# Adaptive Execution on 3D Microprocessors

Koji Inoue

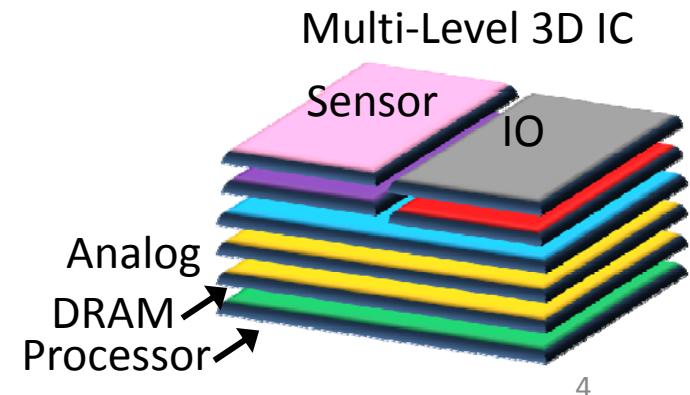Kyushu University

# Outline

- Why 3D?
- Will 3D always work well?
- Support Adaptive Execution!
  - Memory Hierarchy Run-time Optimization
- Conclusions

# Outline

- <span style="color:red">Why 3D?</span>
- Will 3D always work well?
- Support Adaptive Execution!
  - Memory Hierarchy Run-time Optimization
- Conclusions

# From 2D to 3D! (not only TV)

- Stack Multiple Dies
- Connect Dies with Through Silicon Vias

Multi-Level 3D IC

Sensor

IO

Analog

DRAM

Processor

# Chip Implementation Examples from ISSCC'09

- Image Sensors
- SRAM for SoCs
- DRAM
- Multi-core + SRAM connected with wireless TSVs

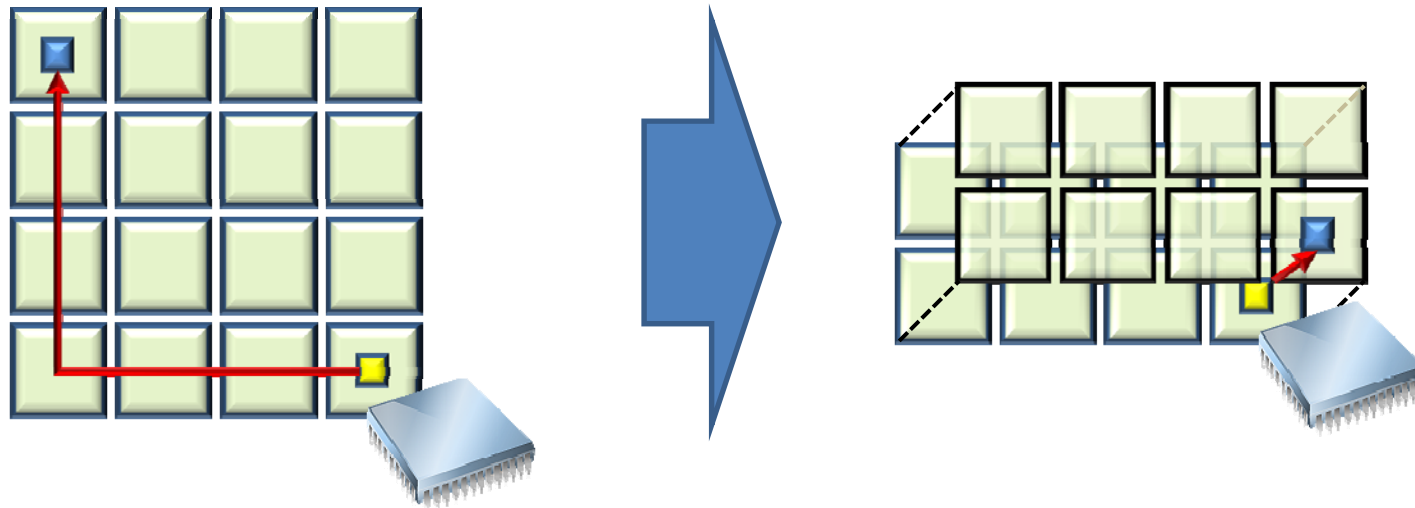U. Kang et al., "8Gb DDR3 DRAM Using Through-Silicon-Via Technology," ISSCC'09.

H. Saito et al., "A Chip-Stacked Memory for On-Chip SRAM-Rich SoCs and Processors, " ISSCC'09.

V. Suntharalingam et al., "A 4-Side Tileable Back Illuminated 3D-Integrated Mpixel CMOS Image Sensor," ISSCC'09.

K. Niitsu et al., "An Inductive-Coupling Link for 3D Integration of a 90nm CMOS Processor and a 65nm CMOS SRAM," ISSCC'09.
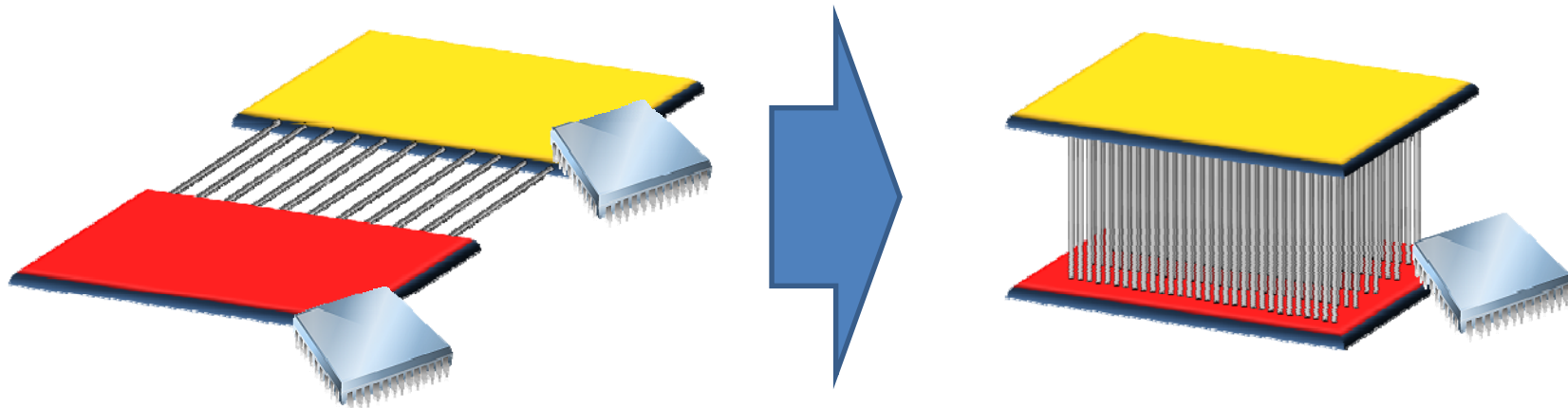
# Why 3D? (1/2)

