



---

# VENTURES INTO POWER-AWARE COMPUTER ARCHITECTURE DESIGN

**R. IRIS BAHAR**  
*SCHOOL OF ENGINEERING*  
*BROWN UNIVERSITY*



**BROWN**  
School of Engineering



# THE PRE-DAWN OF POWER-AWARE PROCESSOR DESIGN

# ARCHITECTURES FOR INSTRUCTION-LEVEL PARALLELISM

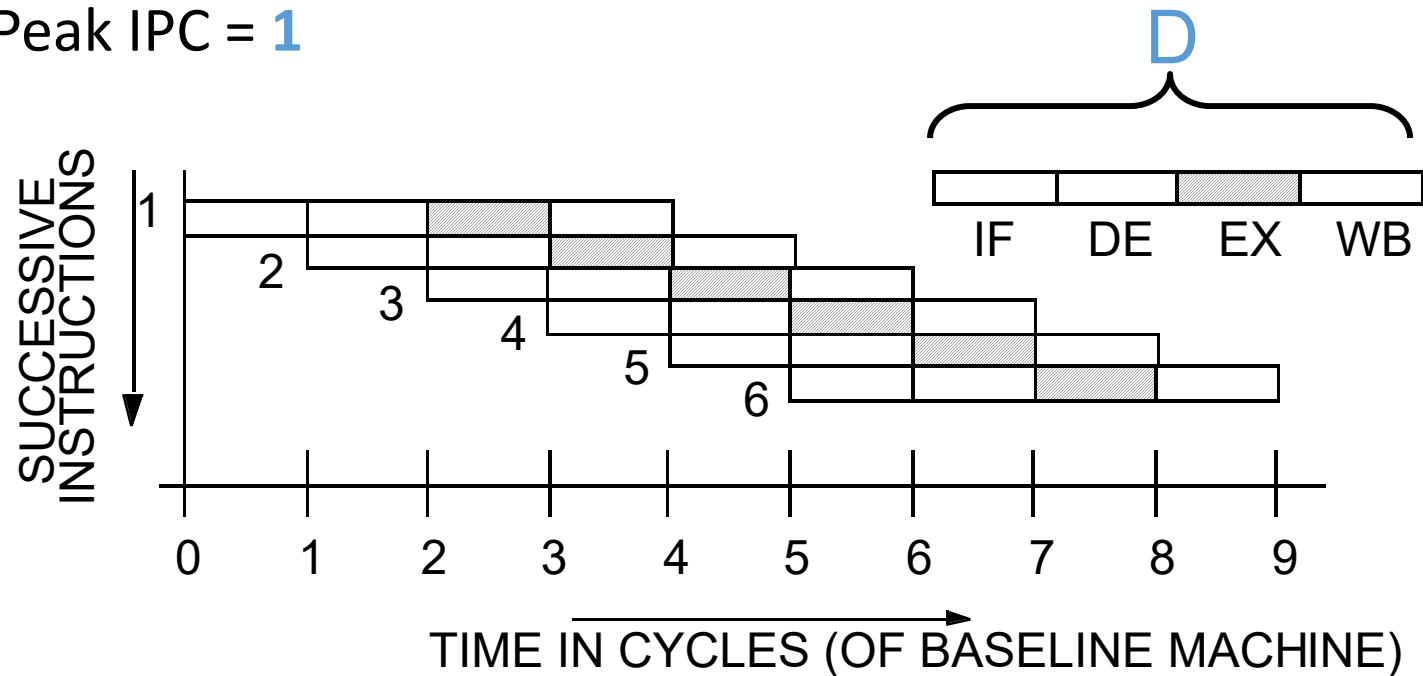


## Scalar Pipeline (baseline)

Pipeline Depth = **D**

Operation Latency = **1**

Peak IPC = **1**



# SUPERPIPELINED MACHINE

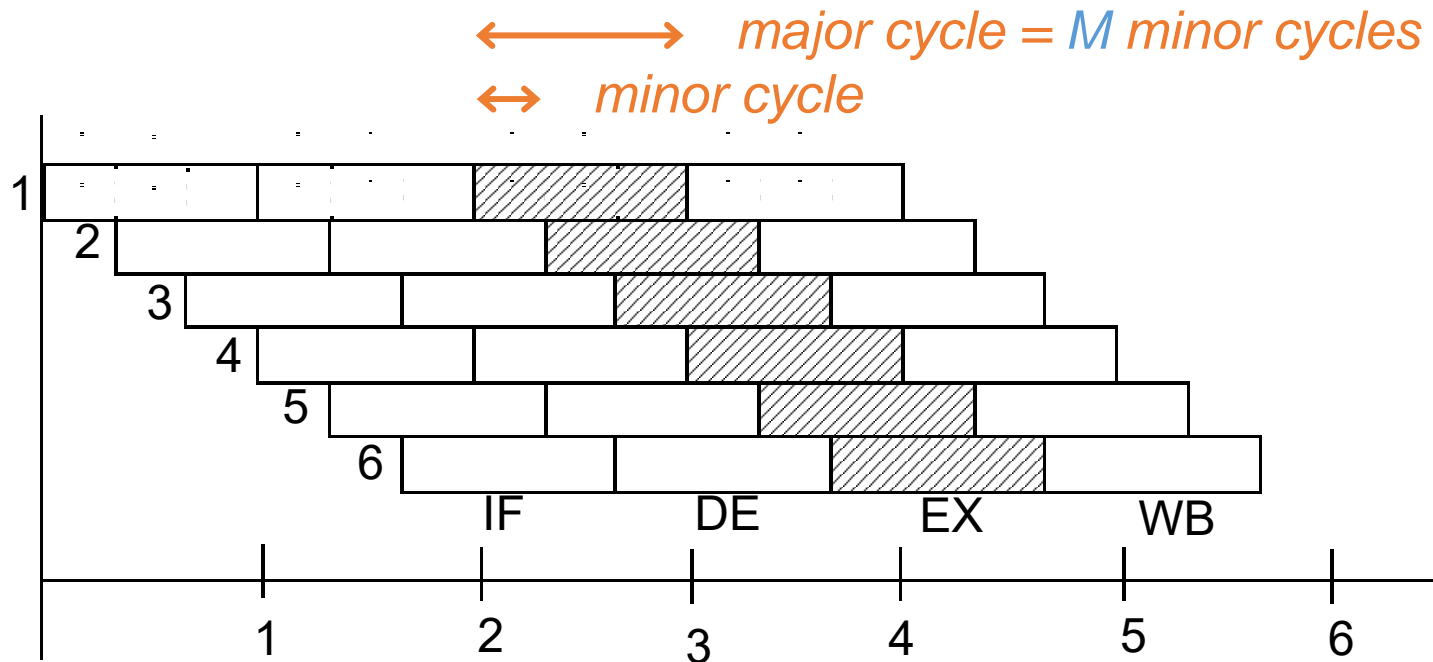


## Superpipelined Execution

Instruction Parallelism =  $D \times M$

Operation Latency =  $M$  minor cycles

Peak IPC =  $1$  per minor cycle ( $M$  per baseline cycle)



# SUPERSCALAR MACHINES

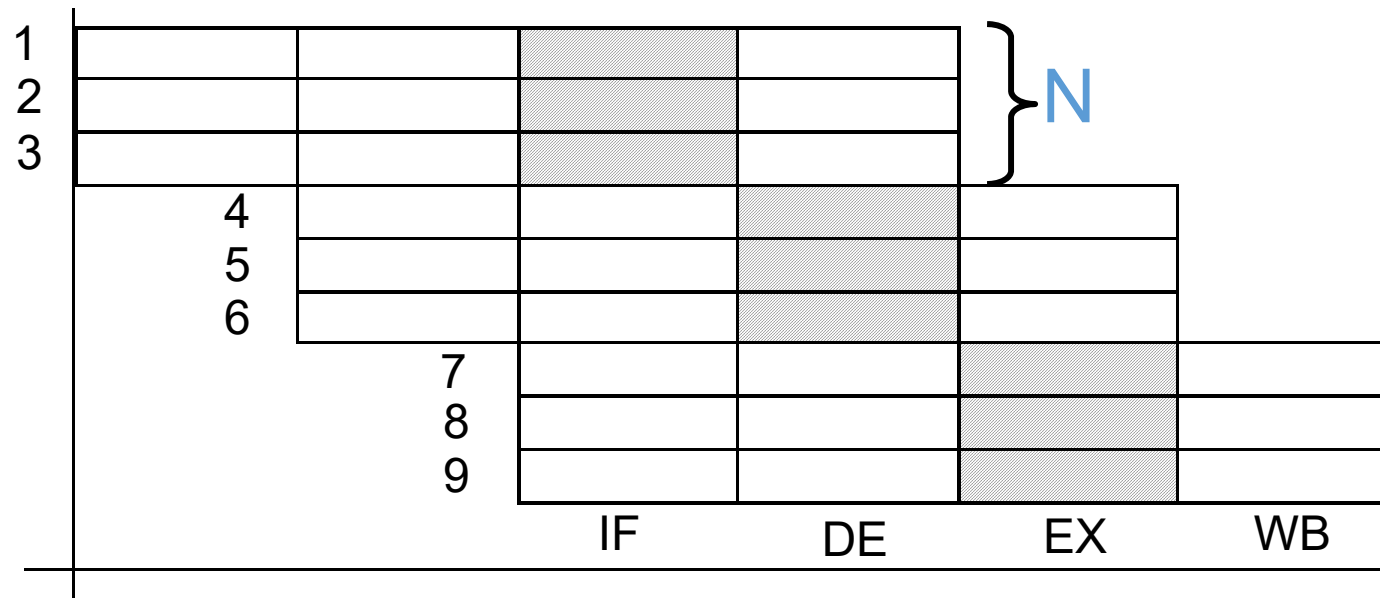


## Superscalar (Pipelined) Execution

Instruction Parallelism =  $D \times N$

Operation Latency =  $1$  baseline cycles

Peak IPC =  $N$  per baseline cycle



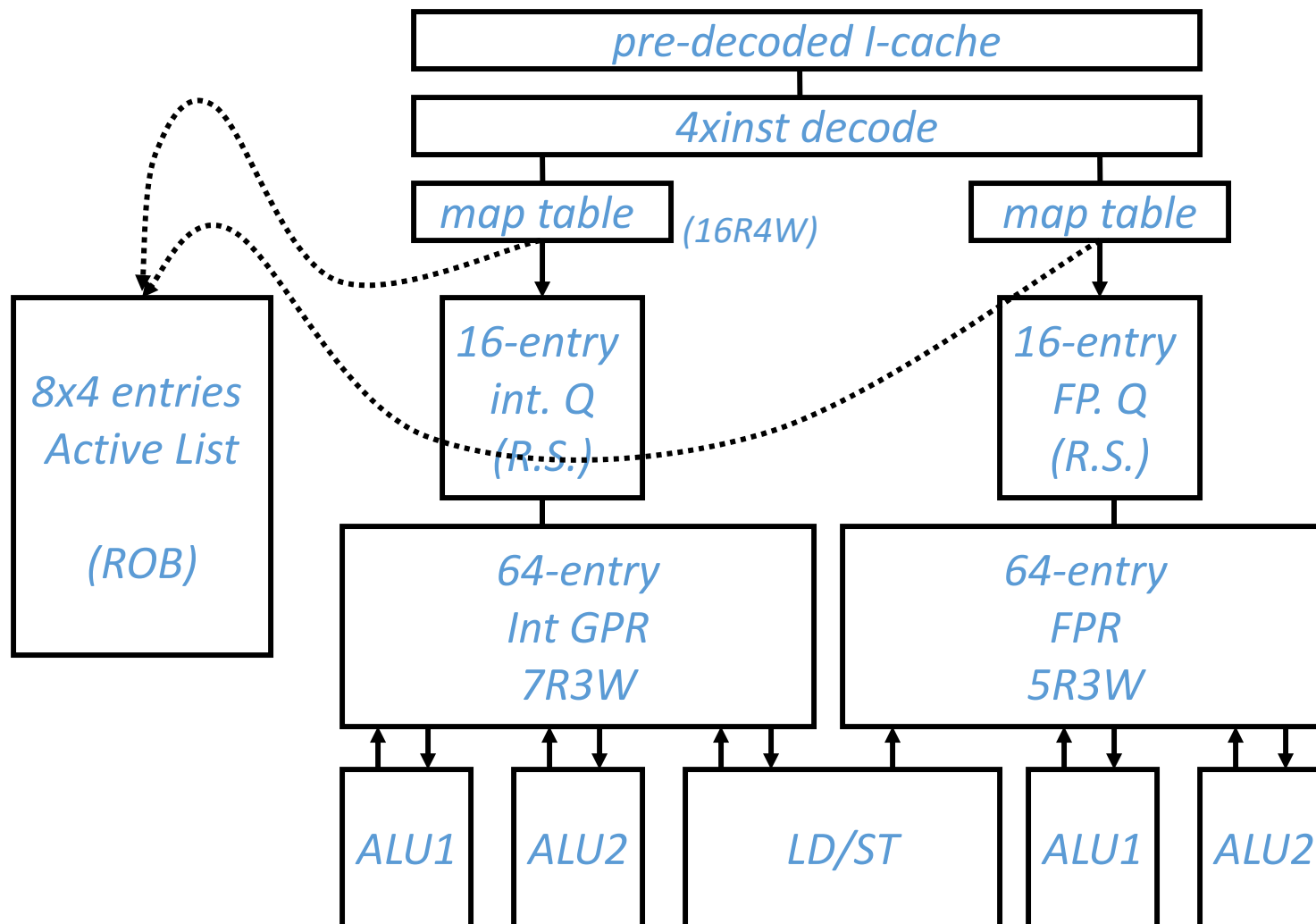


# OUT OF ORDER EXECUTION

---

- Execute instr. based on “data flow” rather than program order.
- **Basic idea:**
  - Fetch fills up a window of instructions.
  - Of all instructions in the window, look for ones ready to execute:
    - All the data on which the instructions are dependent are ready
    - Resources are available.
- As soon as instruction is executed, it needs to signal to its dependent instructions that the input is ready.
  - Triggers “wake-up” and “instruction select” for next cycle
- **Advantages:**
  - Help exploit Instruction Level Parallelism (ILP)
  - Help cover latencies (e.g., cache miss, divide)

