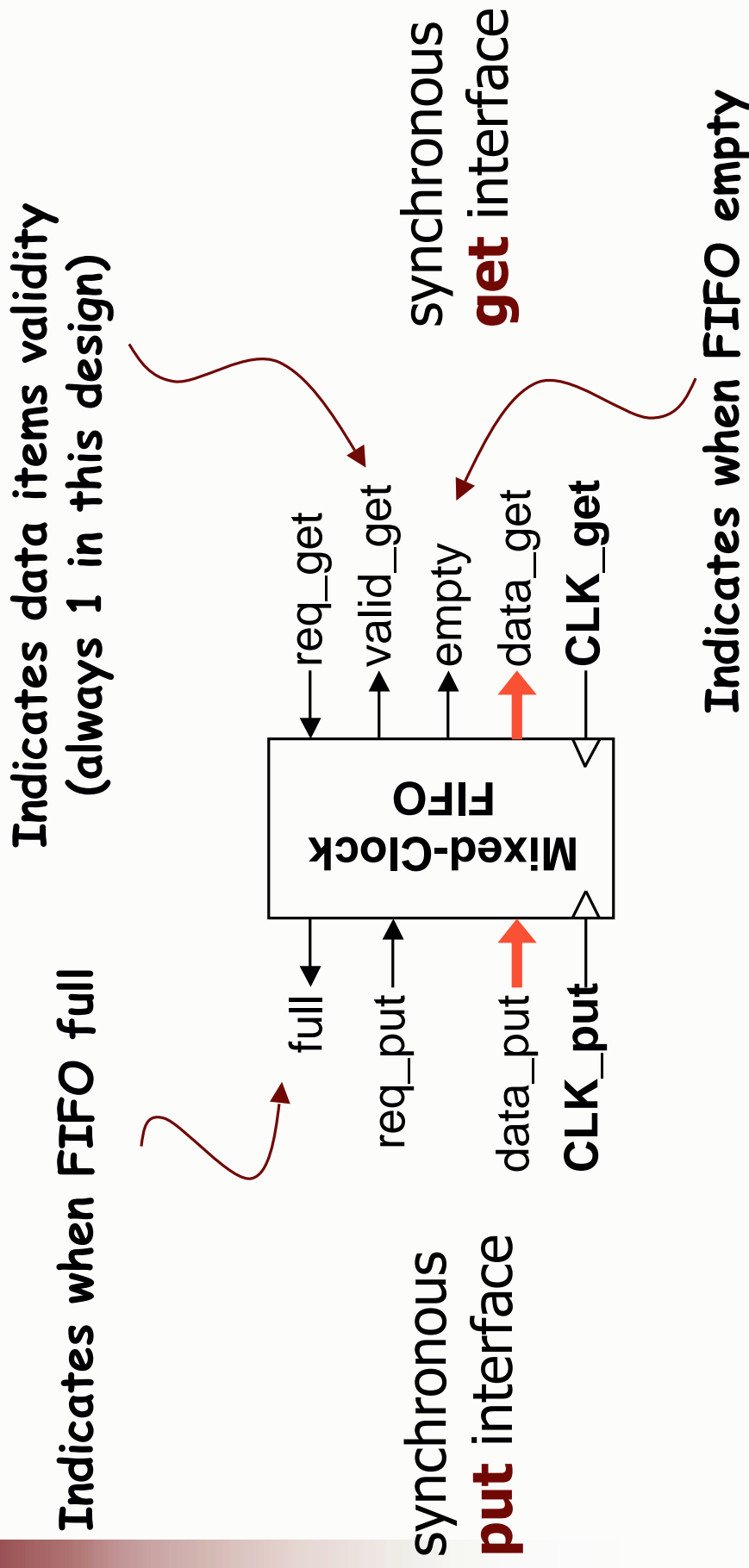
Mixed-timing FIFOs

Initiates put operations



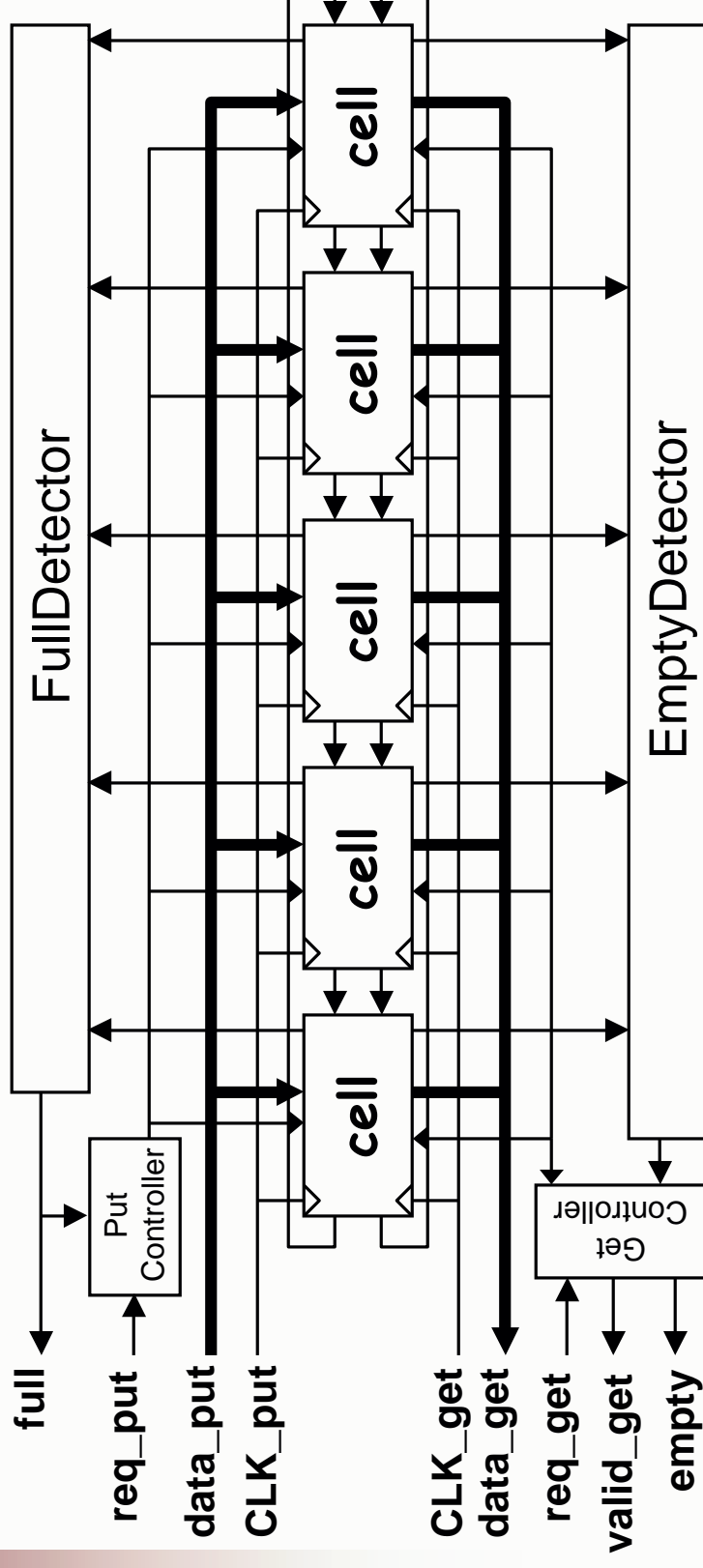
Source: Chelcea and Nowick, 2001, 2002

Mixed-timing FIFOs



Source: Chelcea and Nowick, 2001, 2002

Mixed-timing FIFOs



Source: Chelcea and Nowick, 2001, 2002

Local clock generation

- **Ring oscillators**
 - Pro: Simple
 - Con: Cannot be controlled externally
 - ... But they can be coupled with the “local control” of speed and/or voltage

- **Local PLLs**
 - Pro: Externally controlled
 - Con: More costly
 - ... but they “scale” [IWLS2002, Focus Group on GALs]

Schedule

- **8:30-9:30 (Diana Marculescu)**
 - Trends and issues in clock distribution
 - How much asynchrony do we want?
 - Motivation for GALs, synchronization issues, deadlock prevention, possible inter-clocking domain communication schemes
 - GALs processors: power/performance evaluation