- Wire Length Reduction
  - Replace long, high capacitance wires by TSVs
  - Low Latency, Low Energy
- Small footprint

# Why 3D? (2/2)

- Integration
  - From "Off-Chip" to "On-Chip"
  - Improved Communication
    - Low Latency, High Bandwidth, and Low Energy
  - Heterogeneous Integration
    - E.g. Emerging Devices

# Outline

- Why 3D?
- <span style="color:red">Will 3D always work well?</span>
- Support Adaptive Execution!
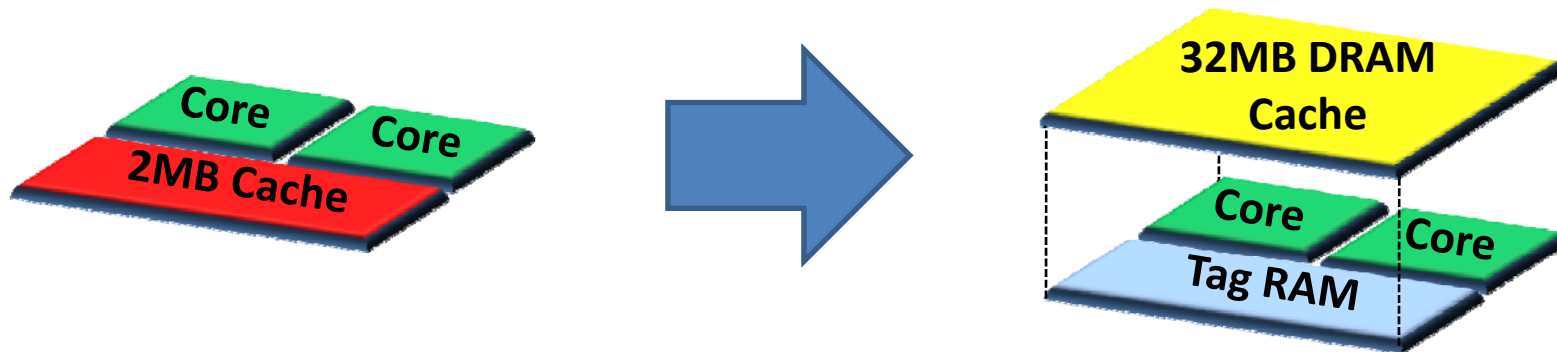  - Memory Hierarchy Run-time Optimization
- Conclusions

# Importance of On-Chip Caches

- Memory-Wall Problem
  - Memory bandwidth does not scale with the # of cores
  - Growing speed gap between processor cores and DRAMs
  - So, Becomes more serious
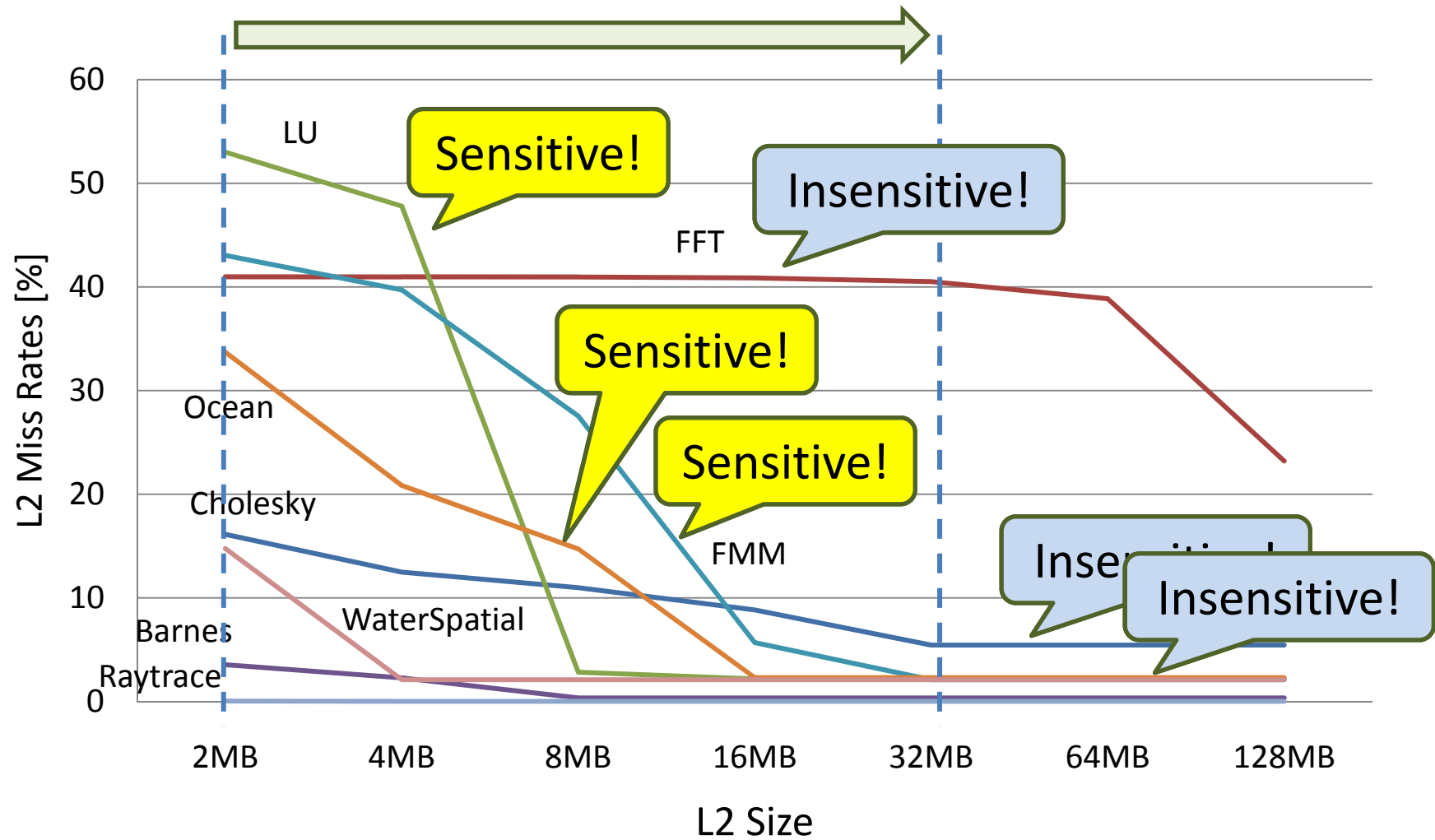- Let's increase on-chip cache capacity, but...
  - Requires large chip area

