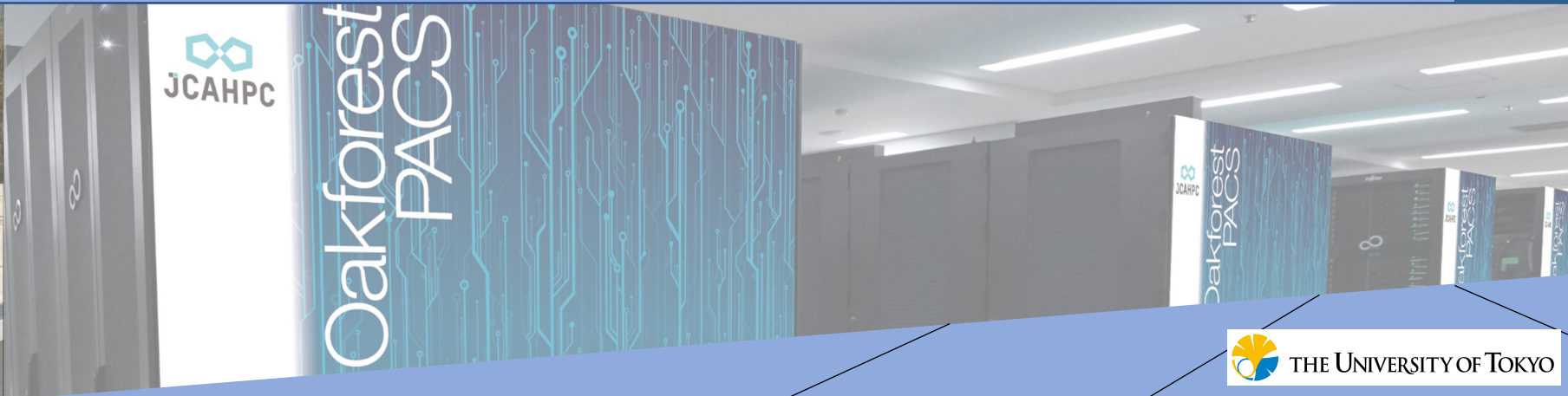


# *SW/HW Optimizations for Next- Generation Supercomputing*

*Eishi Arima Ph.D.*



# Research landscape & our focus

**Goal:** to improve performance and energy efficiency of current/future supercomputing systems

**Approaches:** from microarchitecture to programming models; esp. interactions among them

**Key words:** power/data managements; heterogeneity; new HW devices

## Objectives

Energy efficiency  
Performance  
Resilience  
Productivity  
Security  
Cost  
Predictability

X

## System Stack

Application
Algorithm
Languages
System Software
SW/HW Interface
Microarchitecture
Circuits
Devices

X

## Technology Trends

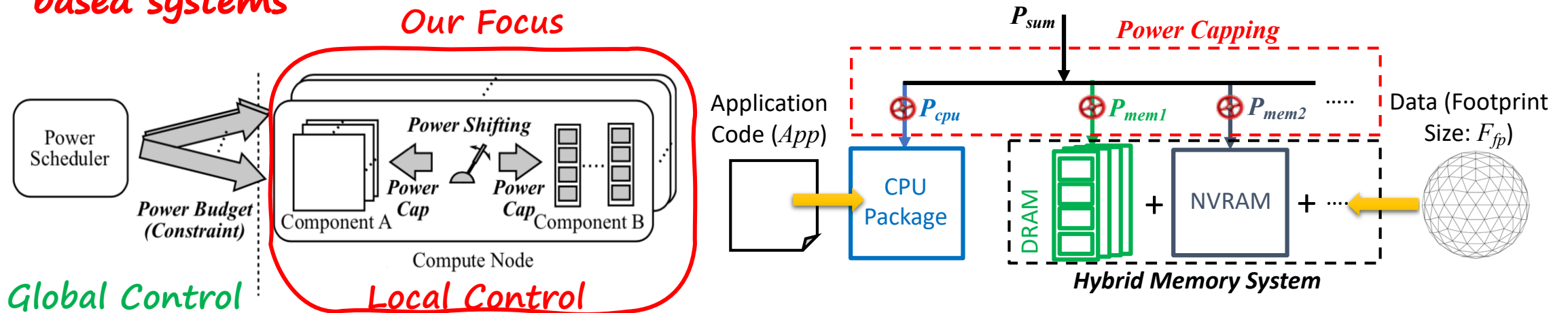
Big data AI  
AMR Par. in time  
Container Portable API  
Power wall Data centric  
ARM Heterogeneity FP16  
NVRAMs 3D stacking  
Quantum Optical NW