- **9:30-9:45 Break**
- **9:45-10:45 (Dave Albonesi)**
 - GALs processors: power/performance evaluation (cont'd)
 - Workload characterization and impact on the use of fine-grained speed/voltage scaling
- **10:45-11:00 Break**
- **11:00-12:00 (Pradip Bose)**
 - Case study - LPX, an IP-CMOS based processor
- **12:00-12:20 (Diana Marculescu)**
 - Looking in the crystal ball: Where will GALs be used?
 - Concluding remarks
- **12:20-12:30 Q&A**

Power-performance evaluation of GALs processors

Architecture definition

- **Automatic partitioning into clock domains [Hemani et al., DAC 1999]**
 - Applied only to random logic, not helpful for high-end processors

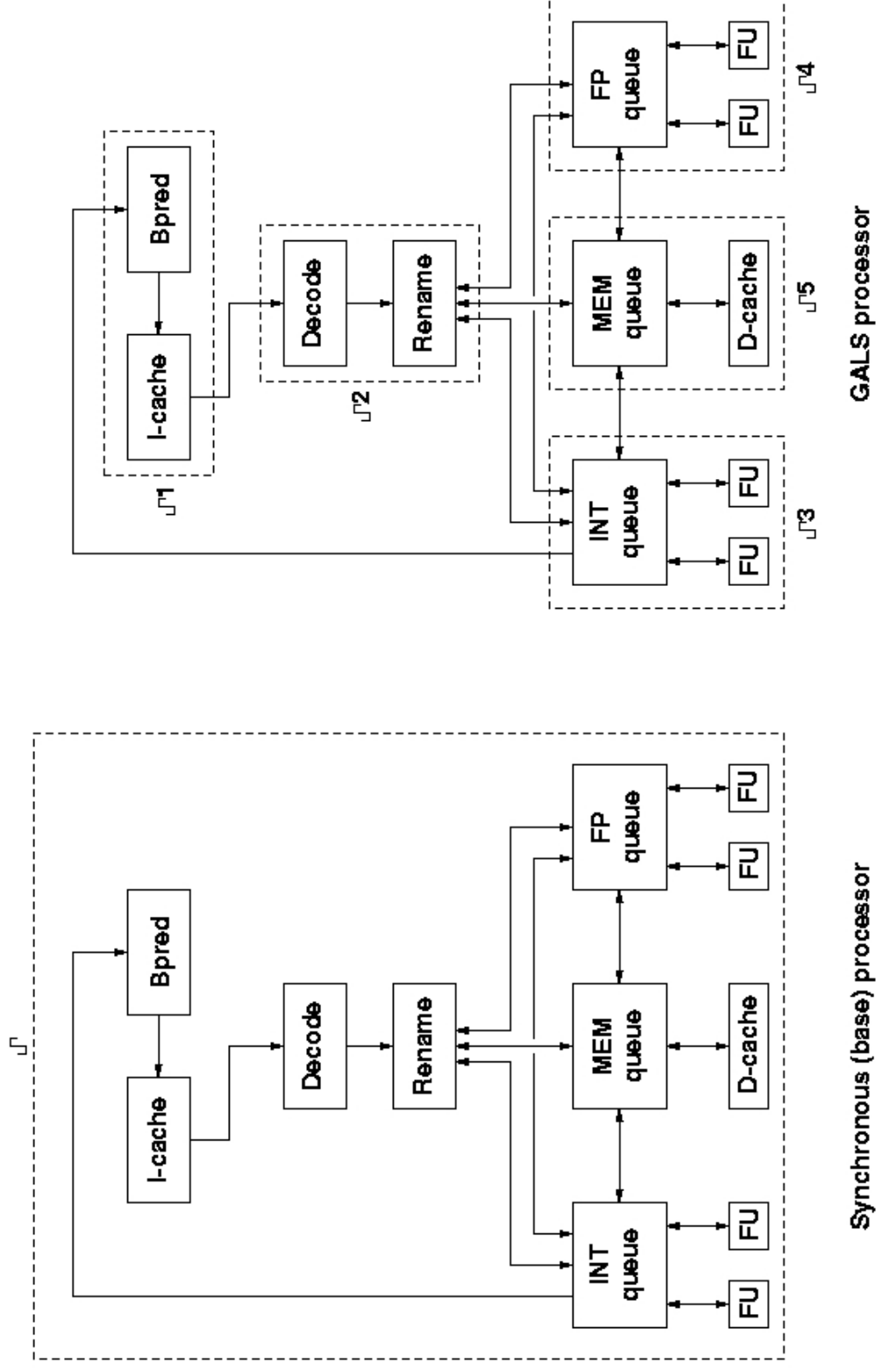
- **A typical superscalar core exhibits “natural decoupling”:**
 - Front-end (I-cache and BPhardware)
 - Decode, Register renaming
 - Integer, FP and memory “domains”