# Will 3D always work well?
## "Stacking a DRAM Cache"

$$AMAT = HT_{L1} + MR_{L1} \times (HT_{L2} + MR_{L2} \times MMAT)$$

L1 Hit Time, L1 Miss Rate, L2 Hit Time, L2 Miss Rate, Main-Memory Access Time

Ave. Memory Acc. Time

| Impact of DRAM Stacking | → | → | ↗ | ↘ ? | → |
|---|---|---|---|---|---|



Core  Core
2MB Cache

→

32MB DRAM Cache
Core  Core
Tag RAM

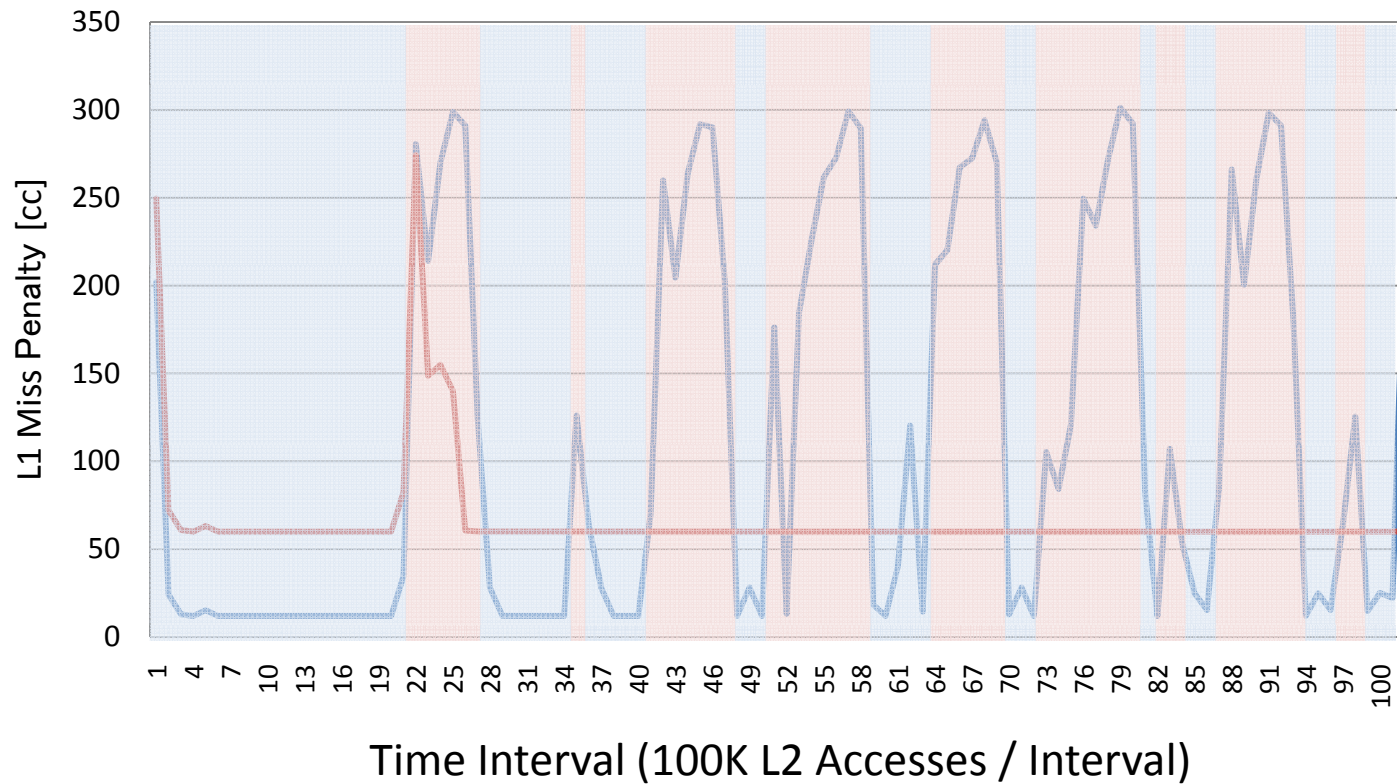# Cache-Size Sensitivity Varies among Programs!