# Footprint-Aware Power Capping [E. Arima+ ISC'20]

Background: Hierarchical power management will be an indispensable feature for future power-constrained supercomputers

- Global control: Power scheduler distributes power budgets accordingly to all compute nodes
- Local control: Each compute node attempts to exploit the allocated power budget

Our insight: Optimal power allocations to components can change considerably when the problem size (= **data footprint size**) is scaled esp. on **hybrid memory based systems**





# Footprint-Aware Power Capping [E. Arima+ ISC'20]

**1st Contribution:** problem formulation & performance modeling

**2nd Contribution:** a profile-based software framework (efficient calibration method + power allocation algorithm)

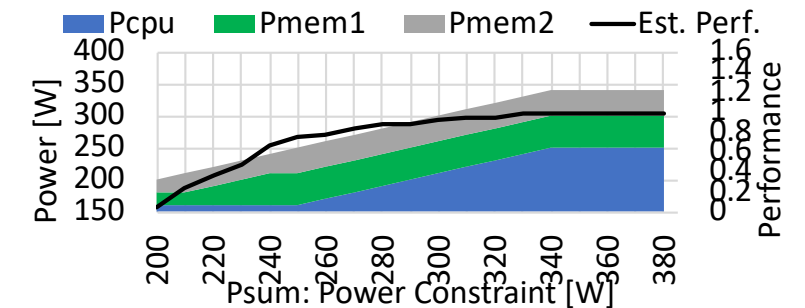
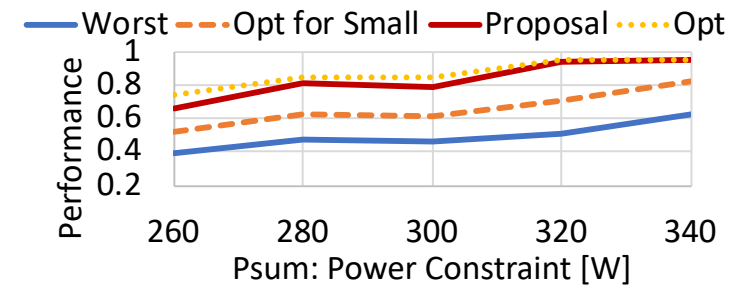
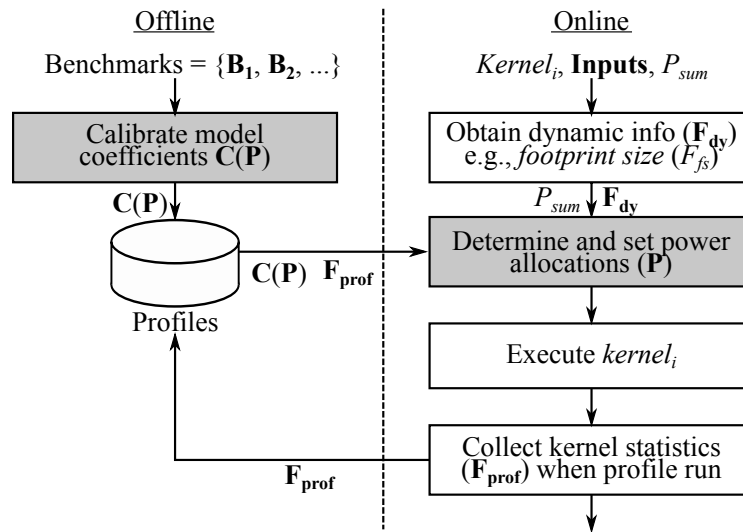
**Result:** near optimal allocations acquired!