- **Local clock grids usually follow the same type of partitioning**

Architecture definition

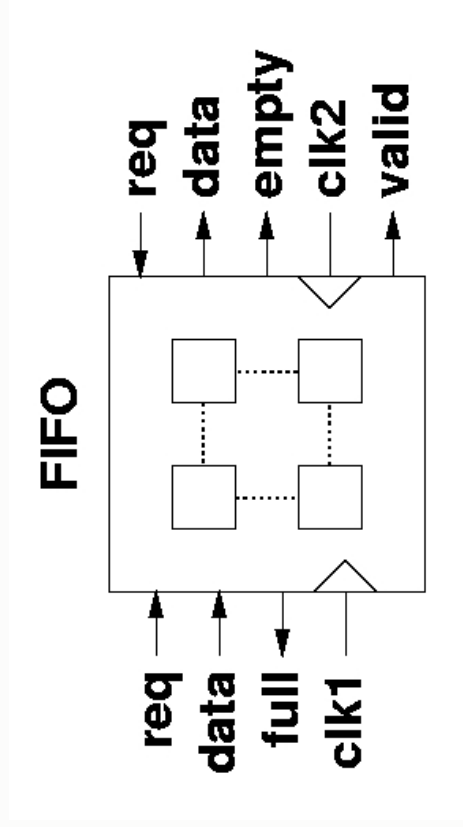
- **Starting point:**
 - A hypothetical asynchronous machine resembling the Alpha 21264
- **GALS architecture: satisfies physical, as well as architectural constraints**
- **Five clock domains**
 - 1. L1-cache, B-pred (FETCH)
 - 2. Decode, rename (DECODE)
 - 3. Integer issue queue, integer units (INT)
 - 4. Memory issue queue, d-cache (MEM)
 - 5. Floating point issue, floating point units (FP)
- **Clock domains are interfaced using asynchronous FIFOs**

Architecture definition



Mixed timing communication

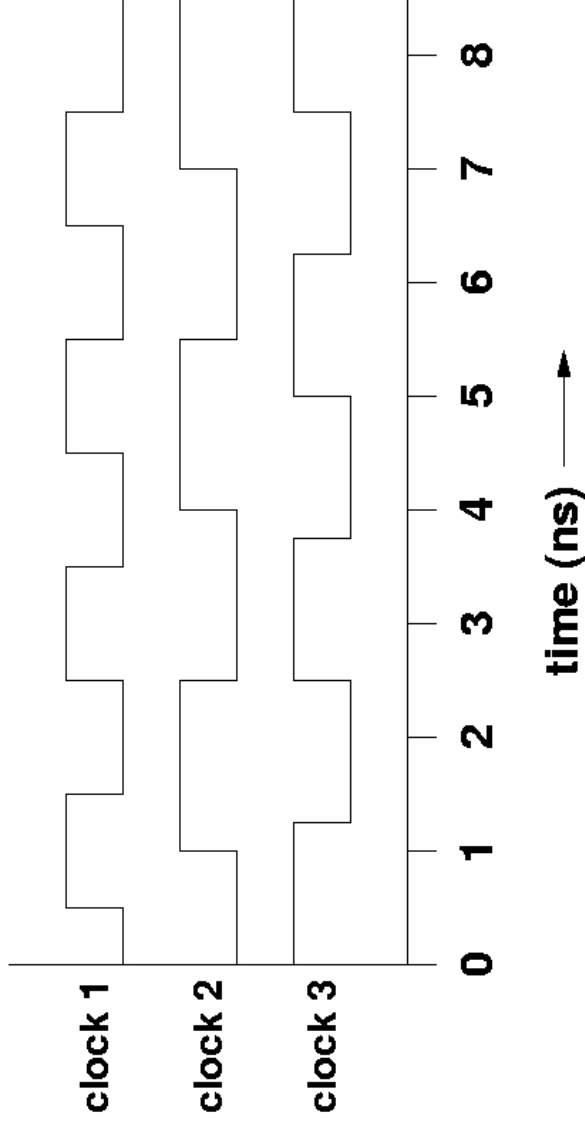
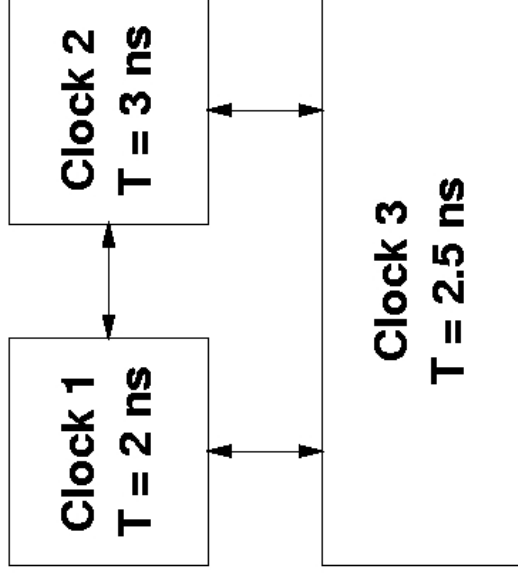
- **Used the token-ring based FIFO design**
 - [Chelcea and Nowick, DAC2001, WCED2002]
 - Interface two clocks by synchronizing control signals
 - Provides full throughput and low latency in steady state
 - Included cycle-accurate model in our simulation environment



Simulation framework

- **Event-driven simulation engine: sim-GALS**
 - Each clock domain is represented by one set of periodic events