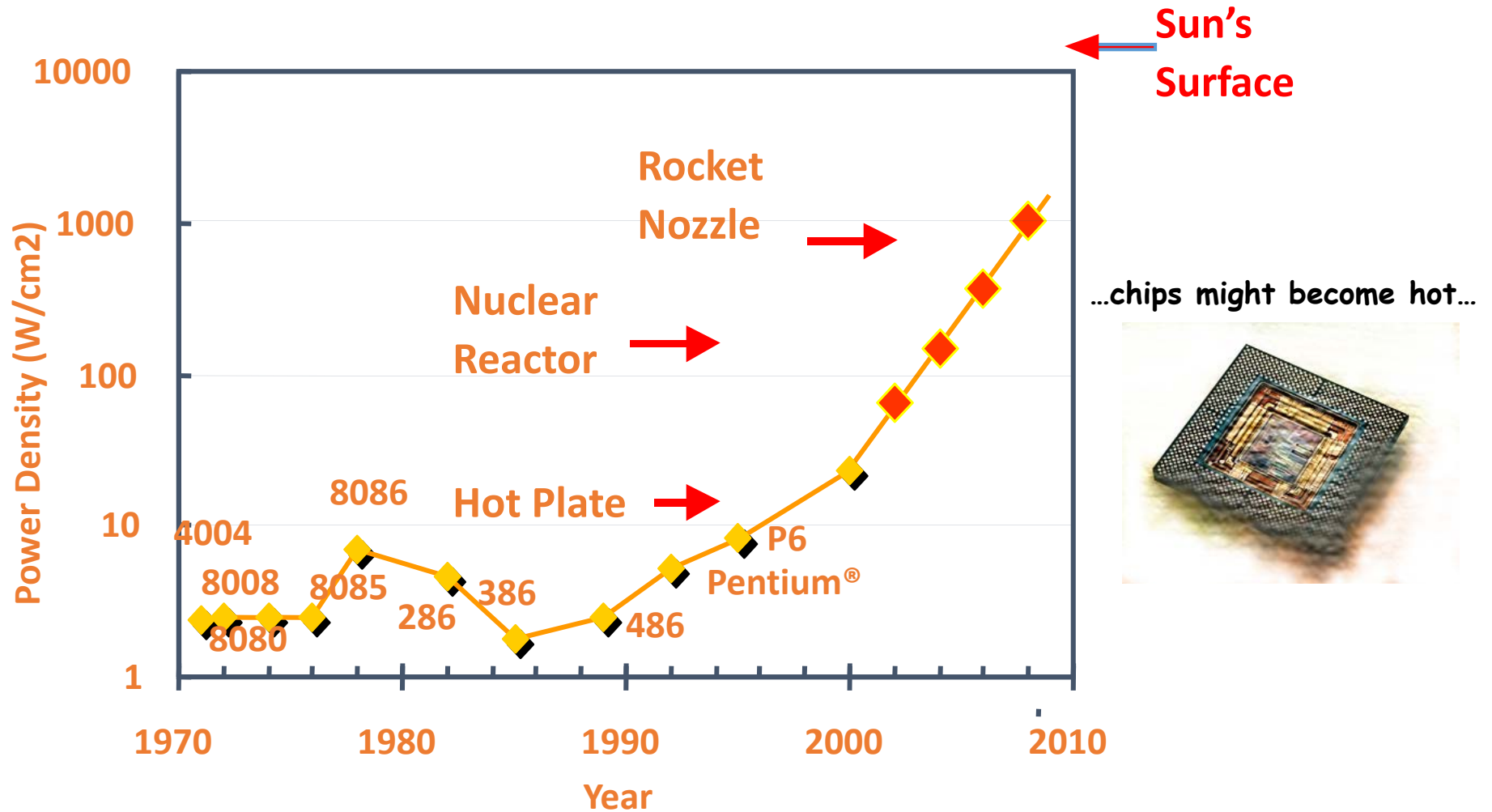
# MIPS R10000 (1996)







# CHIP POWER DENSITY



Source: Borkar, De Intel®

# THE COST OF SPECULATION

---



- How does the processor find enough “ready” instructions?
  - Look beyond branch boundaries
  - Use fairly sophisticated branch prediction mechanisms to guess target and direction.
  - Front end fills instruction window with instructions down path of predicted branch
- What happens if the branch prediction is wrong?
- What is the cost in terms of performance and power of keeping these “wrong path” instructions in the instruction window?

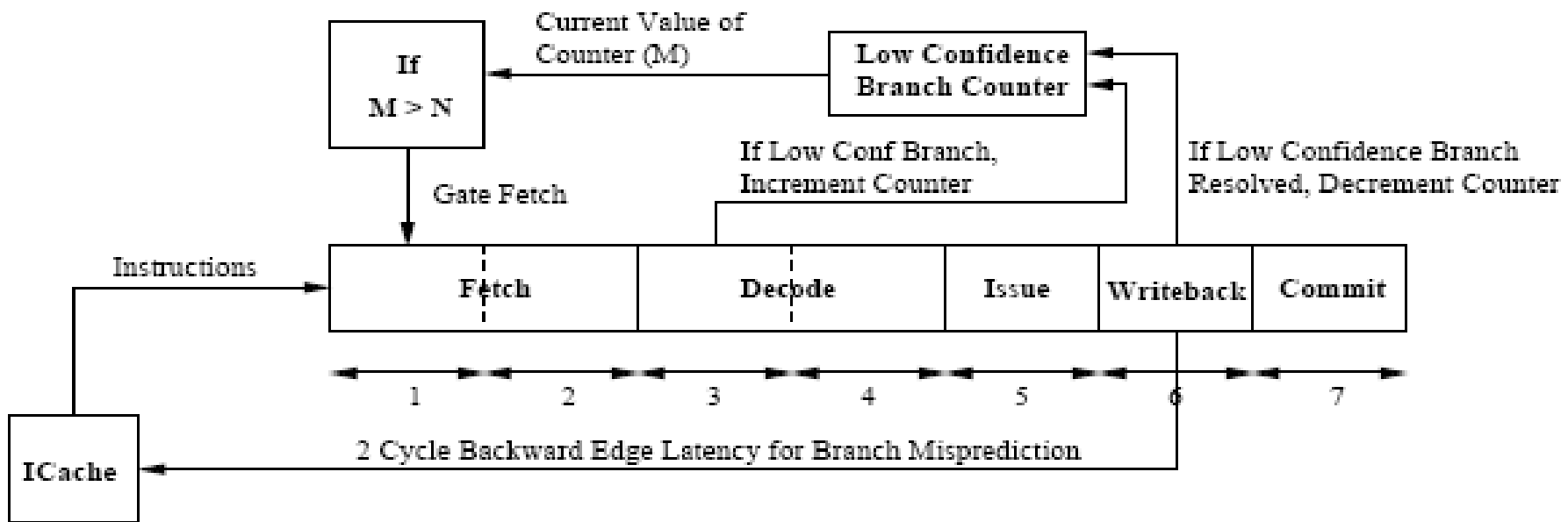
# REDUCING MIS-SPECULATED INSTRUCTIONS TO SAVE POWER

---



- Executing wrong path instruction is a waste of power
  - Does nothing to improve effective IPC either
- **IDEA:** if branch prediction becomes too speculative, don't bother continuing to fetch instructions past branches
  - Fetch unit stops reading new instructions from the cache
  - Instruction window does not take new instructions
  - Instruction execution rate may slow, but only until predicted branches have been resolved
- How do we know if an instruction flow has become “too speculative”?

# PIPELINE GATING (ISCA 1998 MANNE ET AL.)



- Low-confidence branch counter records # of unresolved branches that reported as low-confidence.
- The processor gates off instruction fetch if there are more than  $N$  unresolved low-confident branches in the pipeline.

➔ **ISCA 2013 Influential Paper Award !**

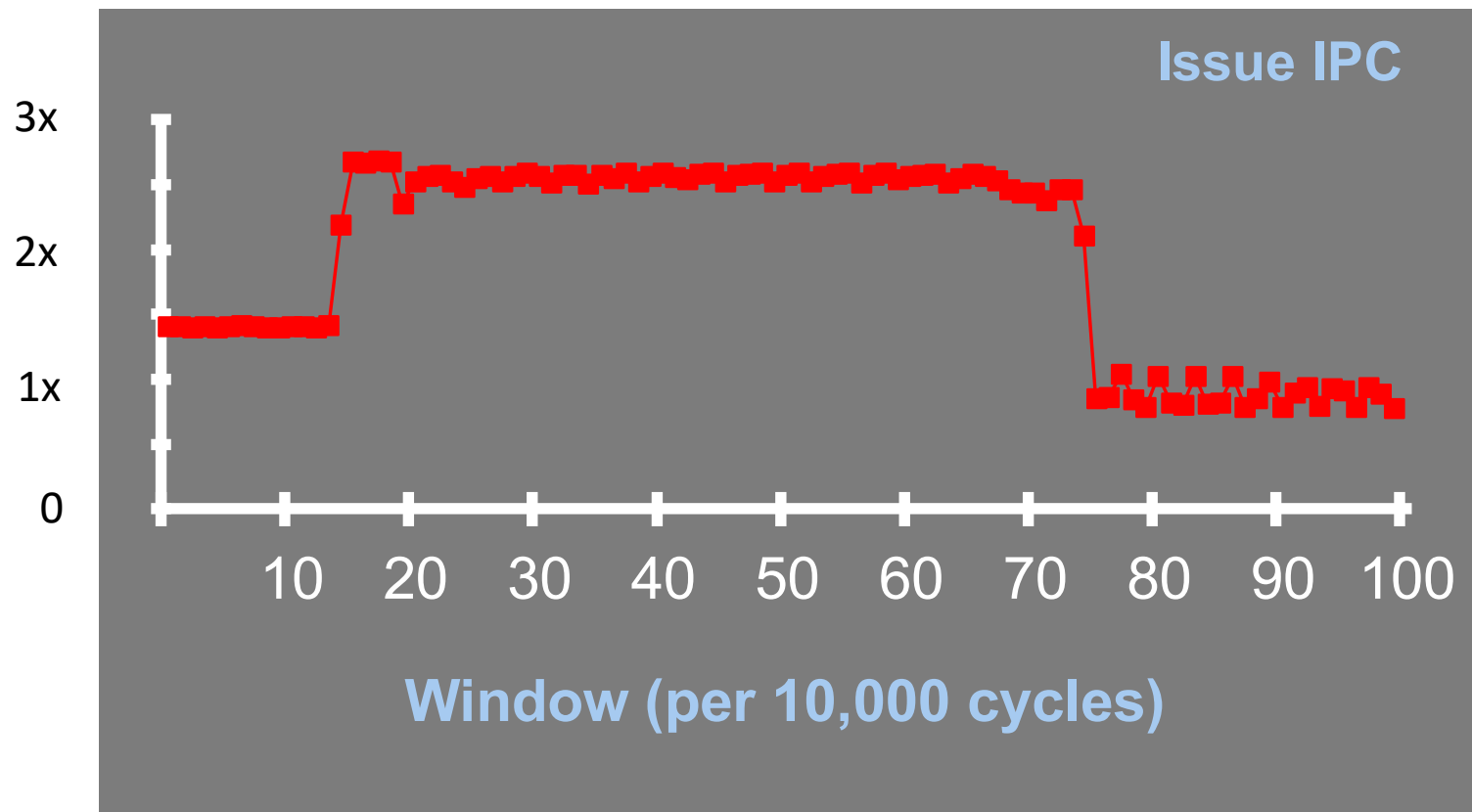
# DO WE NEED WIDE ISSUE MACHINES?

---



- Many programs never achieve IPC values close to the maximum issue of the machine
  - Branch Misprediction
  - Dependency Chains
  - Cache Misses
- Overall IPC is not indicative of superscalar needs of a program

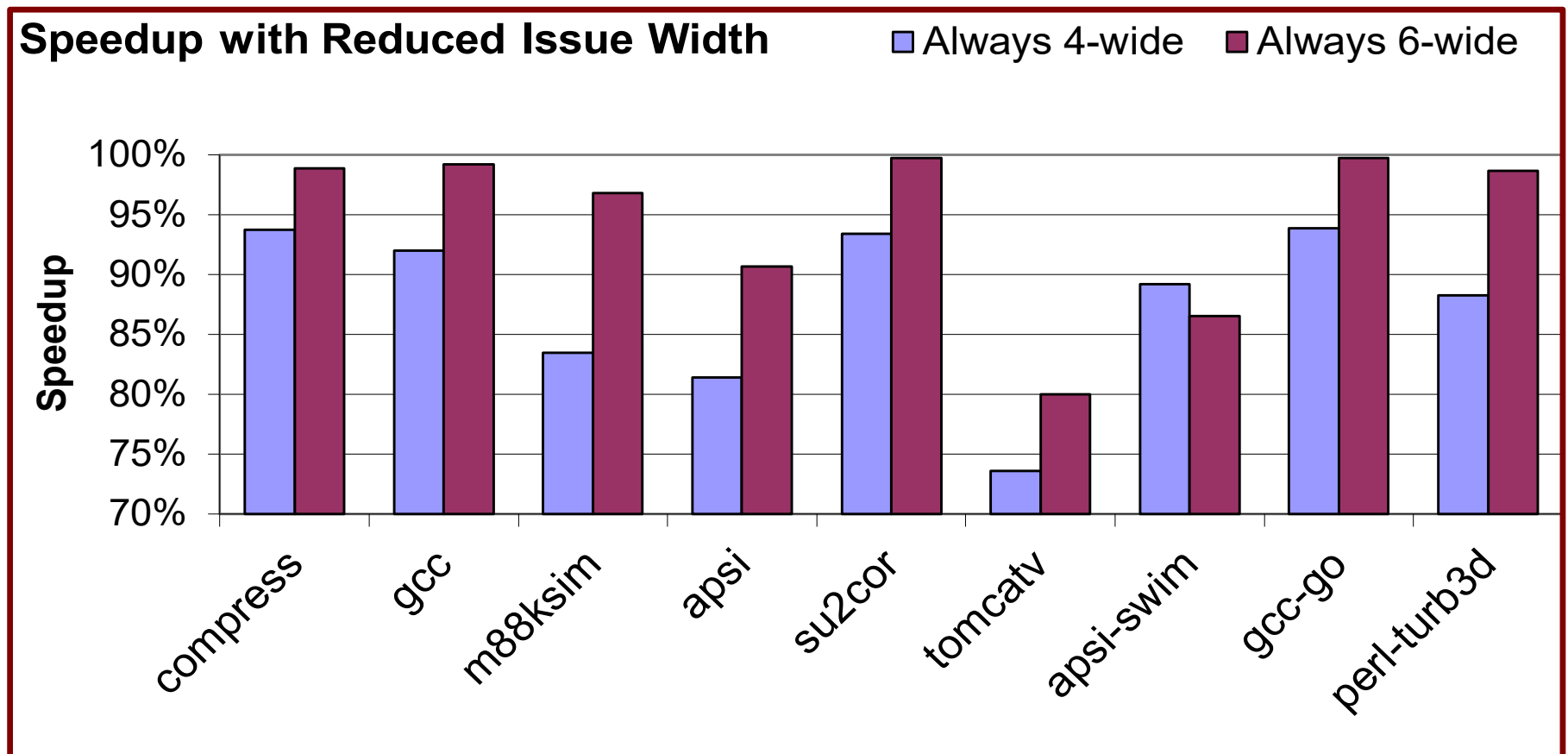
# VARYING PROGRAM NEEDS



- Nearly 3X difference in IPC across successive snapshots of the program execution