Given  $\text{Kernel}, \text{Inputs}, P_{\text{sum}} (\Rightarrow \mathbf{F}, P_{\text{sum}})$   
 Max  $\text{Obj}(\mathbf{P}, \mathbf{F})$   
 s.t.  $\sum P_x \leq P_{\text{sum}}$   
 $P_x \in S_{P_x} (x = \text{cpu}, \text{mem1}, \dots)$

**Kernel:** target kernel region  
**Inputs:** inputs for the app = (arg1, arg2, ...)  
**F:** app feature parameters (e.g., F/B rate)

**Obj(P, F):** objective function

$P_{\text{sum}}$ : given total power budget [W]  
 $\mathbf{P}$ : power allocations =  $(P_{\text{cpu}}, P_{\text{mem1}}, \dots)$   
 $S_{P_x}$ : Set of power caps for component  $x$   
 (e.g.,  $=\{20, 30, 40\}$ )





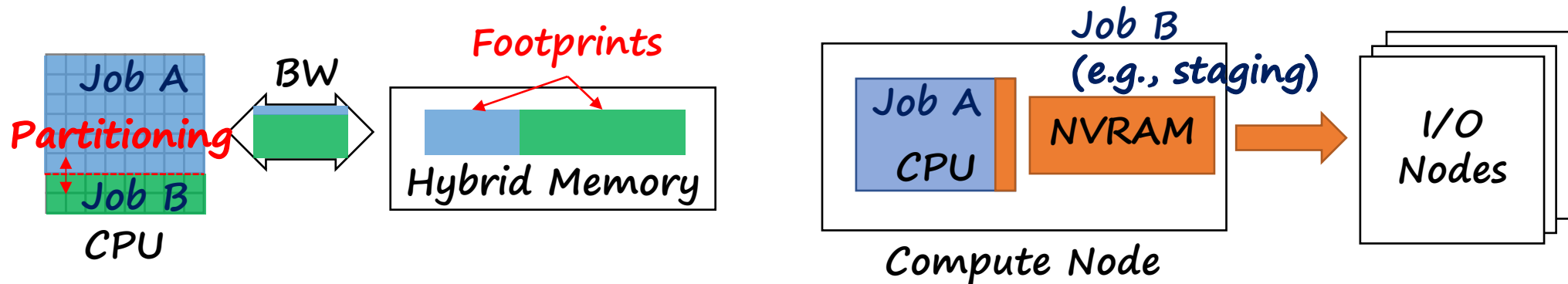
# Next Step: Scheduler-Based Approaches (Ongoing Work)

**Objective:** developing job/power scheduling methods using the **insights** gained through the **power capping work**

- Footprint awareness; performance prediction; system design
- Supported by JSPS grant-in-aid

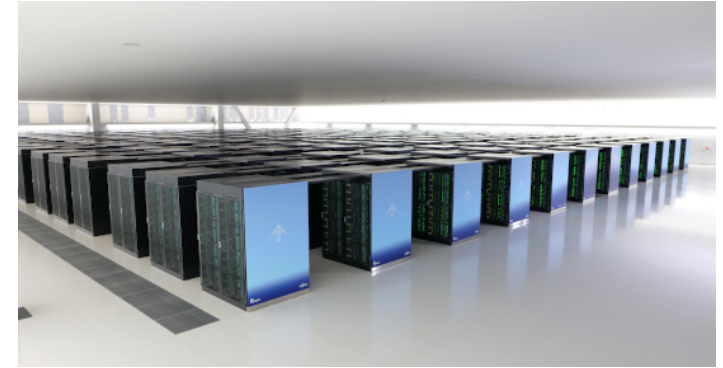
**Our focus:** co-scheduling, i.e., co-running multiple jobs within a node to improve total job throughput or energy efficiency

- Footprint-aware co-scheduling [ICPP'19 poster] – ongoing
- I/O-aware co-scheduling [Not appeared yet] – ongoing