# 2D vs. 3D



$$Profit = \frac{MR_{L2\_REDUCTION} \times MMAT}{HT_{L2\_OVERHEAD}}$$

# Appropriate Cache Size Varies within Programs!



The lower, the better     —— 2MB(12cc)    —— 32MB(60cc)     Ocean

L1 Miss Penalty [cc]

Time Interval (100K L2 Accesses / Interval)

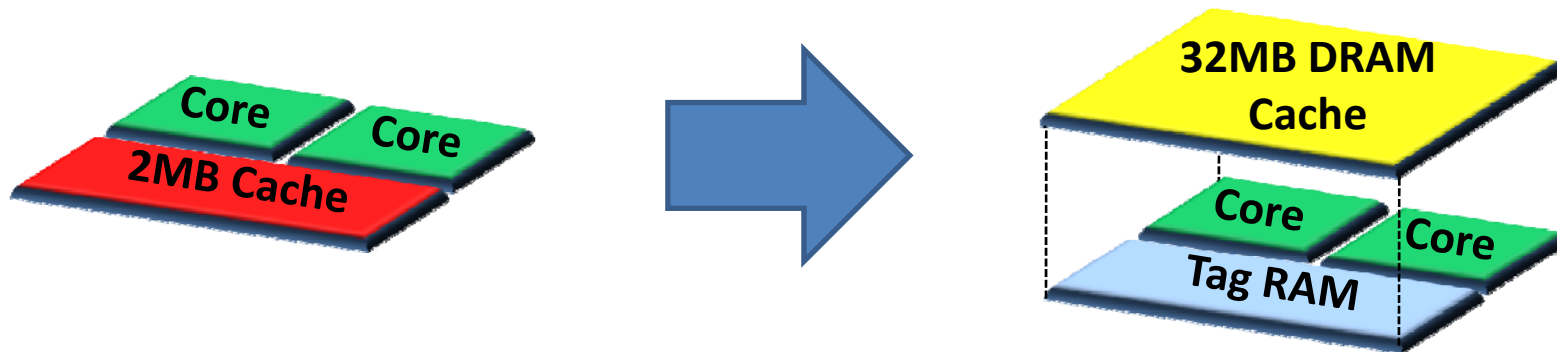# Outline

- Why 3D?
- Will 3D always work well?
- <span style="color:red">Adaptive Execution!</span>
  - <span style="color:red">Memory Hierarchy Run-time Optimization</span>
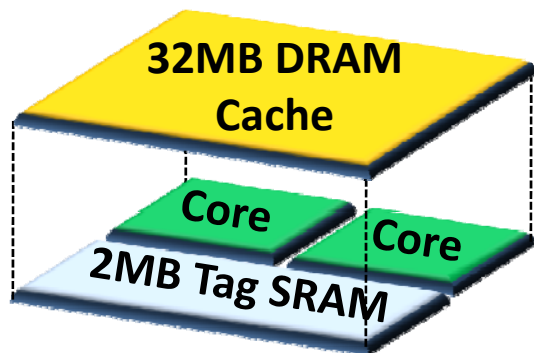- Conclusions

# Will 3D always work well?
# "Stacking a DRAM Cache"

L1 Hit Time  L1 Miss Rate  L2 Hit Time  L2 Miss Rate  Main-Memory Access Time

$$AMAT = HT_{L1} + MR_{L1} \times ( HT_{L2} + MR_{L2} \times MMAT )$$

Ave. Memory Acc. Time

| Impact of DRAM Stacking | → | → | ↗ | ↘ ? | → |
|---|---|---|---|---|---|

Core  Core
2MB Cache
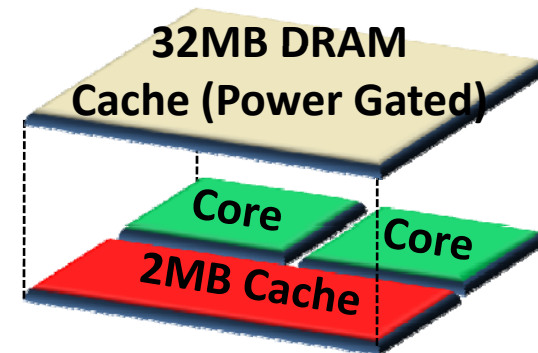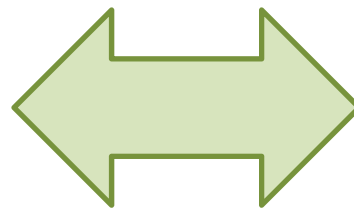
→