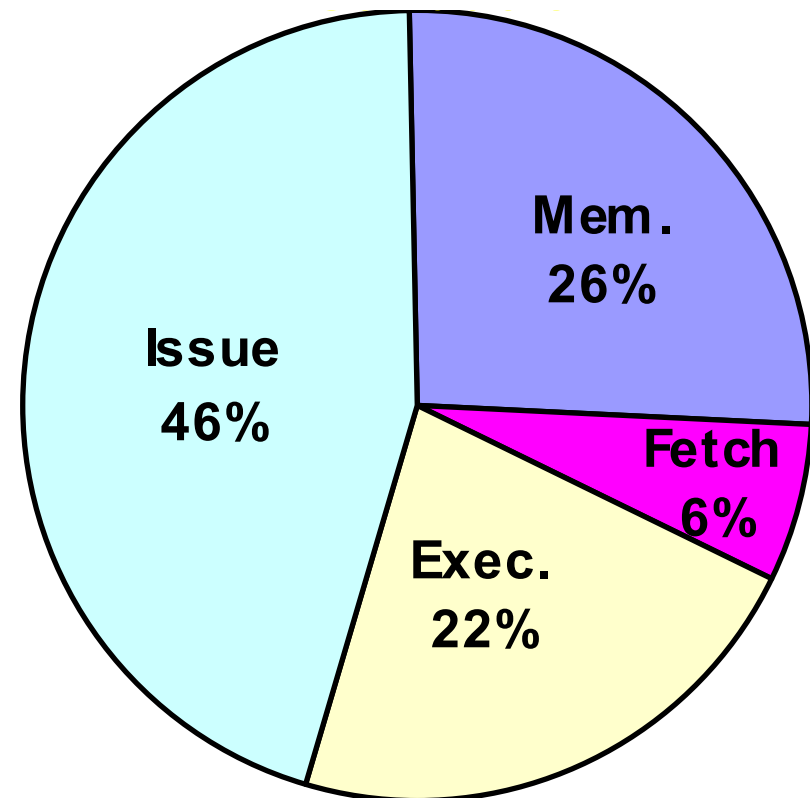
# IMPACT OF REDUCED ISSUE WIDTH



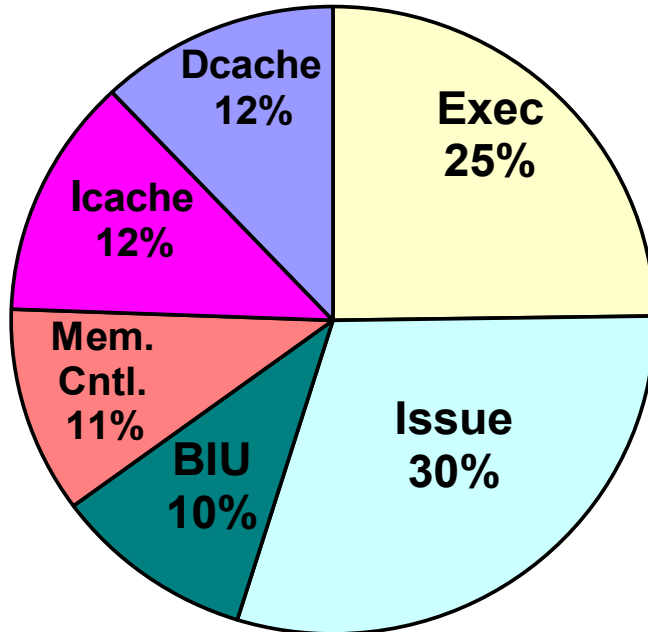
# MICROARCHITECTURAL DECISIONS IMPACT POWER



## 21464 Power Distribution



## 21264 Power Distribution





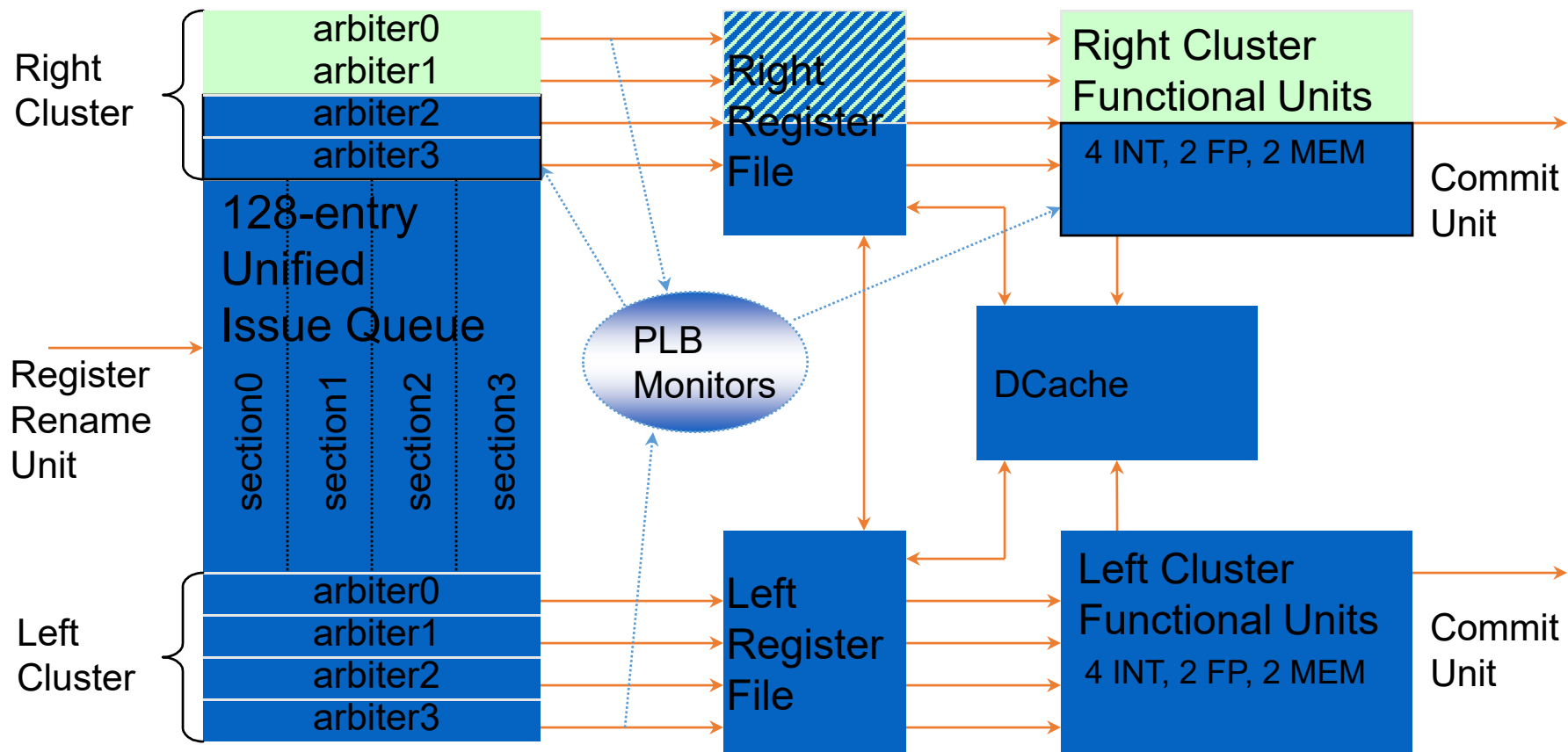
# PIPELINE BALANCING (ISCA 2001, BAHAR, MANNE)

---



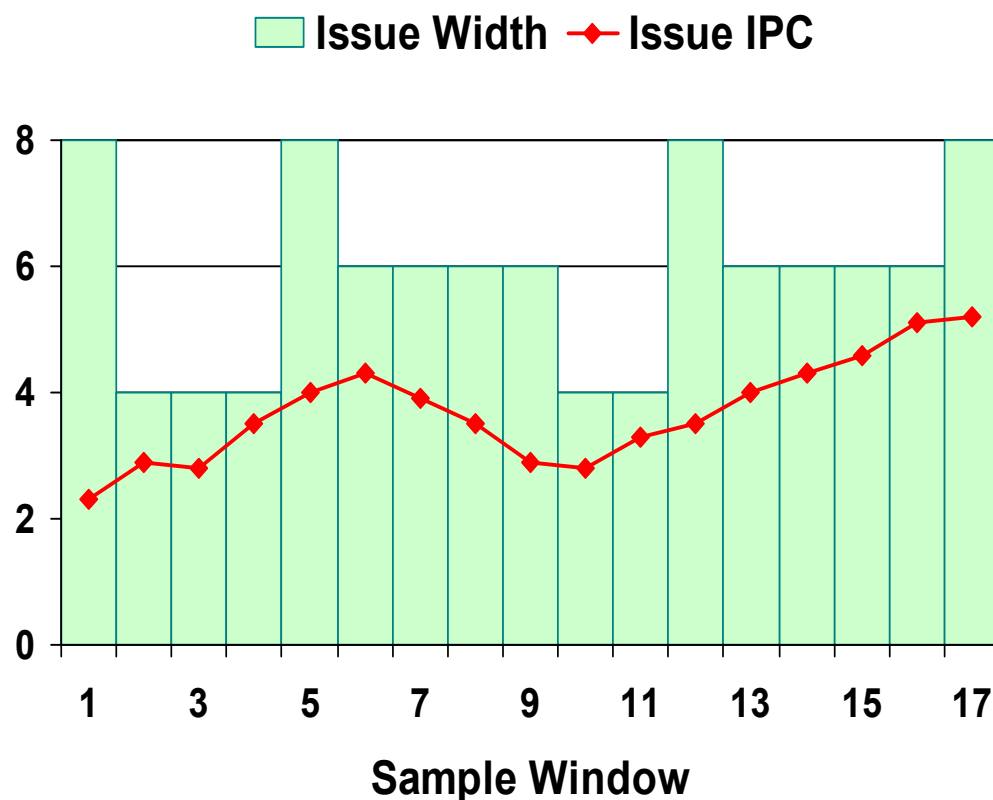
- Monitor the varying issue requirements of each program
  - Overall issue rate
  - Floating point issue rate
  - History of past behavior
- **IDEA:** Tune processor issue and execution resources according to the needs of the program
- **Goal:** Reduce power, retain performance

# OUR 8-WIDE ISSUE PROCESSOR MODEL





# TRIGGERING PLB BY TRACKING ISSUE IPC

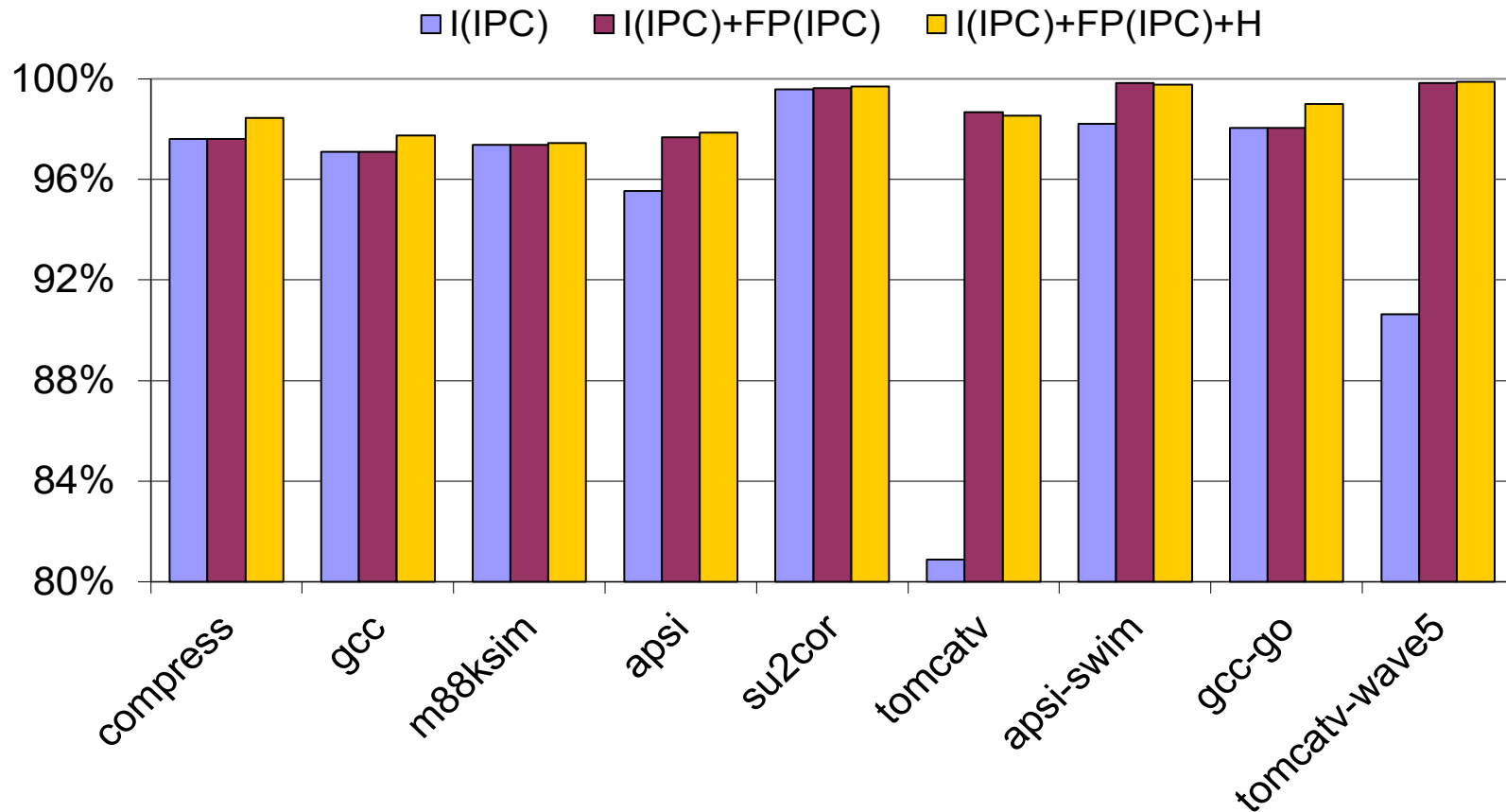


- Switch to 4-, 6-, or 8-wide issue depending on issue IPC of previous sampling window