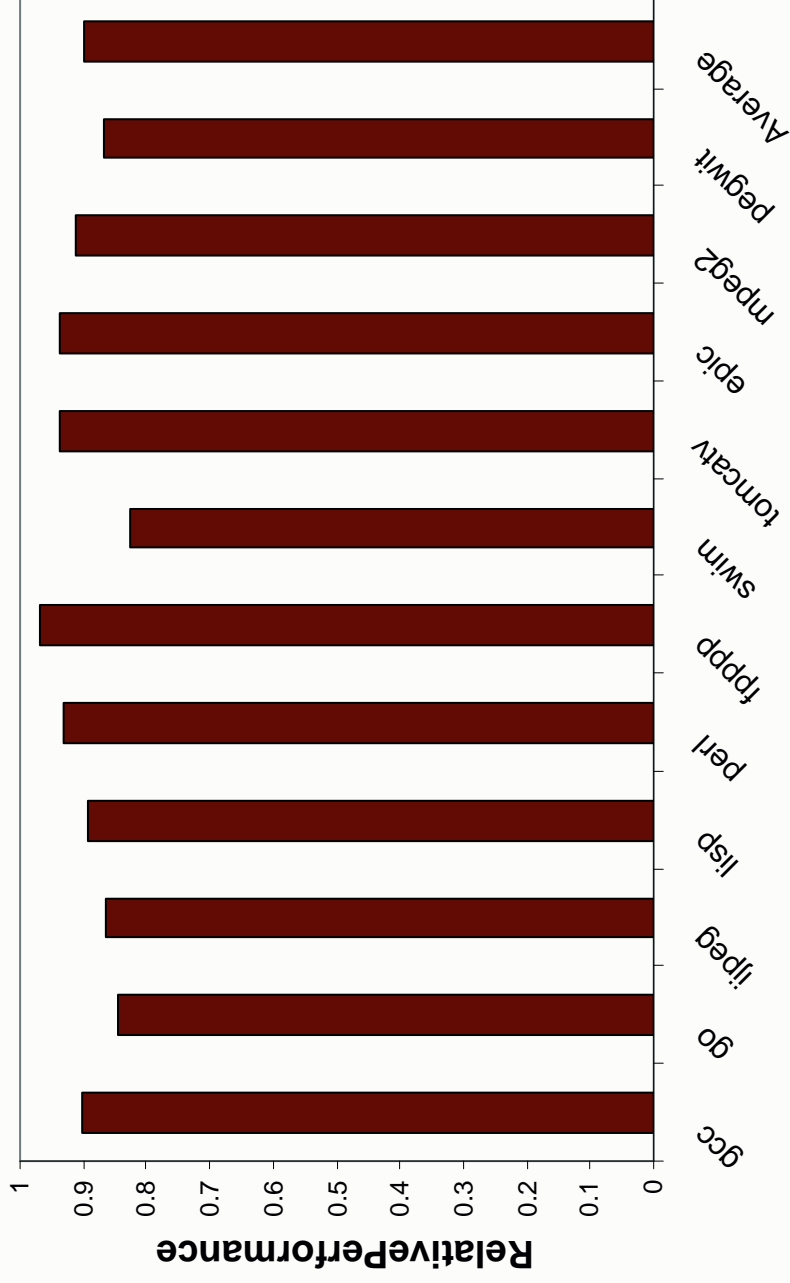
- **Reused some of the primitives in SimpleScalar and Wattch**



Microarchitecture settings

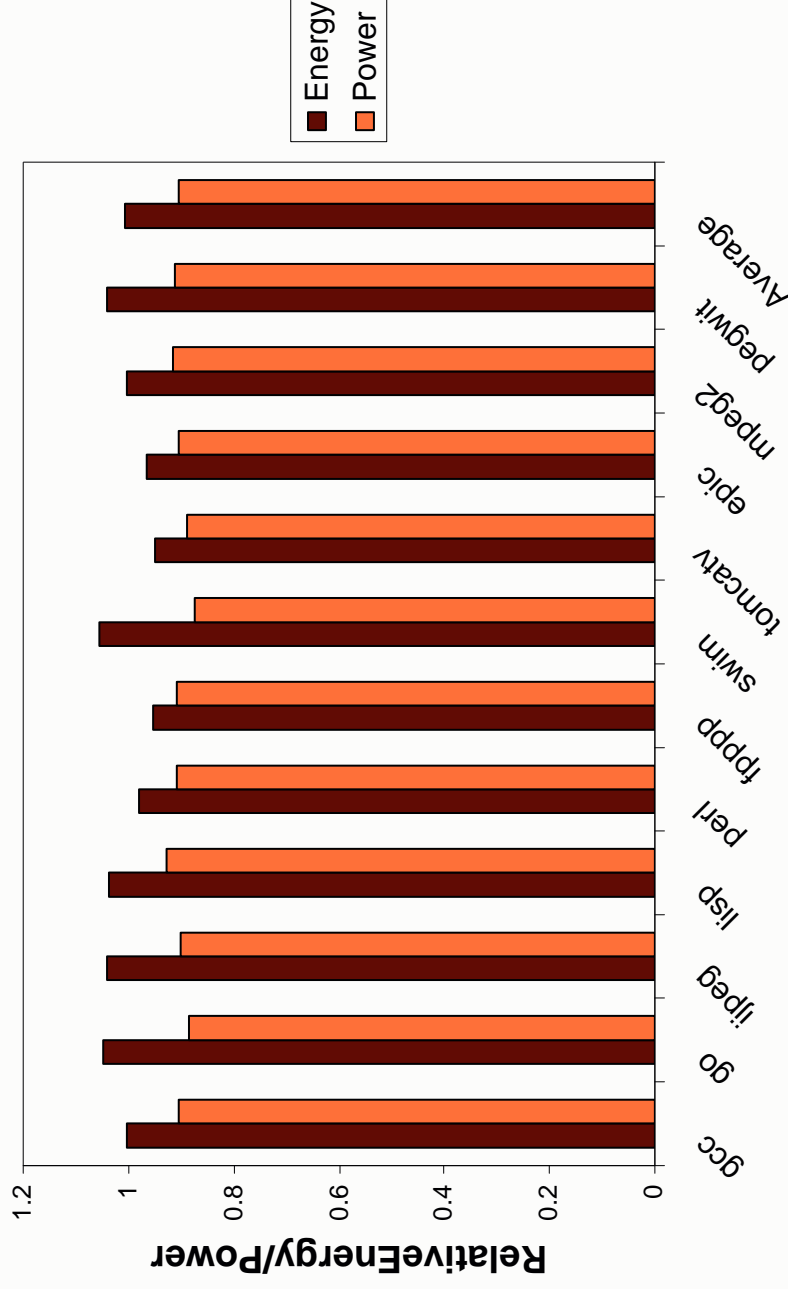
- **Fetch, decode, commit rate** 4 inst/cycle
- **Integer issue queue size** 20
- **FP issue queue size** 16
- **Memory issue queue size** 16
- **Integer registers** 72
- **FP registers** 72
- **L1 data cache** 16KB 4-way
- **L1 instruction cache** 1 cycle latency
16KB direct-mapped
- **L2 unified cache** 1 cycle latency
256KB 4-way
- **ALUs** 6 cycles latency
4 Int, 4 FPs
- **Clock gating** 1-2 cycles for restarting the clock