32MB DRAM Cache
Core  Core
Tag RAM

# SRAM/DRAM Hybrid Cache Architecture

- Support Two Operation Modes
  - High-Speed, Small Cache Mode (or SRAM Cache Mode)
  - Low-Speed, Large Cache Mode (or DRAM Cache Mode)
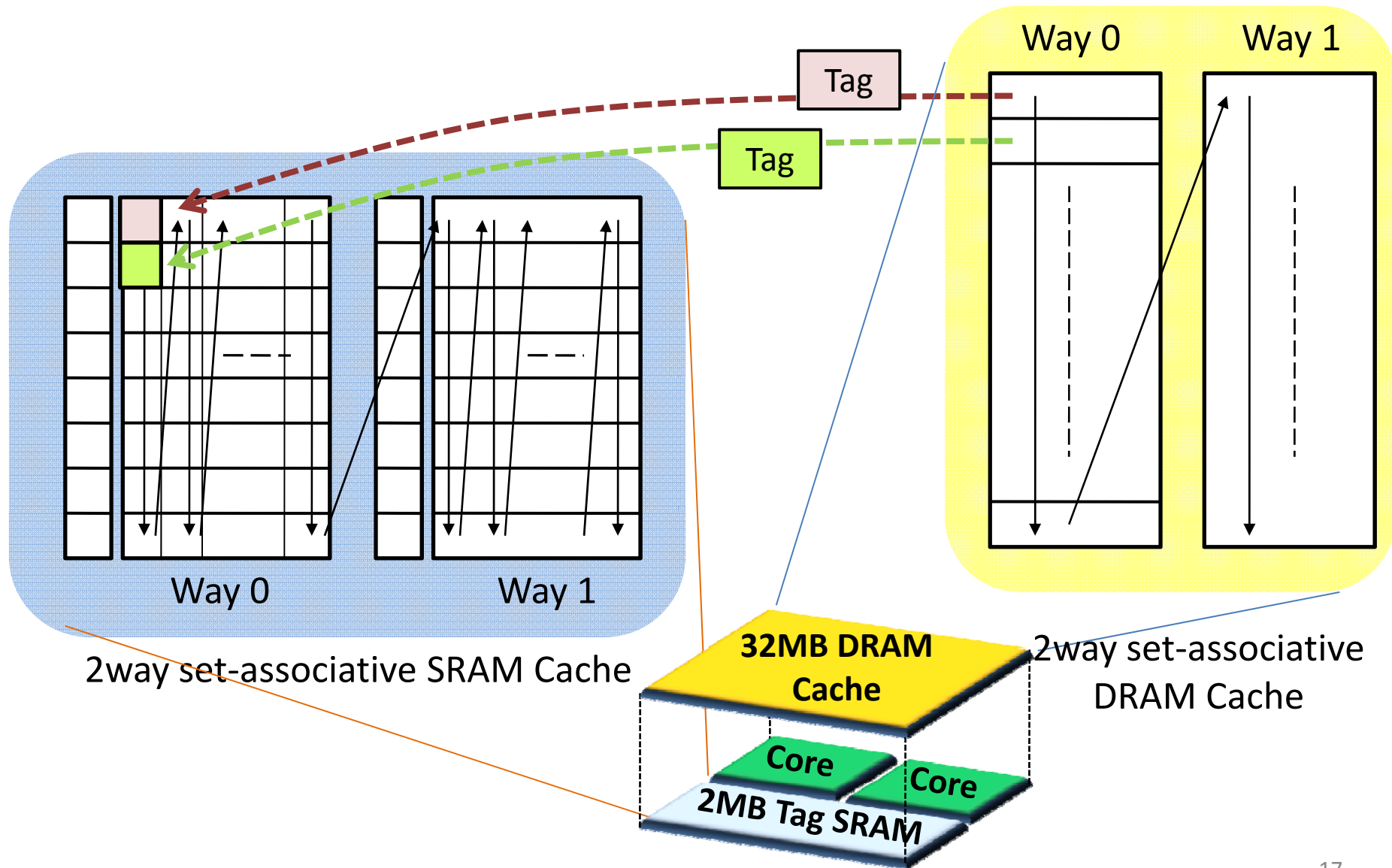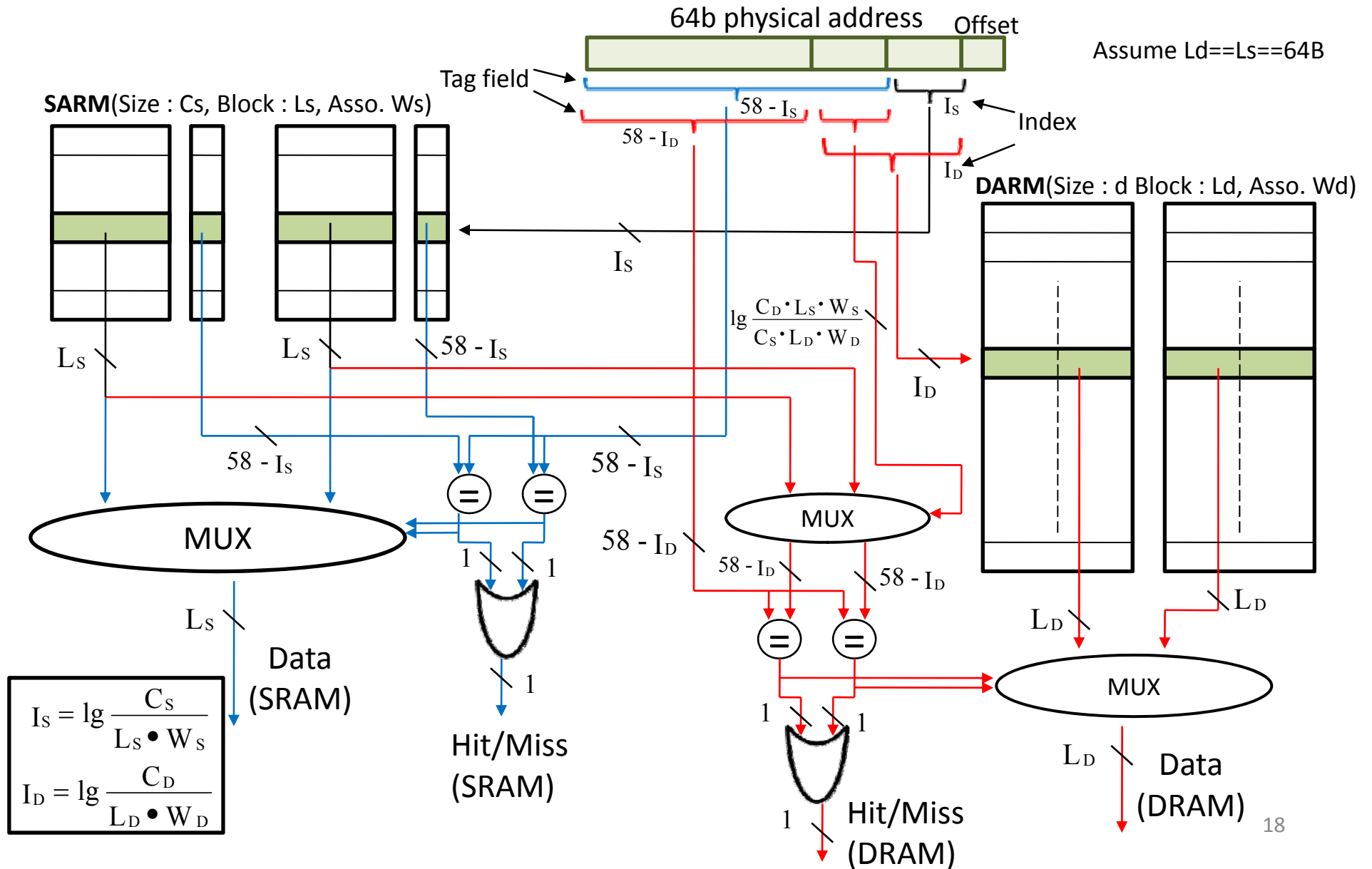- Adapt to variation of application behavior