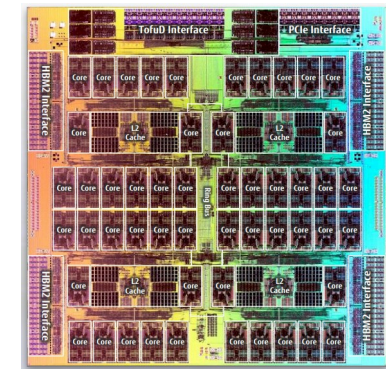


# Power Evaluation using Supercomputer Fugaku

- Evaluation of Power Controls on Supercomputer Fugaku [Y. Kodama+ EEHPC@CLUSTER'20]
  - Large-scale power evaluation using **over 20K nodes** of the Fugaku system
  - Particularly, focusing on some power control features implemented on A64FX processor
- Evaluation for future HPC microprocessors
  - by extending existing performance/power simulators used for the development of the Fugaku system (ongoing work)



Supercomputer  
Fugaku†



A64FX\*

†[https://www.riken.jp/en/news\\_pubs/news/2020/20200623\\_1/](https://www.riken.jp/en/news_pubs/news/2020/20200623_1/)

\*<https://www.nextplatform.com/2019/11/13/a64fx-arm-chip-gets-a-big-push-from-cray/>

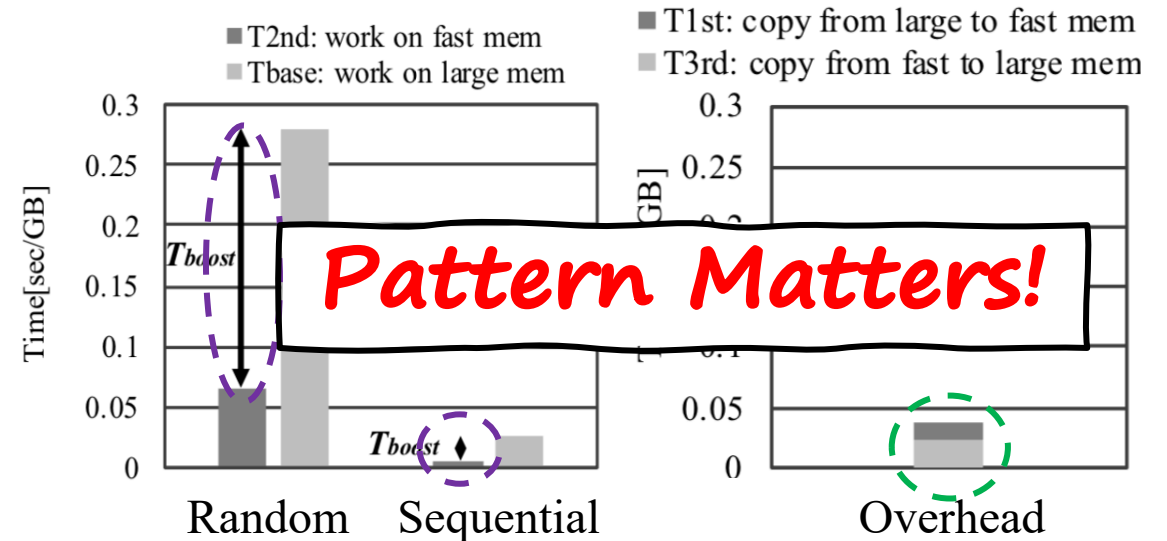
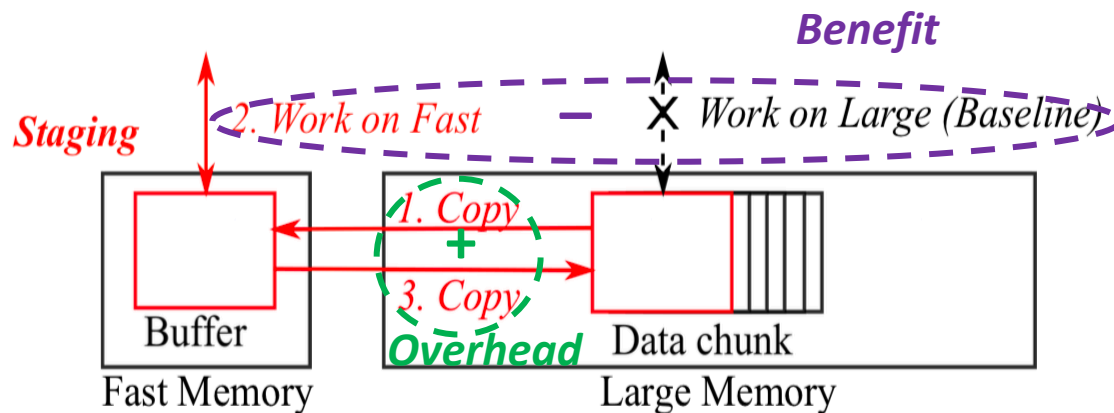
# Pattern-Aware Staging [E. Arima+ISC'20]