Performance



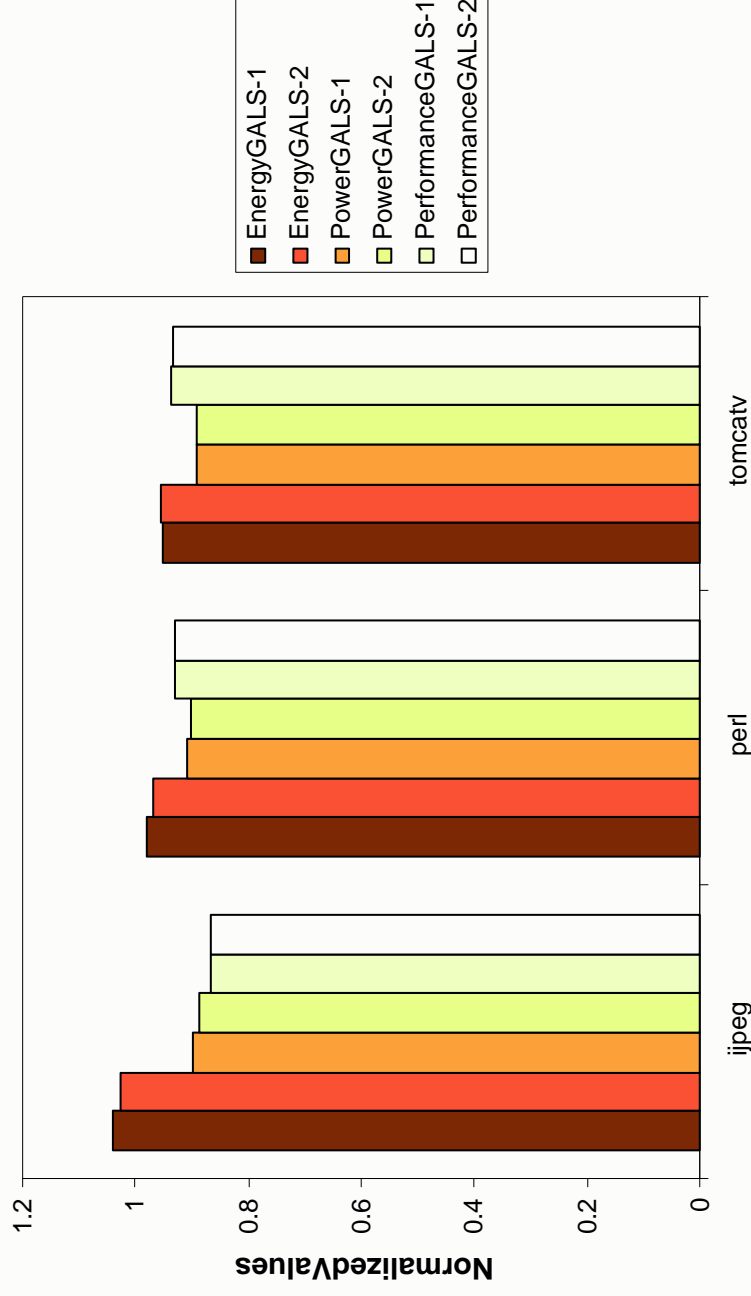
■ 10% average performance penalty (3-18%)

Energy and power



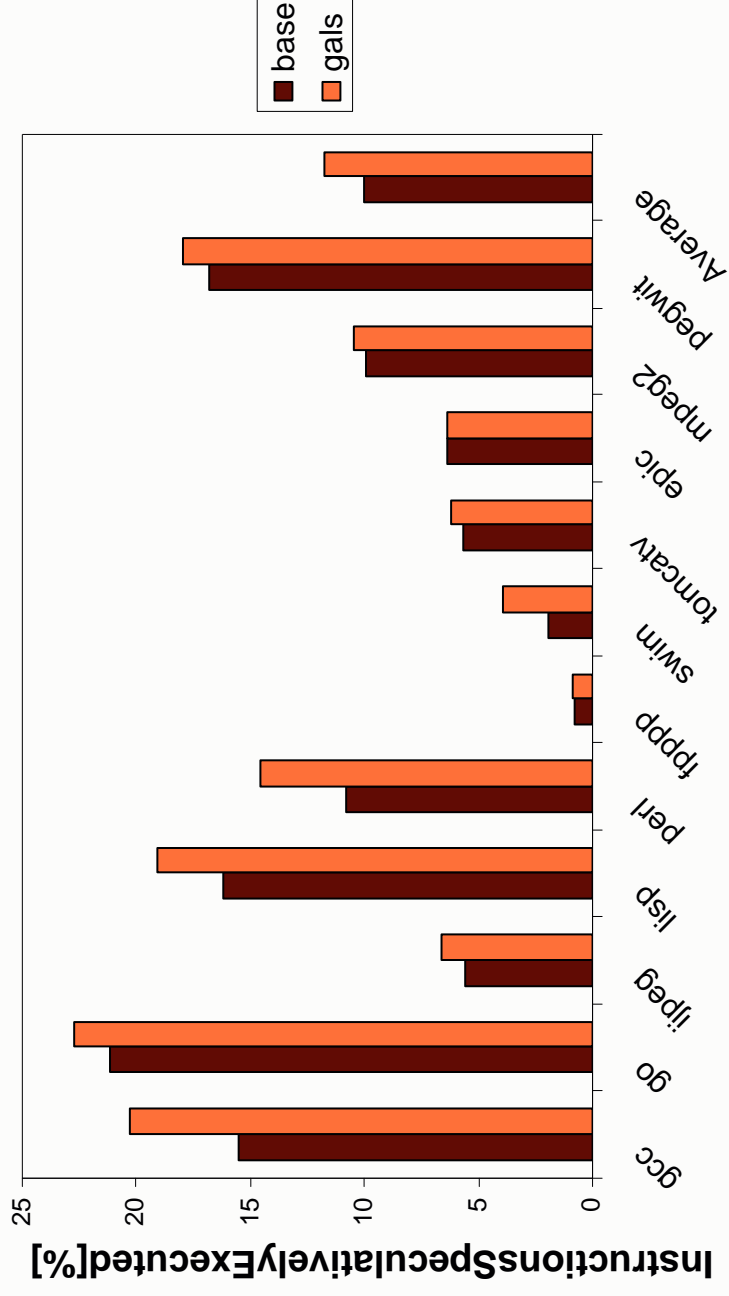
- Elimination of global clock does not necessarily reduce total energy consumption
- Average power consumption is reduced (10% on average)