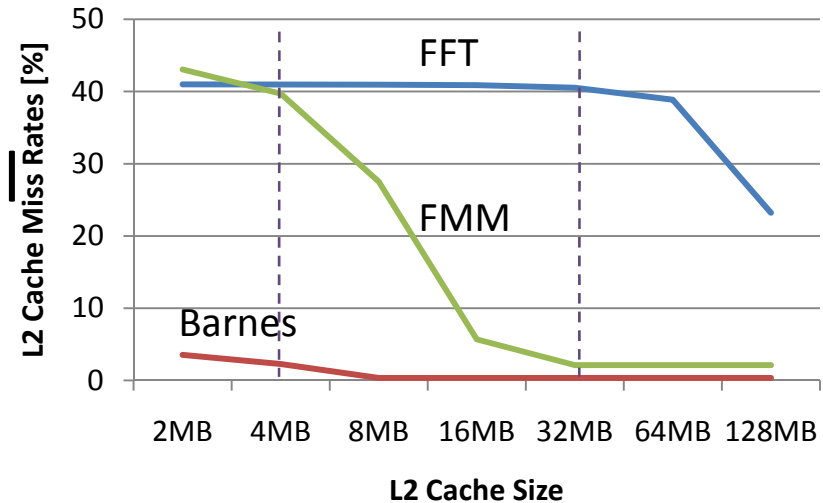DRAM Cache Mode

SRAM Cache Mode

# Microarchitecture (1/2)



Way 0  Way 1

Tag

Tag

Way 0  Way 1

2way set-associative SRAM Cache

**32MB DRAM Cache**

2way set-associative DRAM Cache

**Core**  **Core**

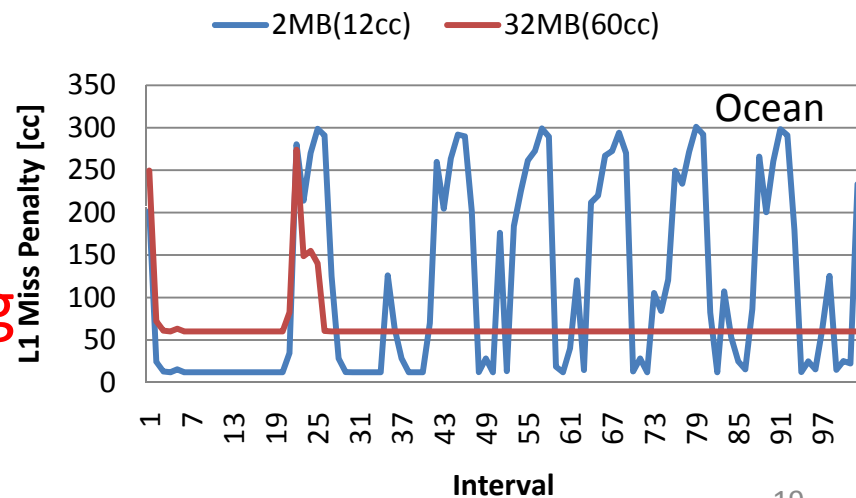**2MB Tag SRAM**

# Microarchitecture (2/2)

# How to Adapt

- ## Static Approach
  - Optimizes at program level
  - Does not change it during execution
  - Needs a static analysis

- ## Dynamic Approach
  - Optimizes at interval level (or phase level)
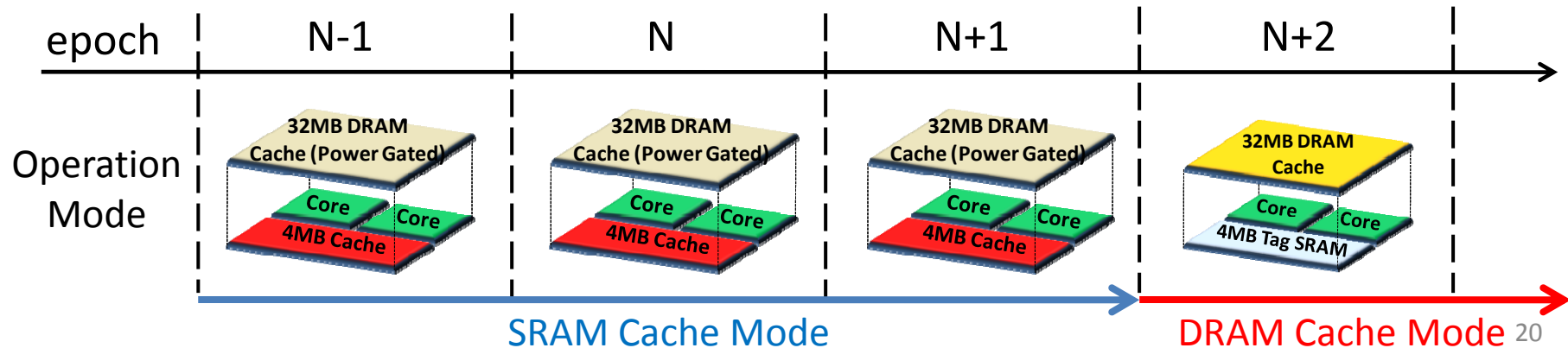  - Needs a run-time profiling