# PRIMARY AND SECONDARY TRIGGERS



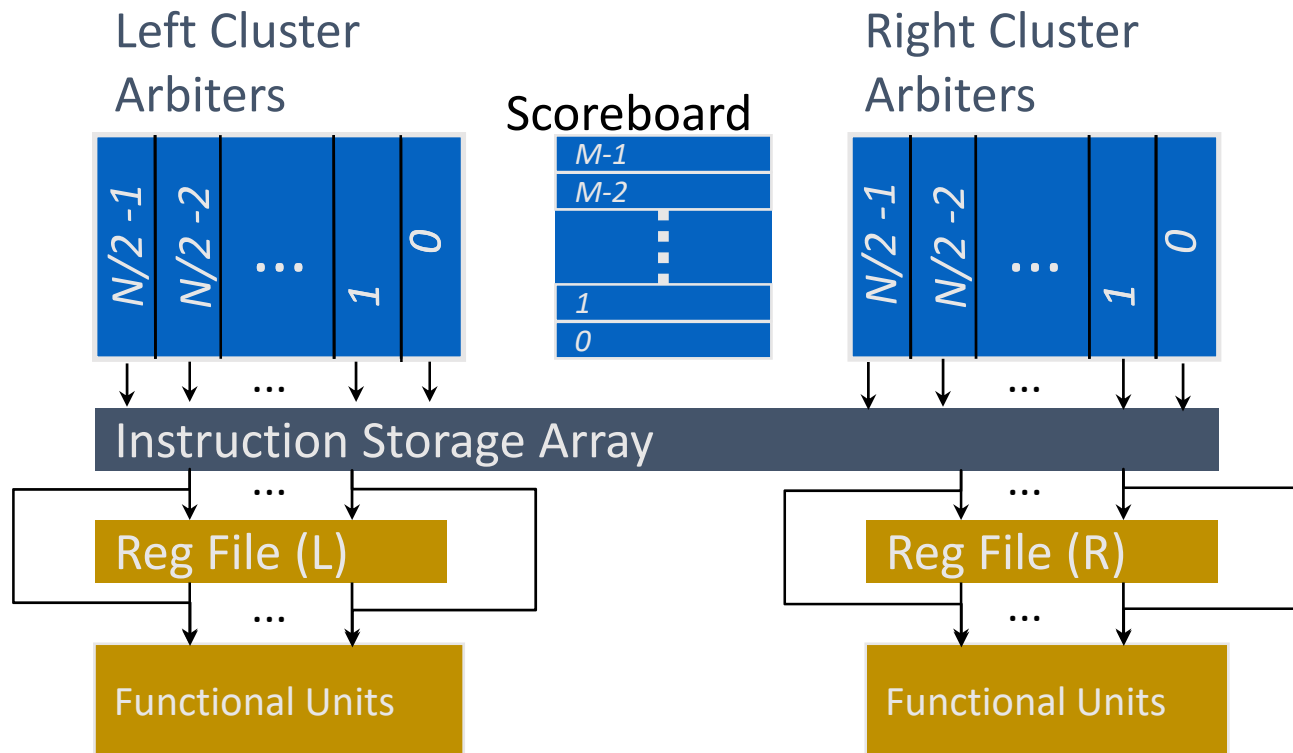
Speedup Using FP IPC and Mode History as Secondary Triggers



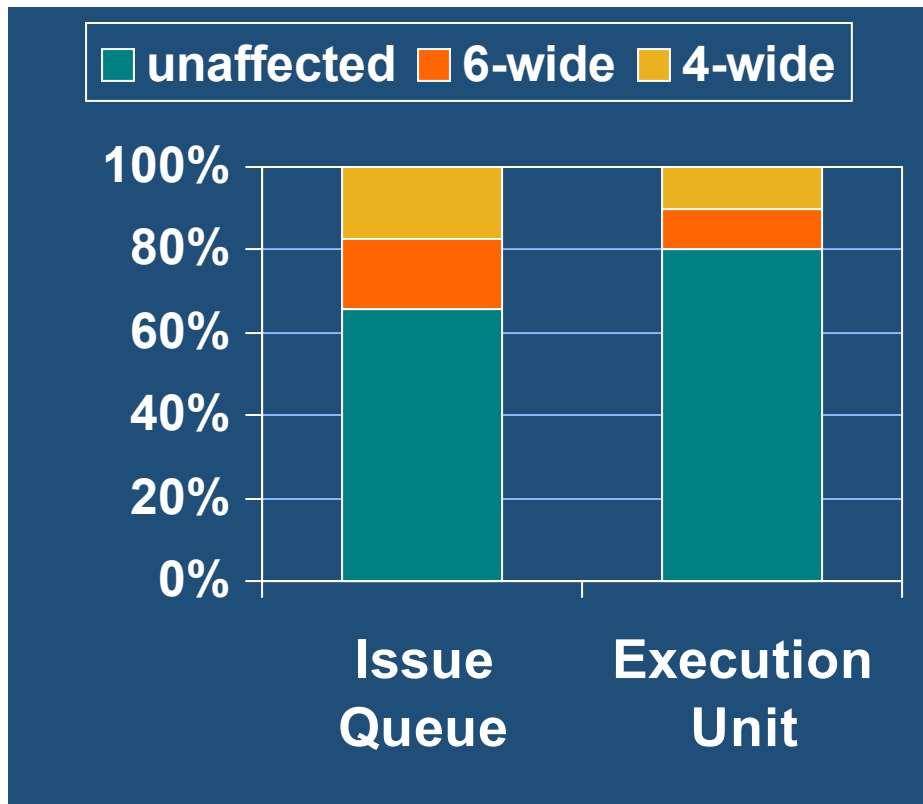


# IQ DESIGN AND POWER DISSIPATION

- Goal: Issue  $N$  instructions out of  $M$  entries in IQ
- IQ power is relative to the size of  $N$  and  $M$



# PER COMPONENT POWER ESTIMATES

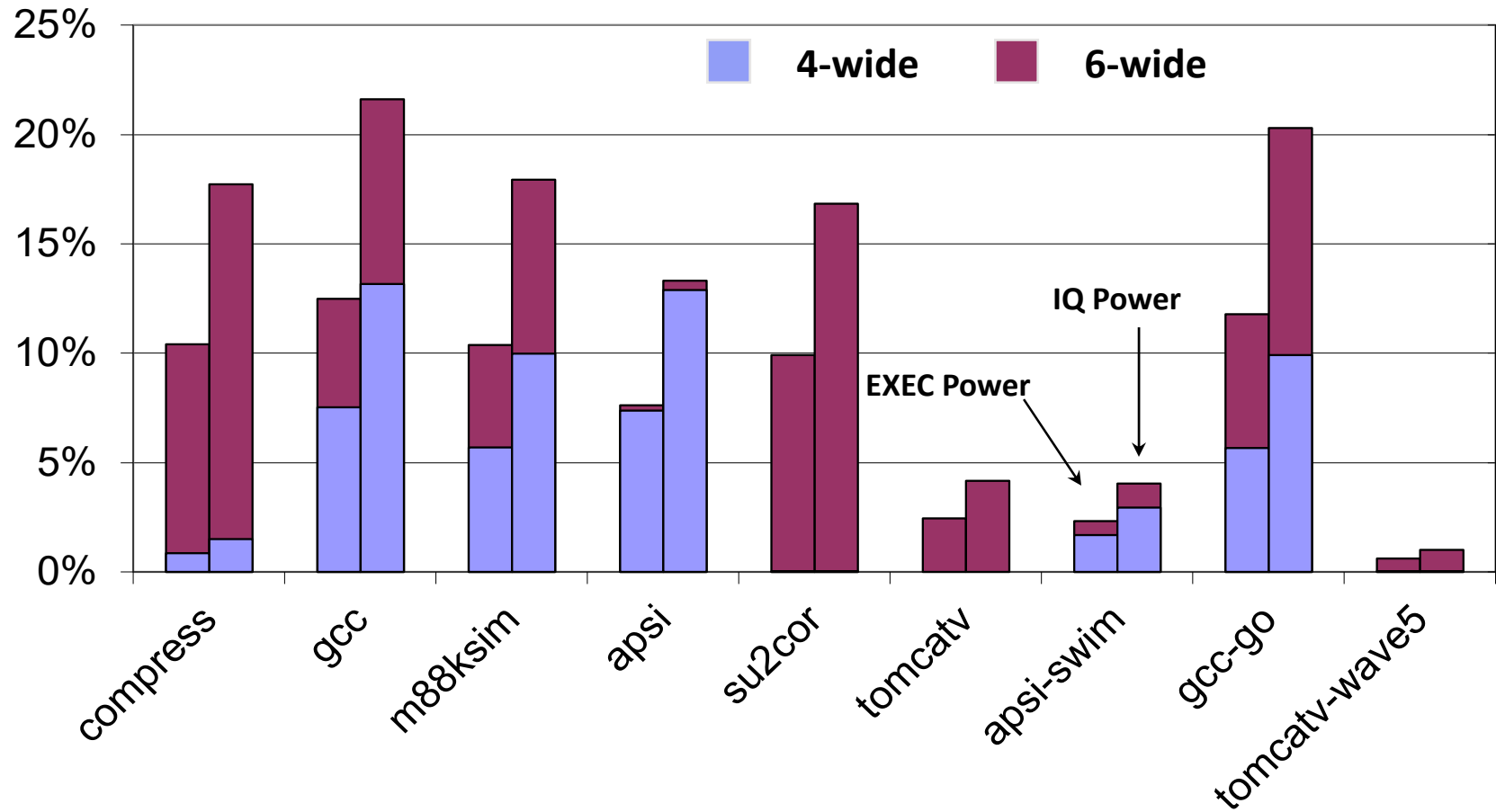


- Up to 35% of Issue Queue power is saved
  - Arbiter enable signals
  - Bid logic
  - Selection logic
- Up to 20% of Execution Unit power is saved
  - Disable clocks in unused portions

# COMPONENT POWER REDUCTION WITH PLB

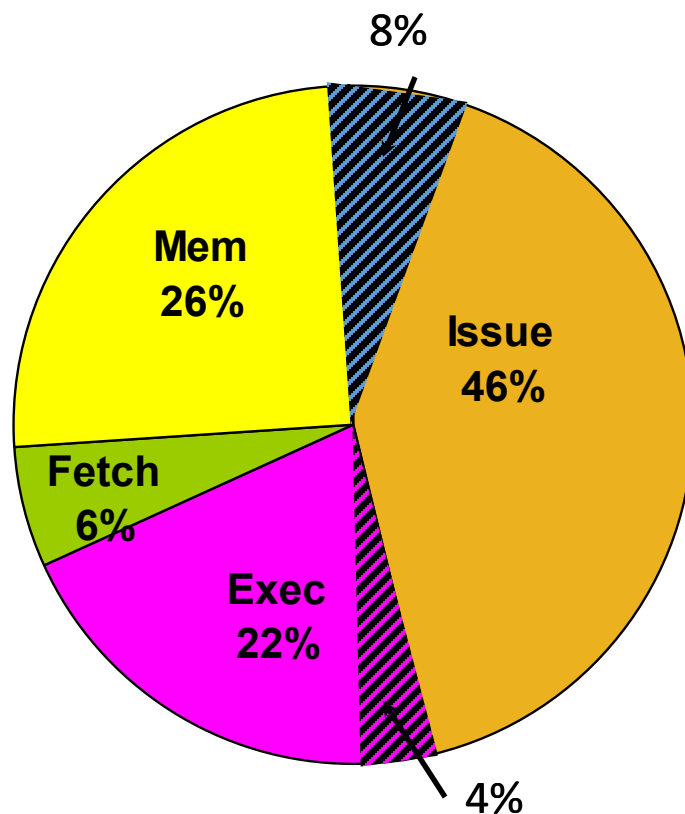


### IQ and Exec Power Savings





# CHIP-WIDE POWER ESTIMATES



- 17-35% of IQ power saved  
➤ % of Total power: 4-8%
- 10-20% of Execute Unit power saved  
➤ % of Total power: 2-4%
- Up to 12% of total power saved using PLB





# ENERGY-EFFICIENT TRANSACTIONAL MEMORY FOR EMBEDDED SYSTEMS

# MULTICORE ARCHITECTURES

---

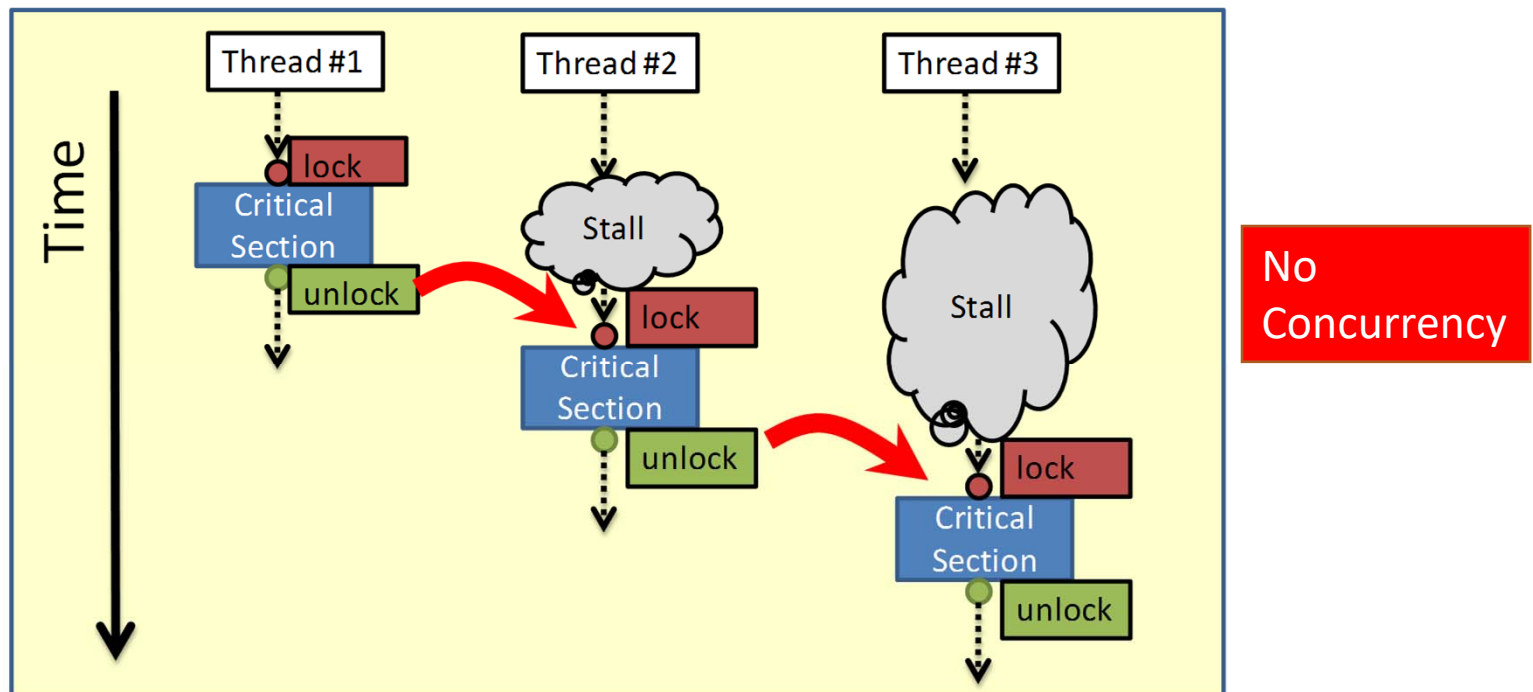


- Like GP systems, embedded platforms have turned to multicore architectures
  - Better scalability and design cycle time
- Issues:
  - Required High Level Of **Parallelism**
  - **Synchronization** Must be Efficient
- What's Special about Embedded Multicore:
  - **Simplicity!** (small caches, no OS support)
  - **Power matters** (as much as throughput)