Impact of clock stopping overhead



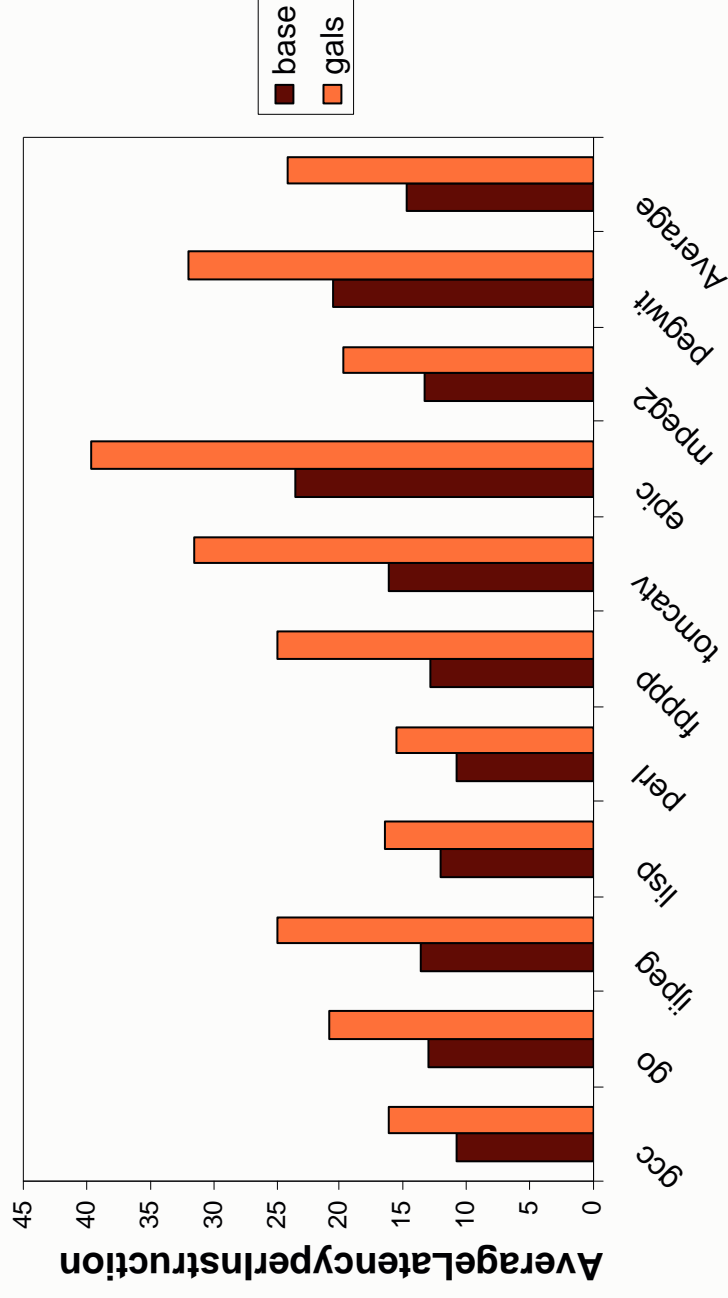
- Shutdown a clock domain after 4 idle cycles
- Restart it when decode stage detects an instruction which “needs” that domain
- Minimal impact of clock stopping/restarting overhead

Speculative execution



- On average, 2% more instructions are speculatively executed
- In the worst case, 5% more instructions are speculatively executed