# Run-Time Mode Selection

- Divide Program Execution into "epochs", e.g. 200K L2 Misses

- Predict an Appropriate Operation Mode for Next Epoch

- On SRAM mode, a small tag RAM which stores sampled tags is used to predict DRAM mode miss rates
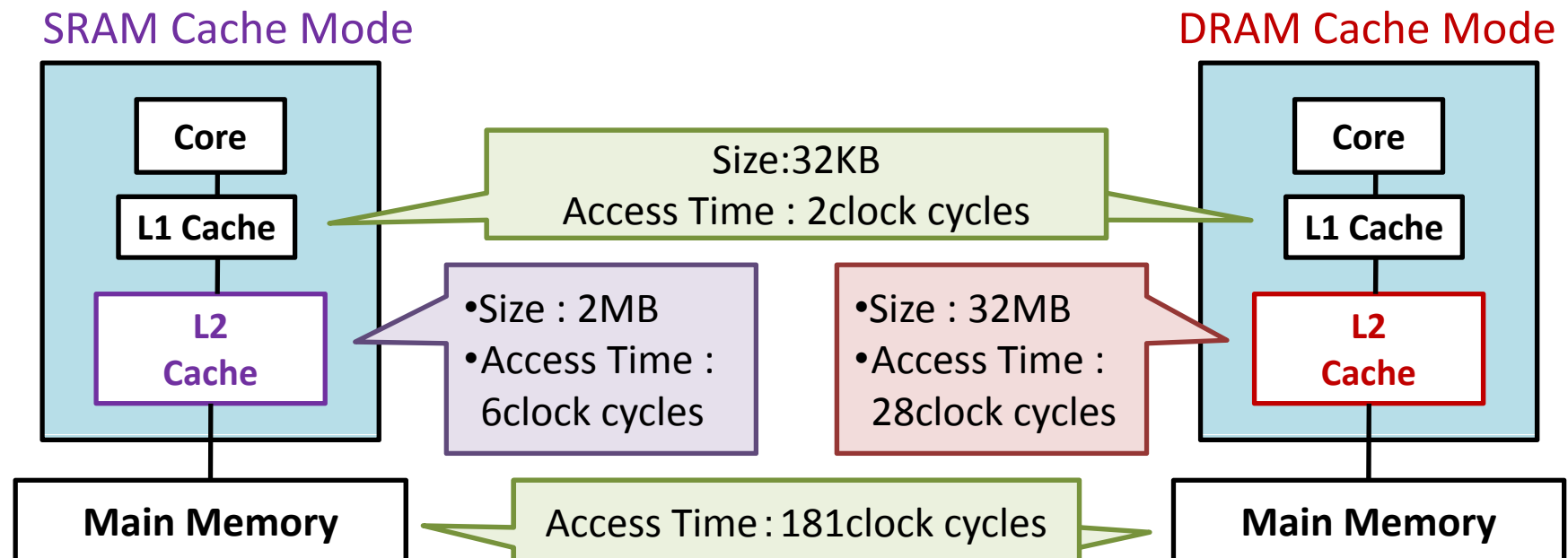
Hardware Support for Measurement

$$if \ MR_{L2SRAM} - MR_{L2DRAM} > \frac{HT_{L2DRAM} - HT_{L2SRAM} + AveOverhead}{MMAT}$$

$$then \ transit \ from \ SRAM \ mode \ to \ DRAM \ mode!$$



epoch · N-1 · N · N+1 · N+2

Operation Mode

32MB DRAM Cache (Power Gated) · Core · Core · 4MB Cache

32MB DRAM Cache (Power Gated) · Core · Core · 4MB Cache

32MB DRAM Cache (Power Gated) · Core · Core · 4MB Cache

32MB DRAM Cache · Core · Core · 4MB Tag SRAM
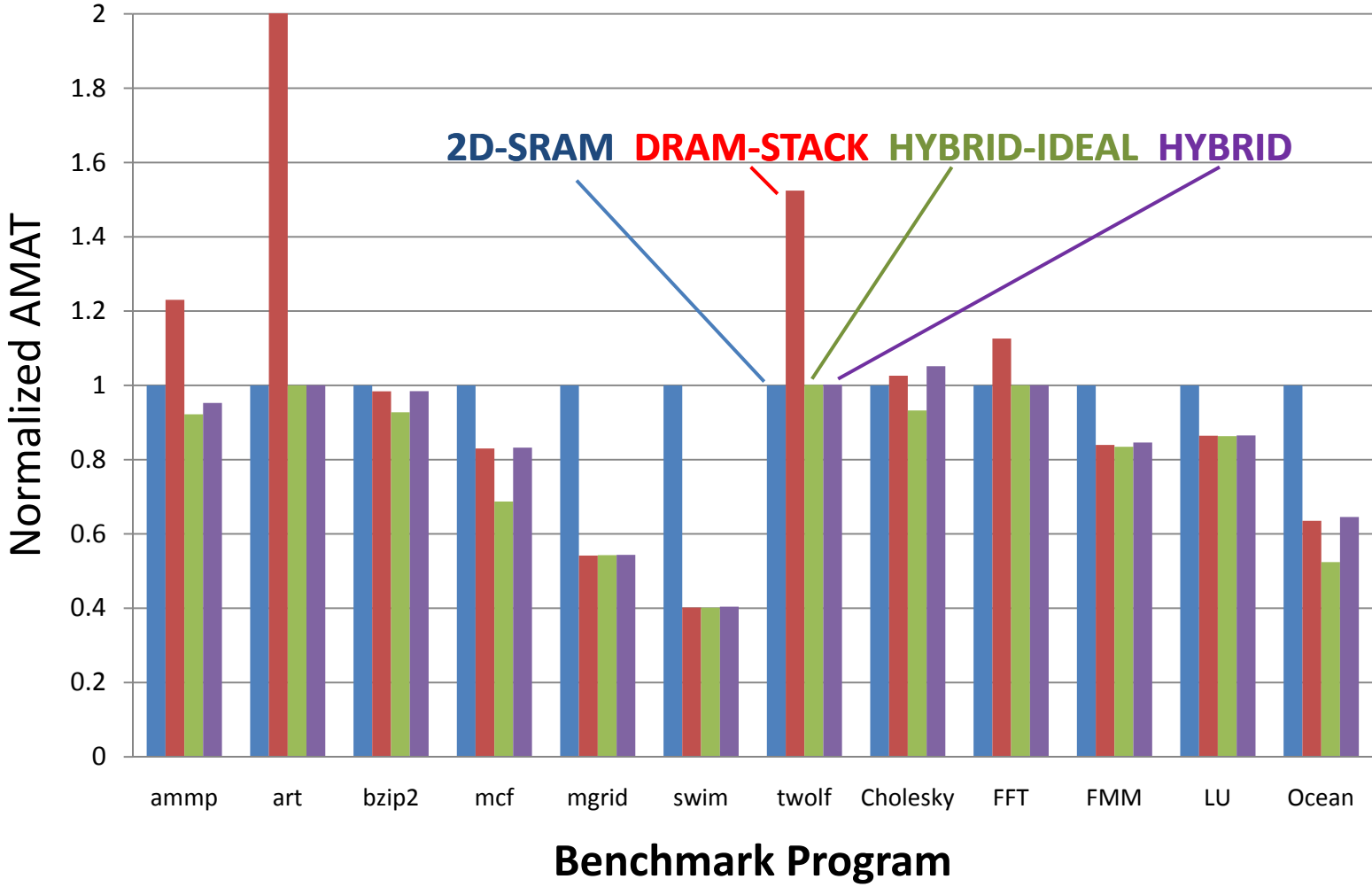
SRAM Cache Mode