# MULTIPROCESSOR SYNCHRONIZATION



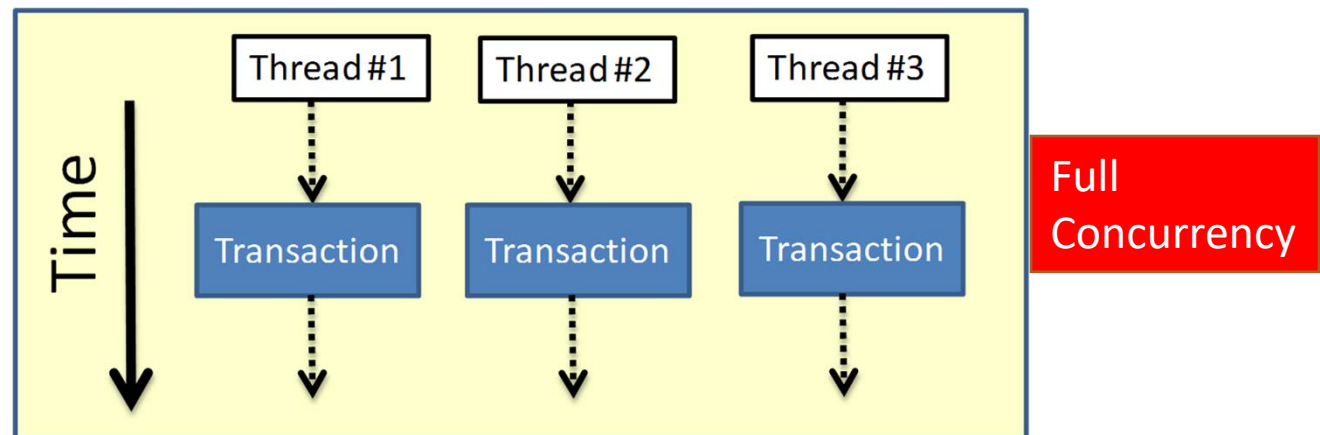
- Focus: **Shared Memory Model** (widely accepted)
- Common Approach: **Lock-based** Synchronization
  - **Idea**: mutual-exclusive access to shared data (i.e., critical sections)



# TRANSACTIONAL PROGRAMMING



- Key Features:
  - Optimistic (i.e., concurrent) execution of critical sections
  - Threads run in isolation
  - Atomicity (roll-back in case of conflict)



- ❑ **Higher Performance** than Fine-grained Locking (*without* reliability issues)

# TM FOR EMBEDDED SYSTEMS



- Software-TM schemes (STM or HW-accelerated STM):
  - Flexible, but heavyweight
  - Require OS support
  - Too slow (and power hungry...)
- Hardware-TM schemes (HTM):
  - Less flexible, but **simple, power-efficient**, and higher performance
  - Transaction sizes may be bounded, but...
  - For embedded systems, resource requirements are well understood

*Full HW-TM seems a natural design choice for embedded systems*

## □ ***New challenges :***

- Peculiar Memory Hierarchy (eg. tightly-coupled private memories)
- Limited resources (eg. smaller caches, no OS)
- Energy Constrained environment

# 1<sup>ST</sup> WORK IN EMBEDDED HTM: SoC-TM

---



- Simplify contention management
  - Disentangle HTM from cache coherence
  - Bloom Module for conflict detection
  - More flexibility for conflict resolution policies
- Create an integrated HW/SW solution for transactional programming on MP-SoCs
  - Expands on Transactional OpenMP (OpenTM)
  - Adds support for speculative task- and data-level parallelism

# SIMPLIFYING CONFLICT DETECTION

---

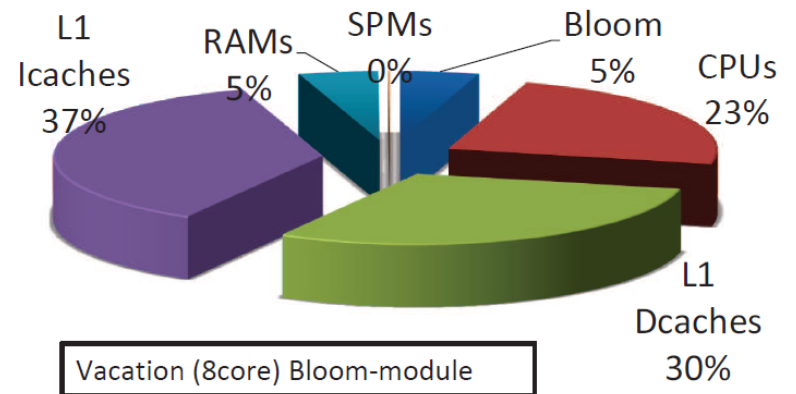
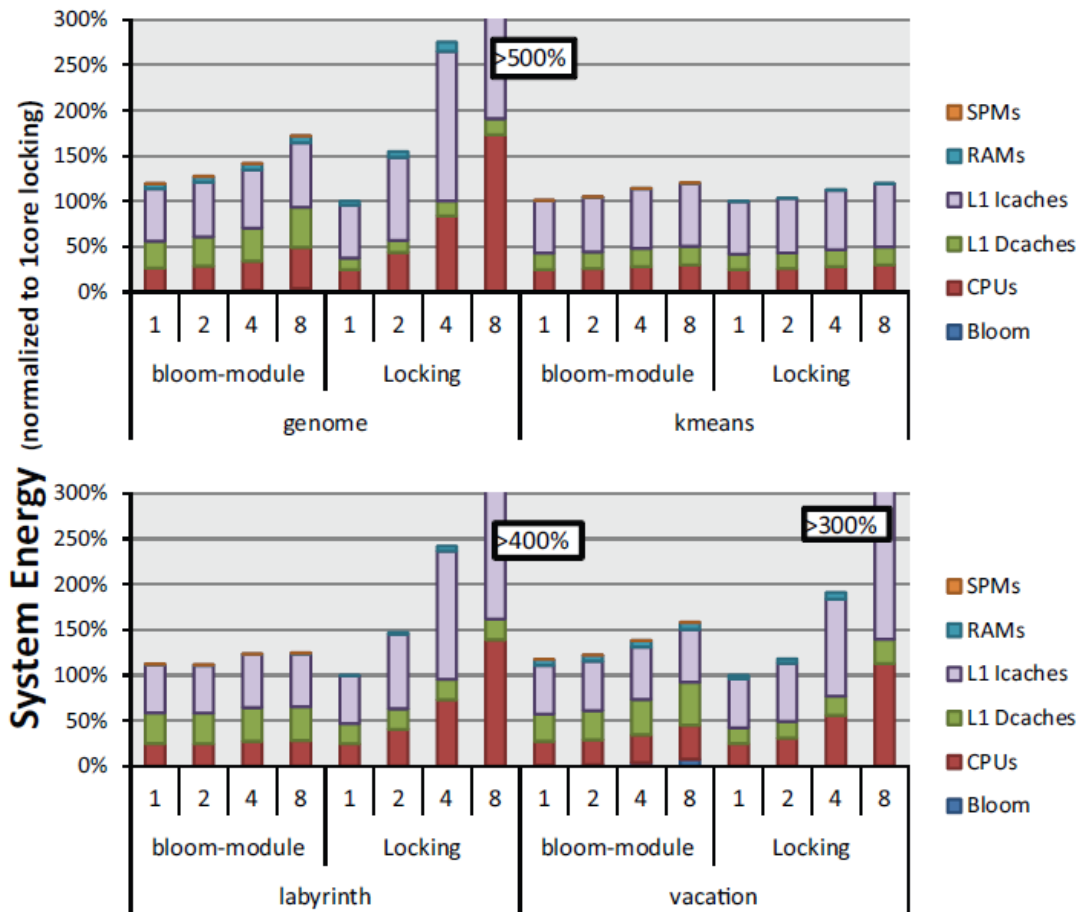


- Originally, conflict detection was through the **cache coherency protocol** of each core
  - Requires modification to MESI protocol
  - Hard to exercise control over conflict resolution policy
- **Bloom Module:** centralized module for read/write conflict detection
- Key Features:
  - Decoupled from Cache Coherency
  - Keeps Tx histories in per-core Bloom signatures
  - Detects conflicts by comparing bits in the signatures
  - Programmable (Any algorithm can be used by the module to decide which core to abort)

# EXPERIMENTAL RESULTS



- STAMP benchmarks
  - energy consumption comparison with locks





# 2<sup>ND</sup> WORK: EMBEDDED-TM

---



- **Limitations of SoC-TM**

- Used a *centralized* unit to snoop on the bus traffic.
- Bus-based architecture: one transaction/cycle.

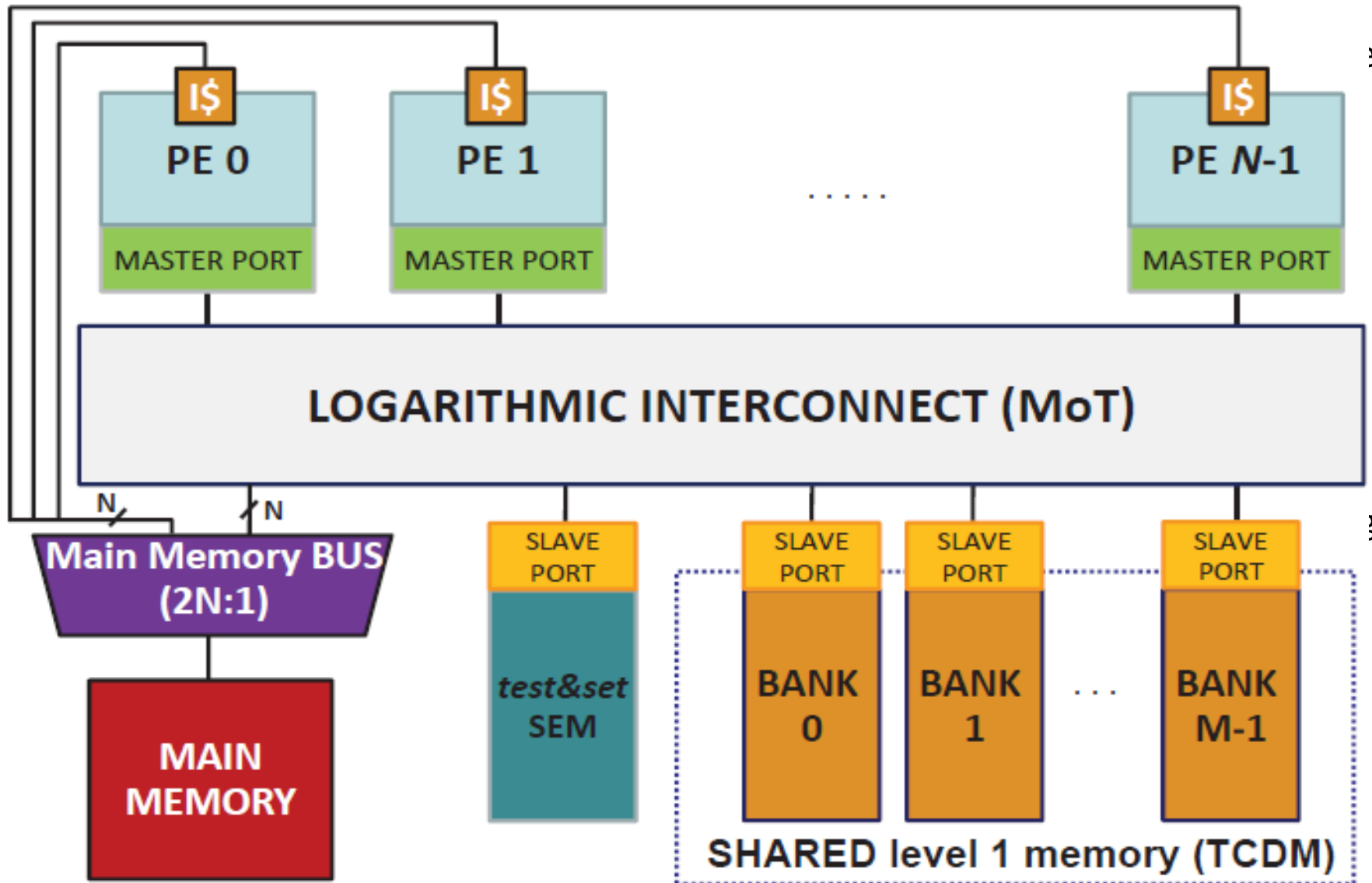
- **Next Work:** First implementation of HTM for cluster-based NUMA embedded systems without cache-coherence.

- Lack of cache-coherence brings major challenges in HTM design.

- **Contributions:**

- Building from scratch a self-contained HTM scheme.
- Propose novel hardware support for tracking memory accesses.
- Single-cluster but extendable to multi-clusters.
- Performance evaluation.

# TARGET ARCHITECTURE





# EMBEDDED-TM IMPLEMENTATION

---

- **In our case:** Logical Interconnect allows multiple transactions/cycle. Snooping → bottleneck, limits scalability.

- Multiple per-bank **Transactional Support Modules (TSM)**.

- Distribute conflict detection/resolution across TCDM banks.



## **Distributed Transactional Memory Control.**

- Our TM support design:
  - a) Transactional Bookkeeping
  - b) Data Versioning
  - c) Rollback Mechanism.

# A) TRANSACTIONAL BOOKKEEPING

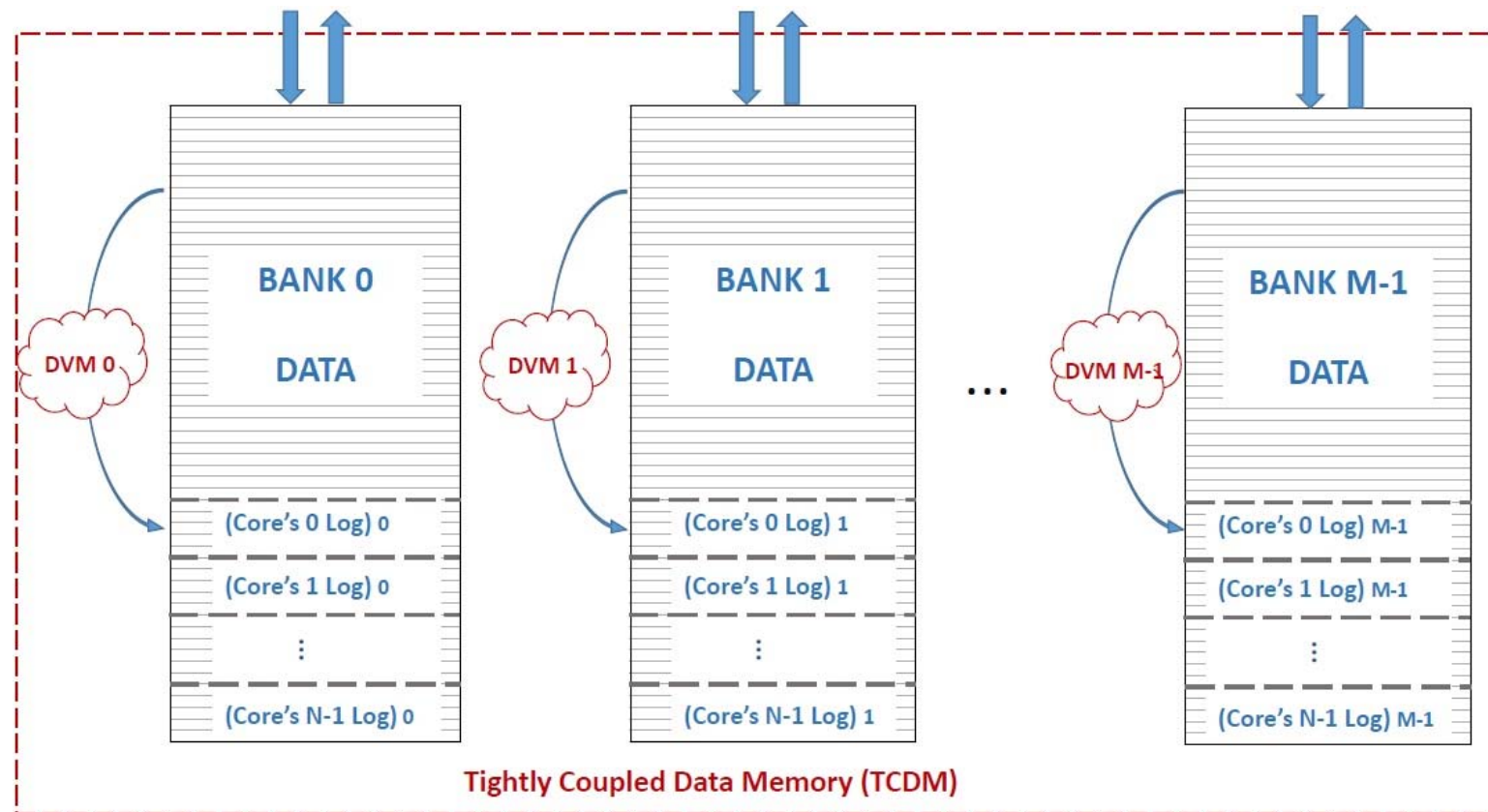


- Mechanism that keeps track of transactional readers/writers in memory.
- Each bank's TSM intercepts all memory traffic to bank. Upon conflict, it notifies conflicted cores.
- For each data line: Multiple Readers, Only one writer/owner.
- Write to a location read by another core → **Conflict!**  
Read to a location read by other cores → **No Conflict!**
- Each bank keeps track of readers/writers per data line through per-bank array of  $k$   $r$ -bit vectors,  
 $r = 1 + \log_2 N + N$ ,  $N$ : # cores in cluster,  $k$ : # data lines per bank,