Average latency per instruction

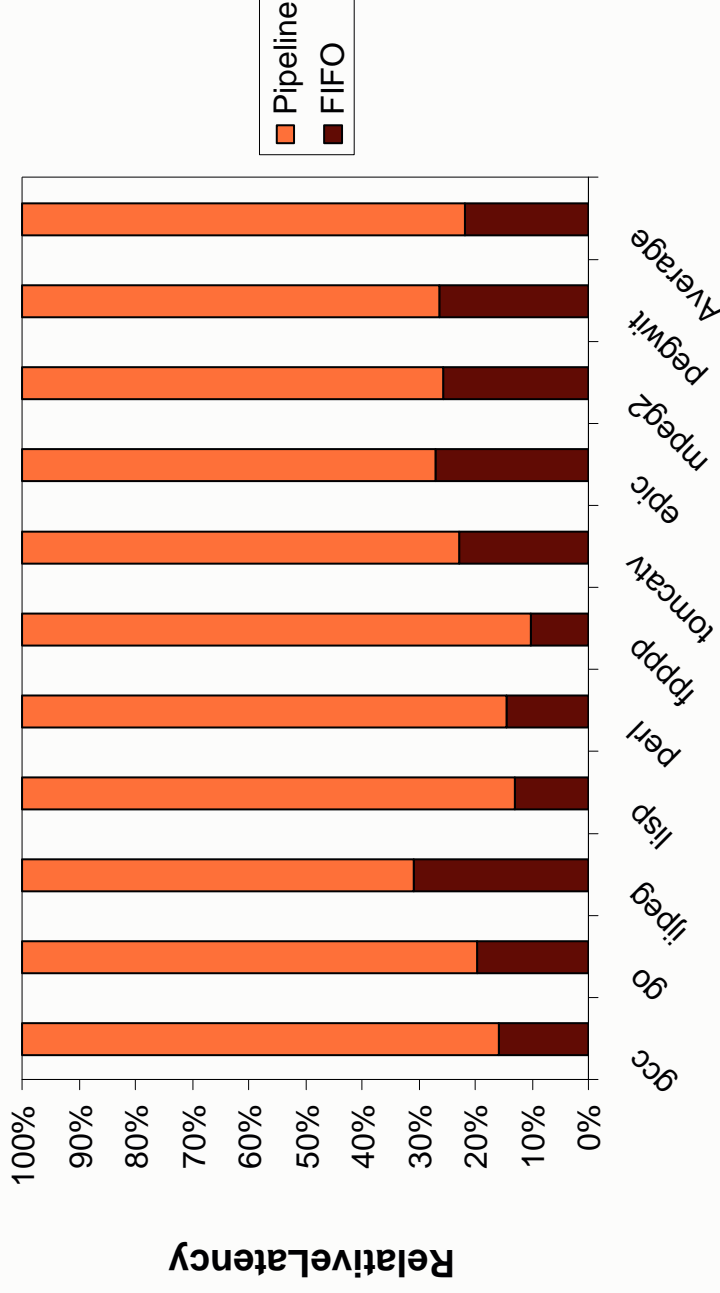


■ End-to-end latency increases from 15 to 24 cycles, on average

Energy and performance trends

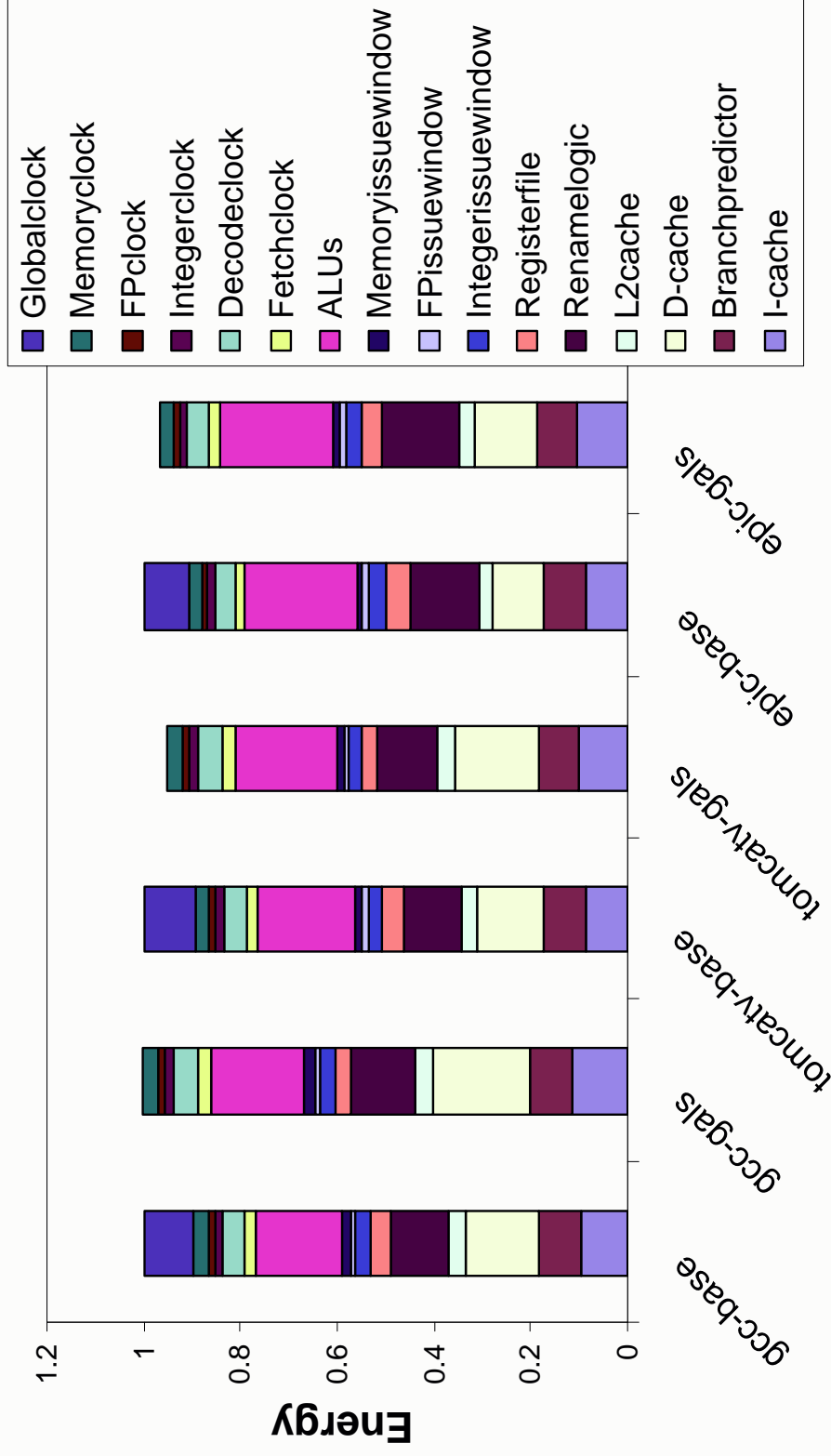
- **Average power decreases, but overall energy may increase**
- **Reasons:**
 - Performance drops due to asynchronous communication
 - Hence longer execution time and more overhead for unused modules
 - Longer branch recovery pipeline
 - Consequently, more speculative execution
 - Higher fetch to commit time for each instruction
 - Higher occupancies of rename tables and issue queues

Where does performance go?



- Asynchronous FIFOs introduce 22% latency overhead on average
- The rest of latency is spent in other stages of the pipeline

Where does energy go?



■ Longer execution times translate into larger energy costs

Delay - voltage dependency

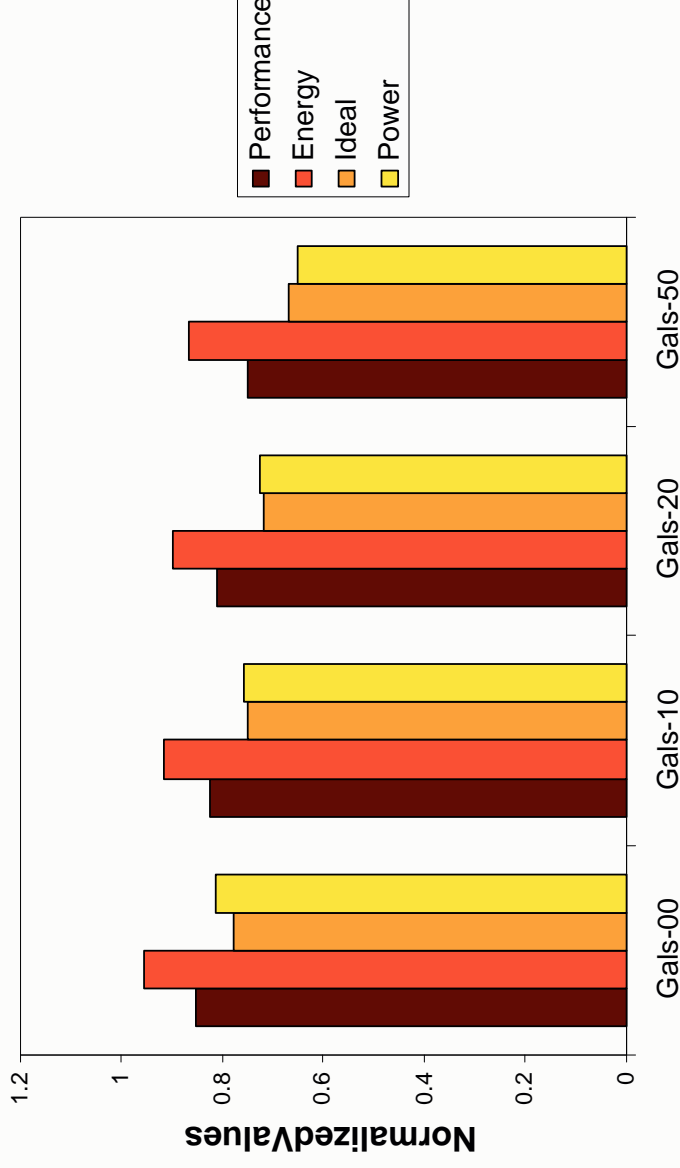
■ Fine-grain voltage reduction can be beneficial in GALS systems