DRAM Cache Mode

# Experimental Set Up

- Processor: In-Order

- Benchmarks: SPEC CPU 2000, Splash2



SRAM Cache Mode

DRAM Cache Mode

Core

L1 Cache

L2 Cache

Size:32KB
Access Time : 2clock cycles

- Size : 2MB
- Access Time : 6clock cycles

- Size : 32MB
- Access Time : 28clock cycles

Core

L1 Cache

L2 Cache

Main Memory

Access Time : 181clock cycles
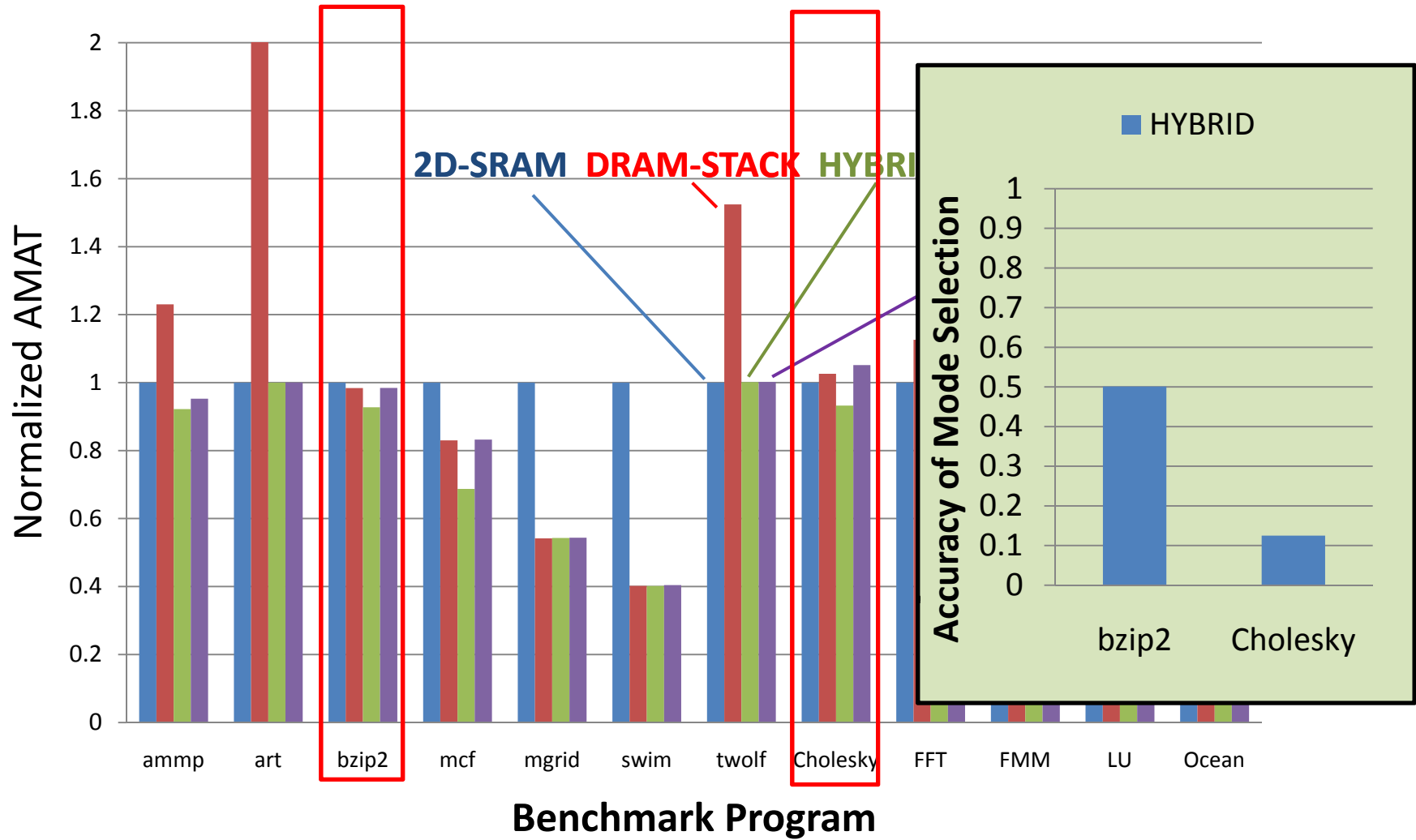
Main Memory

# Results

# Results

# Results

# Conclusions

- The 3D solution is one of the most promising ways to achieve…

  – High performance

  – Low energy

- It does not ALWAYS work well!

- Run-time adaptive execution by considering memory access behavior

# Acknowledgement