	<b>Owner Bit</b>	<b>Writer ID</b>	<b>Readers (1 bit per core)</b>
<b>Time <math>t_1</math>:</b>	<b>0</b>	<b>0 0 0 0</b>	<b>0 0 0 0 0 0 0 1 1 0 0 0 0 0 1 0</b>
<b>Bit #:</b>	<b>20</b>	<b>19 18 17 16</b>	<b>15 14 13 12 11 10 9 8 7 6 5 4 3 2 1 0</b>



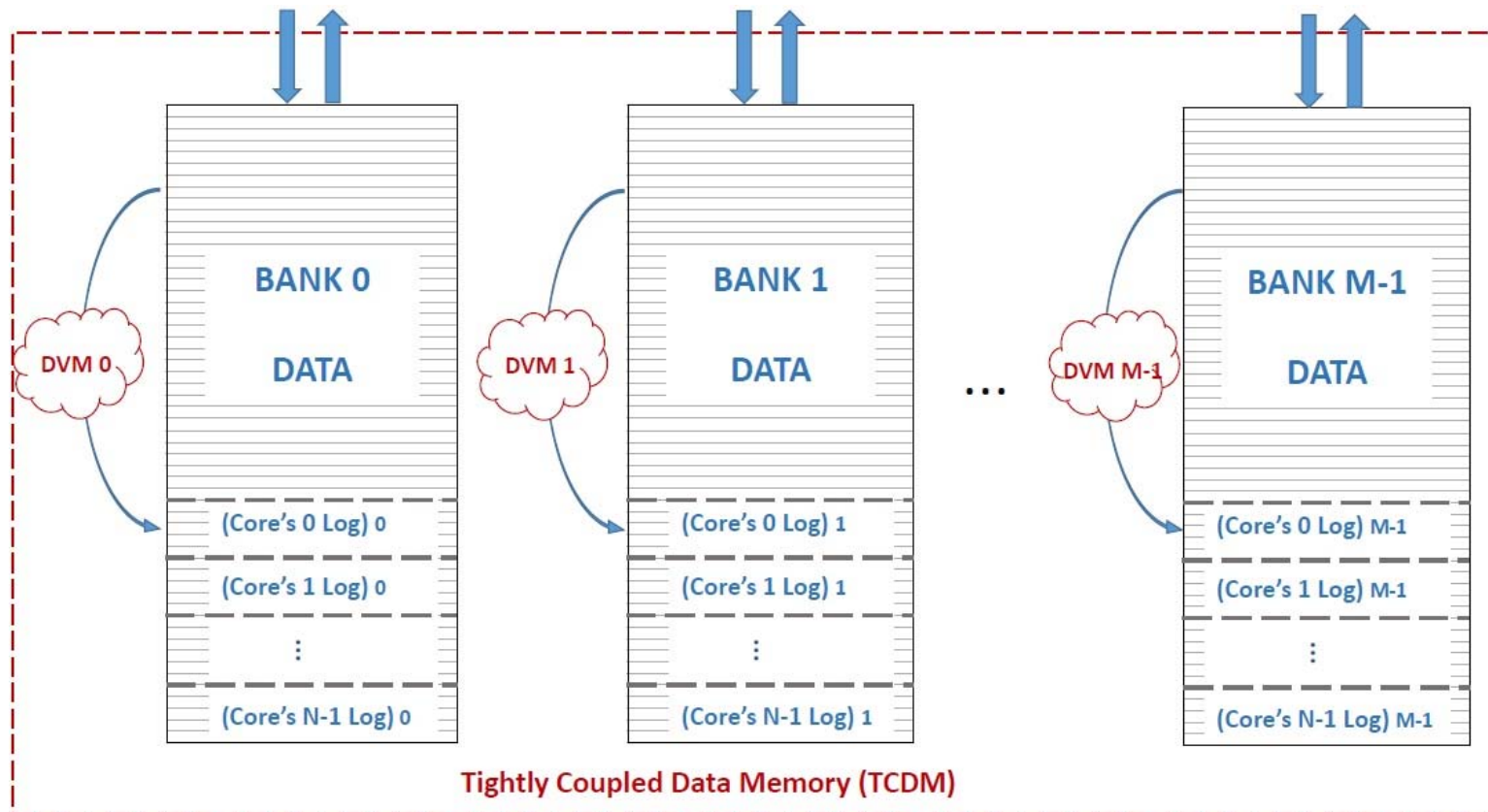
## B) DATA VERSIONING



- Original data saved in **Logs**, distributed among TCDM banks.
- In each bank: fixed-size log space per core.
- Each core's log includes entries for modified data for that bank only
- **Data Versioning Modules (DVM)** monitor accesses for its bank only.



## B) DATA VERSIONING



- **Transaction begins:** save state of core (PC, registers, SP, stack)
- **On a write:** DVM checks log for entry, creating new entry only once.
- **On a data conflict:** Each bank's log is traversed. Data restored.
- **Transaction completes with no error:** logs are cleared.

## C) ROLLBACK MECHANISM

---

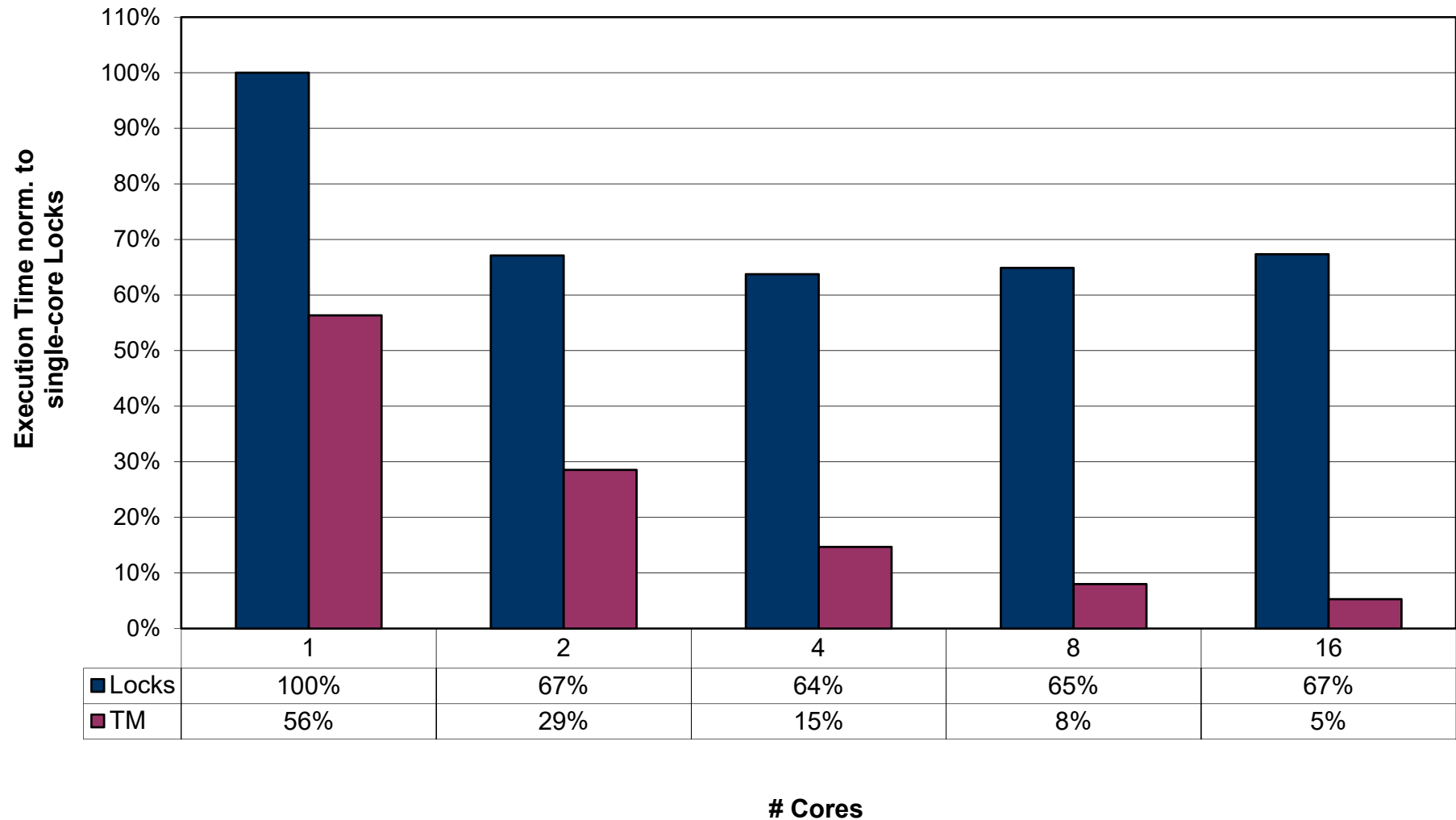


- Conflicts are detected while transaction executes
  - We choose to abort the “requester”
  - Triggers a call to *abort\_transaction()* function
  - Restores saved state from logs of all banks in TCDM for requester transaction
  - Aborted transaction tries again after some backoff period