- Each clock domain can be run at a different speed

$$D \propto \frac{V_{dd}}{(V_{dd} - V_t)^\alpha}$$

- V_{dd} is the supply voltage
- V_t is the threshold voltage
- α is a technology - defined constant between 1 and 2;
- α is 1.2 - 1.6 for present generation [Chen et al, JSSC 1998]

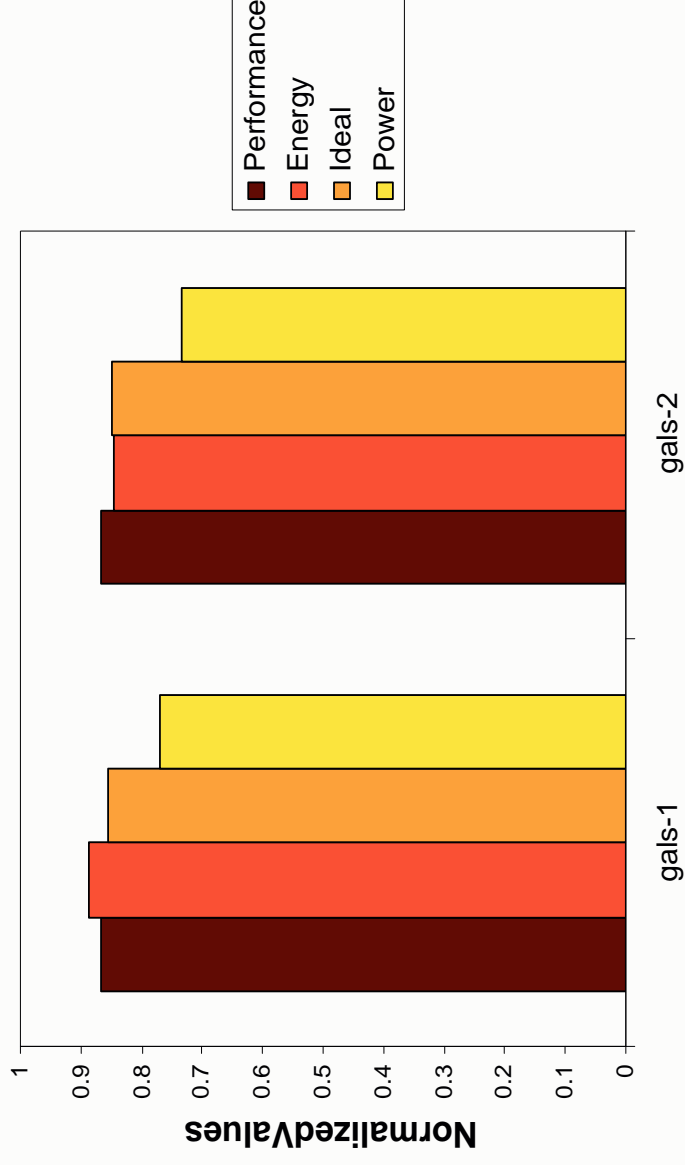
Selective voltage reduction for jpeg



■ GALS version gives some energy savings

- Fetchclocksloweddown10%
- FPclocksloweddown20%
- Memoryclocksloweddown0 -10-20-50%
- Ideal:base,fullysynchronousversionwithreducedvoltage

Selective voltage reduction for gcc



- 1: FP clock: 1.5; 2: FP clock: 3.0
- GALS version is as good as the fully synchronous baseline

Related reading list

- D.M. Chapiro, Globally Asynchronous Locally Synchronous Systems. PhD thesis, Stanford University, 1984.
- Ivan E. Sutherland. Micropipelines. (Turing award lecture) Communications of the ACM, 32(6):720-738, June 1989.
- Robert F. Sproull, Ivan E. Sutherland, and Charles E. Molnar. The counterflow pipeline processor architecture. IEEE Design & Test of Computers, 11(3):48-59, Fall 1994.
- K.Y. Yun, R.P. Donohue. Pausible Clocking: A First Step Toward Heterogeneous Systems. Computer Design: In Proc. Intl. Conference on VLSI in Computers and Processors (ICCD), 1996, pages: 118 - 123.
- A. Hemani, T. Meincke, S. Kumar, A. Postula, T. Olsson, P. Nilsson, J. Oberg, P. Ellervee, and D. Lundqvist. Lower Power Consumption in Clock By Using Globally Asynchronous Locally Synchronous Design Style. In Proc. Design Automation Conference (DAC), 1999.
- J. Muttersbach, T. Villiger, and W. Fitchner. Practical Design of Globally Asynchronous Locally Synchronous Systems. In Proc. Intl. Symposium on Advances in Research in Asynchronous Circuits and Systems (ASYNC), 2000.
- I.E. Sutherland, S. Fairbanks. GasP - A Minimal FIFO Control. In Proc. Intl. Symposium on Advanced Research in Asynchronous Circuits and Systems (ASYNC), 2001.
- T. Chelcea and S.M. Nowick. Robust Interfaces for Mixed - Timing Systems with Application to Latency-Insensitive Protocols. In Proc. Design Automation Conference (DAC), 2001.
- A. Iyer and D. Marculescu. Power and Performance Evaluation of Globally Asynchronous Locally Synchronous Processors. In Proc. Intl. Symposium on Computer Architecture (ISCA), 2002, pages: 158 - 168.
- <http://www.cs.man.ac.uk/amulet/index.html>
- <http://www.async.caltech.edu/>

Schedule

- **8:30-9:30 (Diana Marculescu)**
 - Trends and issues in clock distribution
 - How much asynchrony do we want?
 - Motivation for GALs, synchronization issues, deadlock prevention, possible inter-clocking domain communication schemes
 - GAL processors: power/performance evaluation
- **9:30-9:45 Break**
- **9:45-10:45 (Dave Albonesi)**
 - GAL processors: power/performance evaluation (cont'd)
 - Workload characterization and impact on the use of fine-grained scaling
- **10:45-11:00 Break**
- **11:00-12:00 (Pradip Bose)**
 - Case study - LPX, an IP-CMOS based processor
- **12:00-12:20 (Diana Marculescu)**
 - Looking in the crystal ball: Where will GALs be used?
 - Concluding remarks
- **12:20-12:30 Q&A**