# 1) REDBLACK



### Execution Time vs # Cores

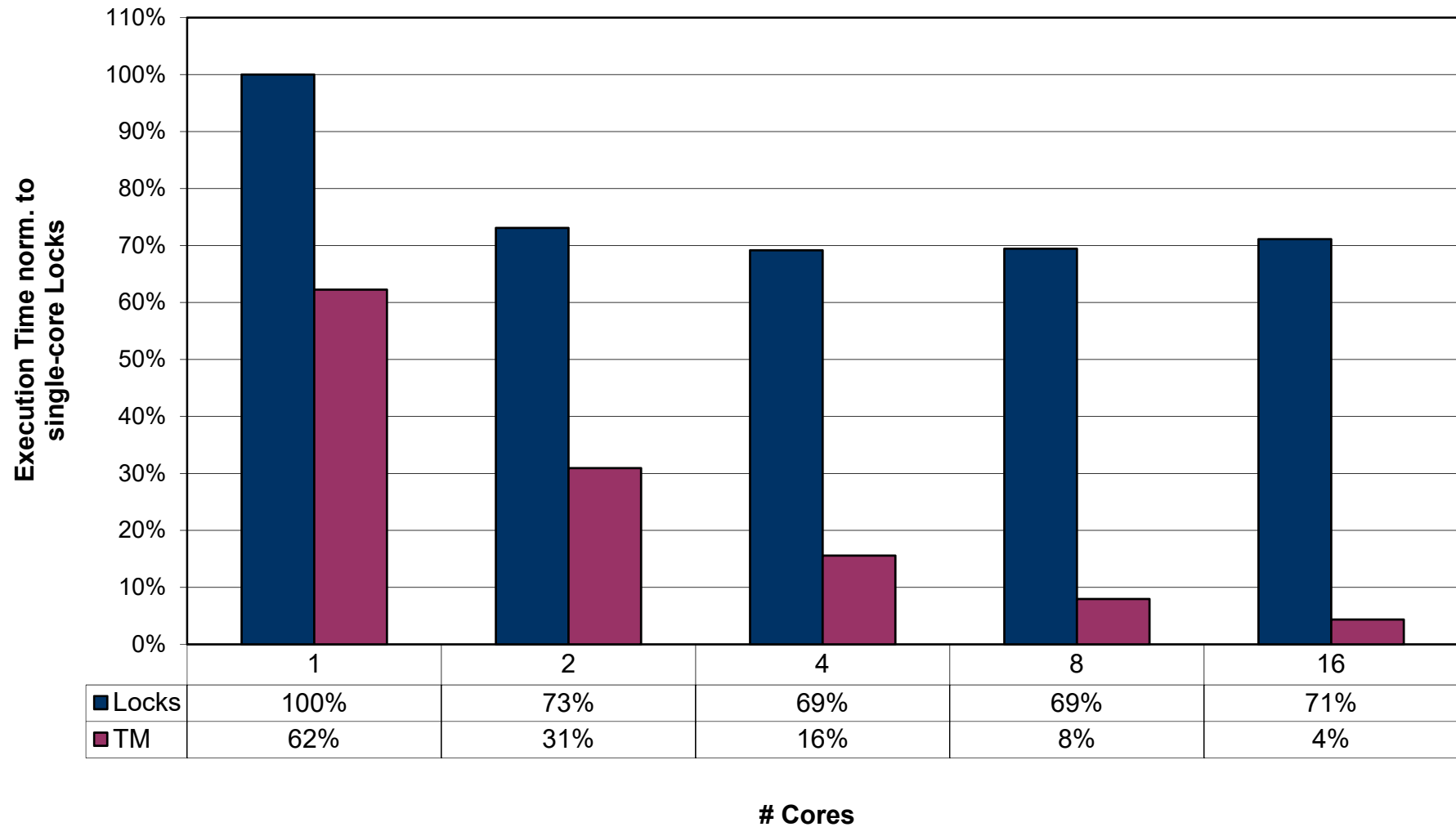




## 2) SKIPLIST



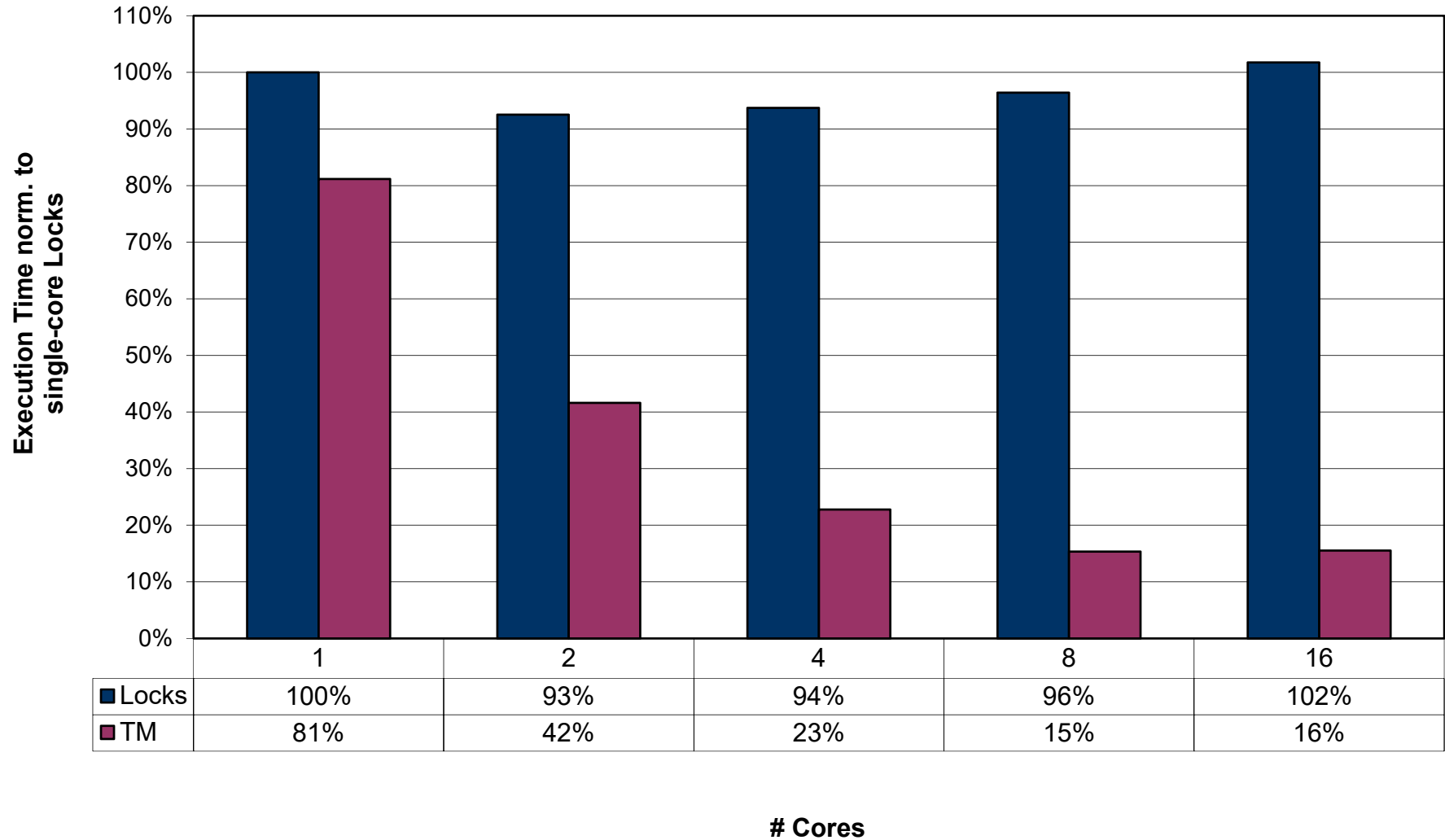
Execution Time vs # Cores



# 3) GENOME



### Execution Time vs # Cores





# EXPLOITING HARDWARE TRANSACTIONAL MEMORY FOR ERROR TOLERANCE AND ENERGY EFFICIENCY

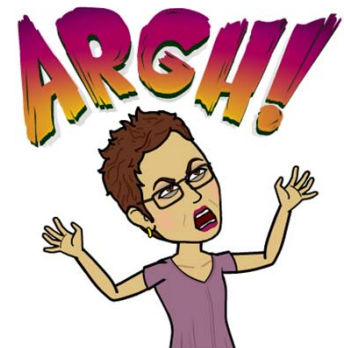
# VARIABILITY-INDUCED ERRORS



- Shrinking transistor sizes has left devices more susceptible to variability
  - **Static Variations:** lead to performance/power mismatches
  - **Dynamic variations:** aging, voltage drops, temperature fluctuations
- **Solution:** Conservative guardbands on the operating frequency/voltage

*Performance* ↓ *Energy*

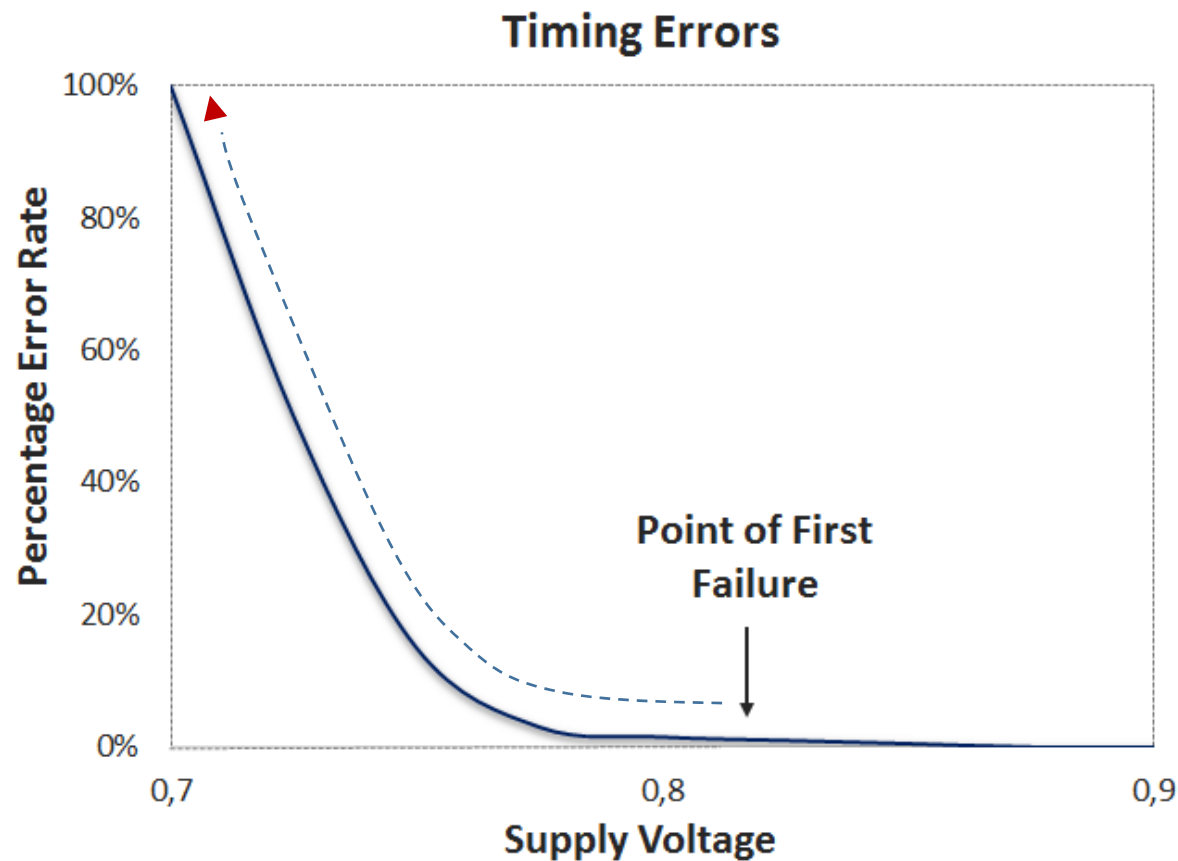
- **Reduce guardbands?**
  - Intermittent timing errors
  - Critical operating point (COP)



# A) INTERMITTENT TIMING ERRORS



**Intermittent Timing Errors** cause erroneous instructions with wrong outputs or incorrect control flow.

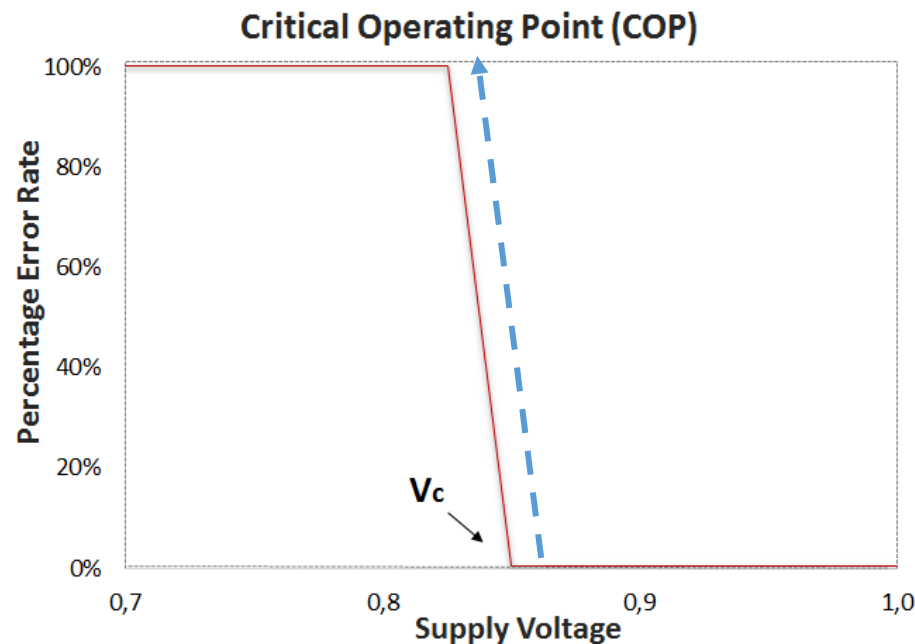


# B) CRITICAL OPERATING POINT



**Critical Operating point (COP):** ( $V_c$ ,  $F_c$ ) pair that if surpassed:

⇒ **Massive instruction failures.**



**COP Hypothesis:** In large CMOS circuits, there exists a critical operating frequency  $F_c$  and a critical voltage  $V_c$  for a fixed temperature  $T$ , such that:

- Any  $F > F_c$  causes massive errors.
- Any  $V < V_c$  causes massive errors.
- Any  $F < F_c$  and  $V > V_c$ ,  
no errors occur.

COP might change over space/time.

In a many-core environment:

- COP might differ among cores
- Become unsafe.

# CONTRIBUTIONS

---

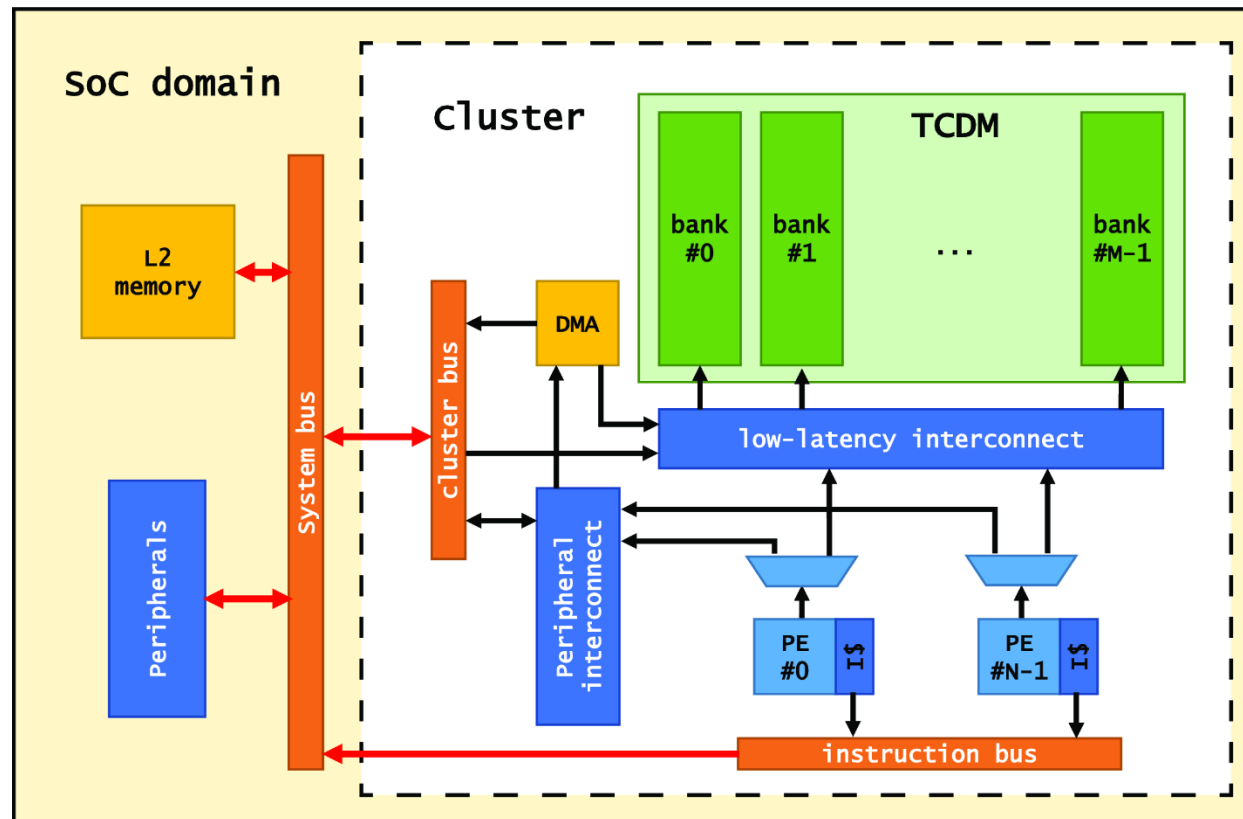


- We need a recovery mechanism to **avoid guardbands** and operate at lower voltage to **save energy**.
- **Edge-TM**: HW/SW scheme for **error resilience** and **energy efficiency**:
  - Scales down the voltage to save energy
  - Monitors errors
  - **Hardware Transactional Memory (HTM)** for error recovery (Intermittent timing errors and COP).
- Our approach:
  - Enables operation at **highly reduced supply voltage margins**.
  - **Forward progress**.
  - **Low overhead**.

# EDGE-TM ARCHITECTURE DESIGN



- Parallel Ultra-Low Power Platform (PULP): Cluster-based shared memory embedded architecture.



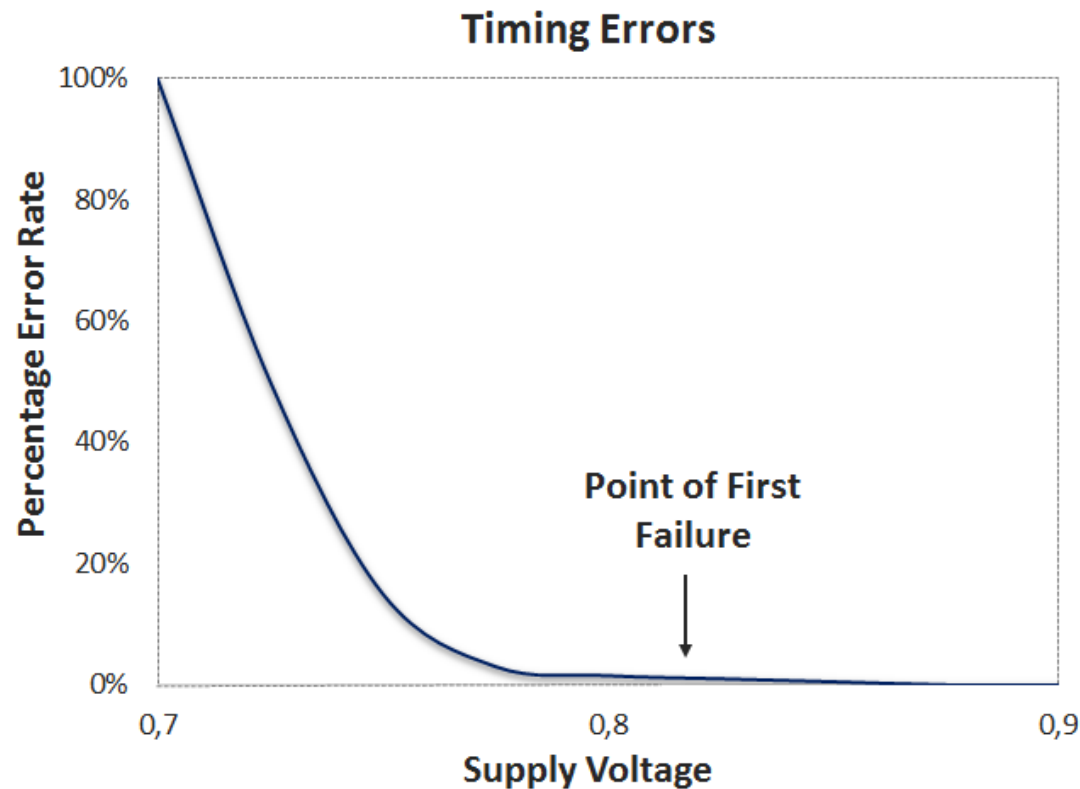
- Edge-TM:
  - Hardware extensions for error tolerance: a) Error Detection, b) Error Correction.
  - Software directed DVS error management policies for energy efficiency.



# A) ERROR DETECTION



- **Error Detection Sequential (EDS):** Error detection circuitry for each core.
- Probabilistic Error Model (Fojtik et al.): Expected error rate vs. Supply Voltage for intermittent timing errors.

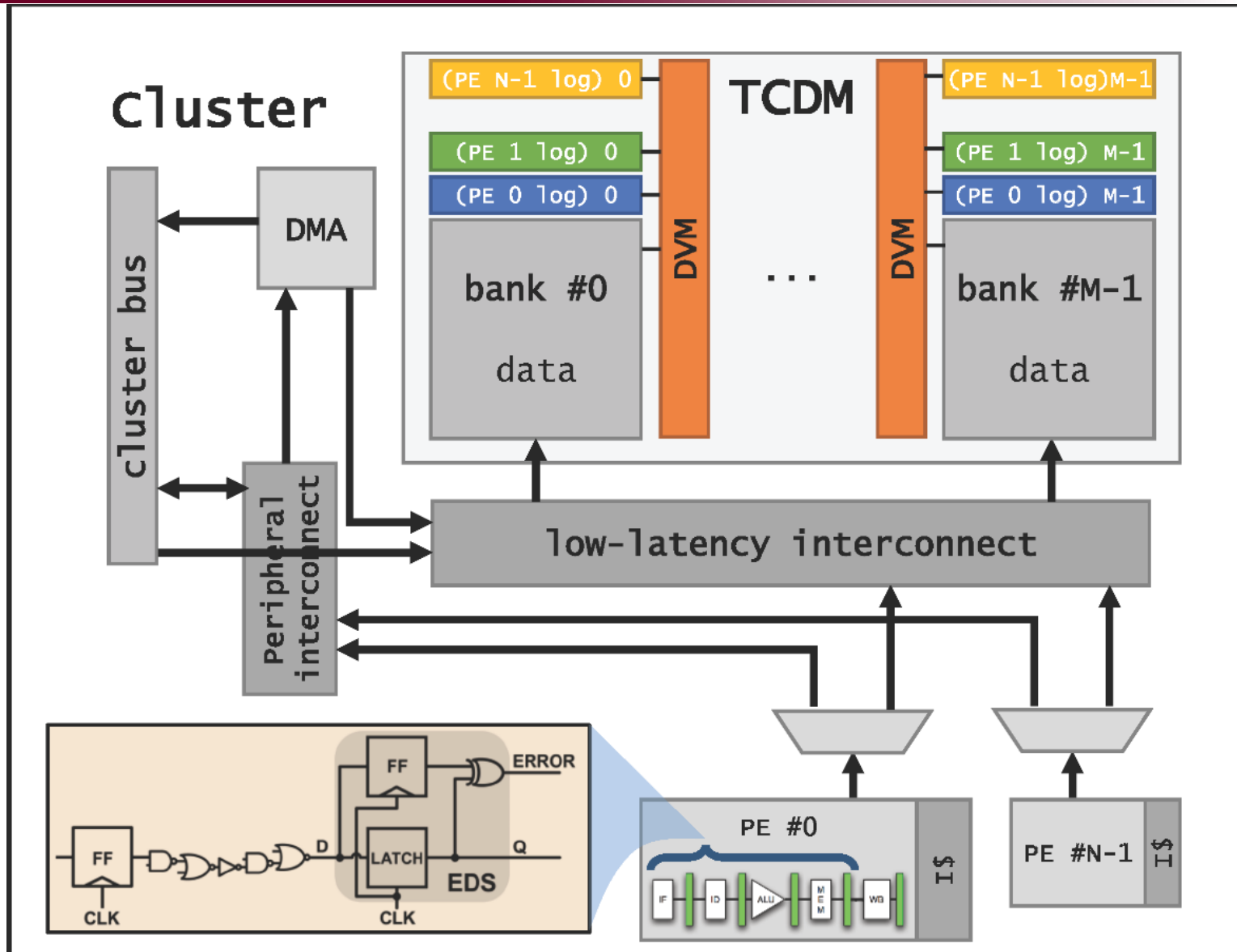


# B) ERROR CORRECTION



- **HTM:** A speculative execution mechanism, used for synchronization in multicore environments.
- Traditionally HTM requires:
  1. Bookkeeping
  2. Data versioning
  3. Checkpointing/Rollback
- **Data Versioning:** Fast distributed logging scheme saves original data.
- **Checkpointing/Rollback:**
  - Transaction boundaries are defined using OpenMP constructs for parallelism.
  - Transaction Starts: A snapshot of the processor state is taken.
  - If error occurs: **Transaction Aborts**. State and Logs are restored.
  - If no error occurs: **Transaction Commits**.

# PULP CLUSTER EXTENDED FOR ERROR TOLERANCE



# EDGE-TM DVS ERROR MANAGEMENT POLICIES



a) **Point of First Failure (PoFF) policy:** Operates just above the edge of failure.

- Lowers voltage until the first failure. Then voltage is step increased.

b) **Thrifty Uncle- Reckless Nephew (TURN) policy:**

- Lowers voltage beyond the point of first failure.
- Adjusts voltage based on aborts and commits.



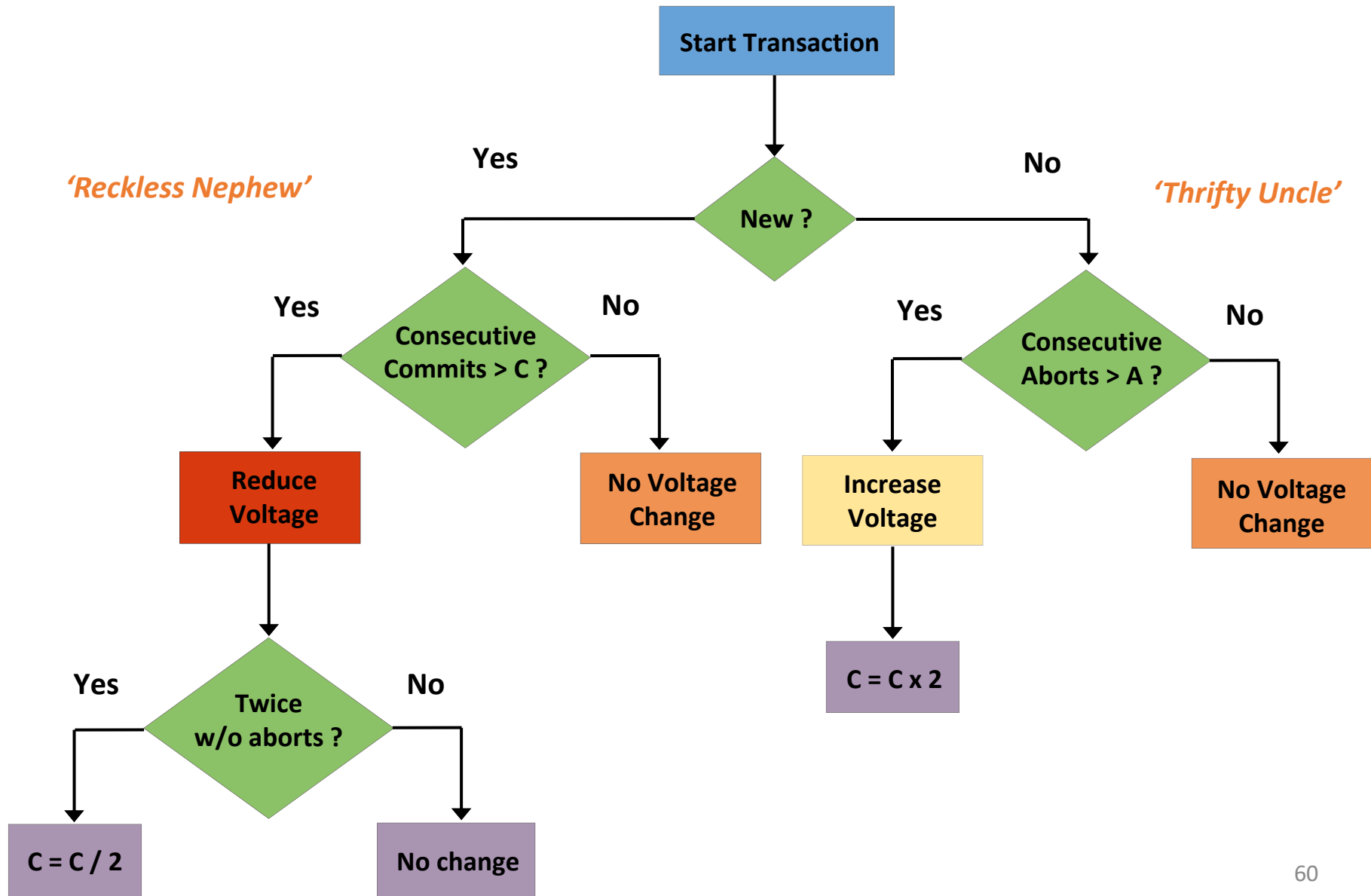
**Trade-off:** Further voltage scaling saves energy but increases abort rate  
➡ Energy and time is wasted in recovery/re-execution.

▪ Two parts:

1) **Reckless Nephew:** Scales down voltage for further energy savings.

2) **Thrifty Uncle:** Moderates energy loss due to aborts by setting up a threshold for voltage scaling based on **number of consecutive aborts: A** and **commits: C**.

# THRIFTY UNCLE/RECKLESS NEPHEW POLICY



# EXPERIMENTAL SETUP

---

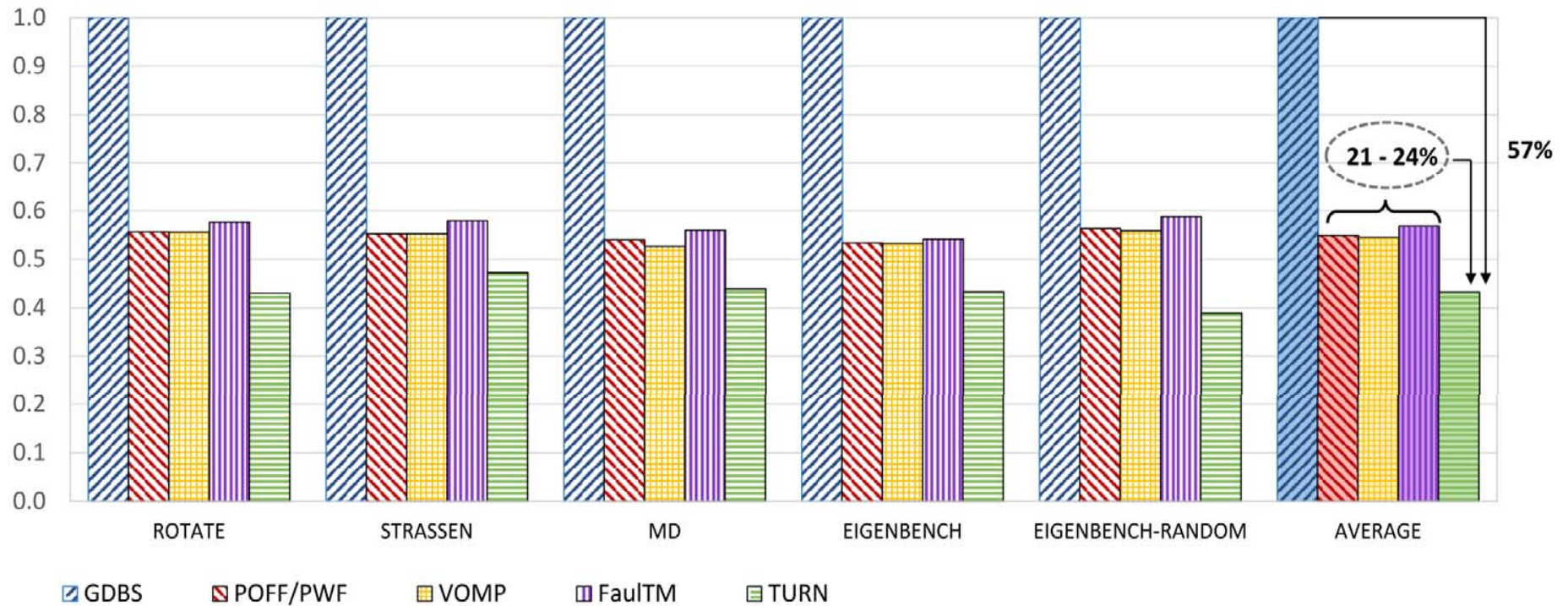


- Power/Performance numbers from STMicroelectronics 28nm UTB FD-SOI implementation of PULP.
- Comparison of Edge-TM policies (**POFF** and **TURN**) with related work:
  - **VOMP**: Vulnerability-aware, error-tolerant task scheduling (Rahimi et al.)
  - **FaultM**: Revisiting transactional memory for fault tolerance (Yalcin et al.)
  - **PWF**: Revisiting transactional memory for timing error-tolerance (Papagiannopoulou et al.)
- **GDBS**: Uses guardbands and keeps voltage steady (Baseline technique).

# RESULTS: ENERGY CONSUMPTION



## Energy Consumption



- **TURN** vs **POFF/PWF/VOMP**: 21%
- **TURN** vs **FAULTM**: 24%
- **TURN** vs **GDBS**: 57%
- **POFF** vs **GDBS**: 45%