**Open question:** How can we exploit both the high BW and the large capacity on a hybrid memory system?

- E.g., HBM + DDR4 or DRAM + NVRAM

**Observation:** **The access pattern matters** when we decide whether or not move a chunk of data

- **Performance benefit** v.s. **data movement overhead**





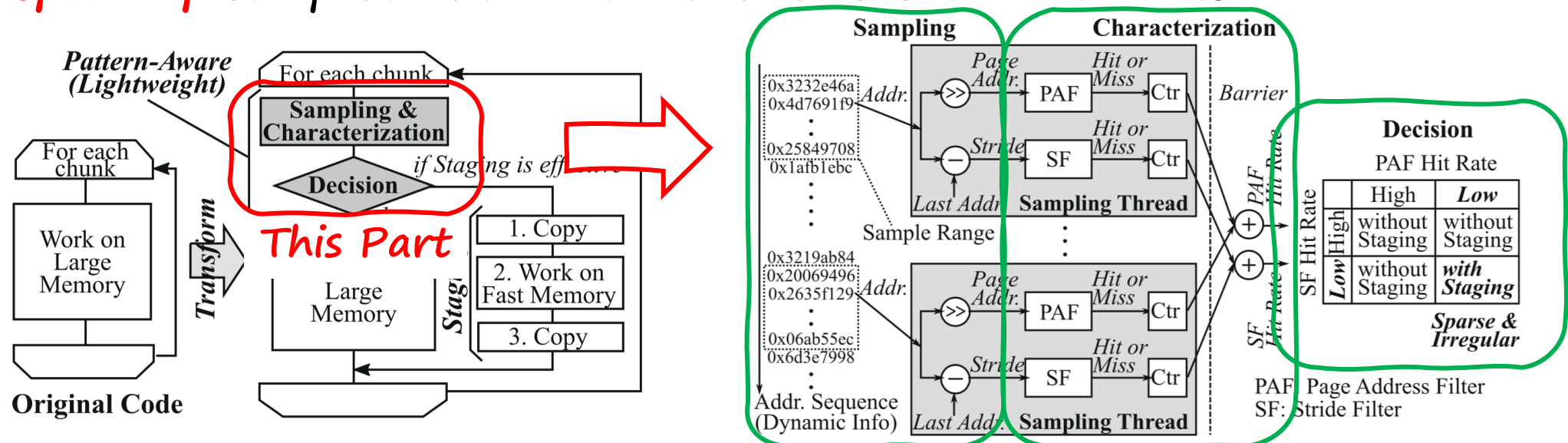
# Pattern-Aware Staging

## [E. Arima+ISC'20]

**Pattern-aware staging:** a lightweight software mechanism that can be ultimately automated by source-to-source compilers

1. **Sampling:** helper threading inspired dynamic address sequence sampling
2. **Characterization:** Bloom filter based runtime access pattern detection
3. **Decision:** criterium to determine whether or not move a chunk of data

**Evaluation results:** **300% speed-up** compared with the large memory only and **41% speed-up** compared with the hardware cache mode at best



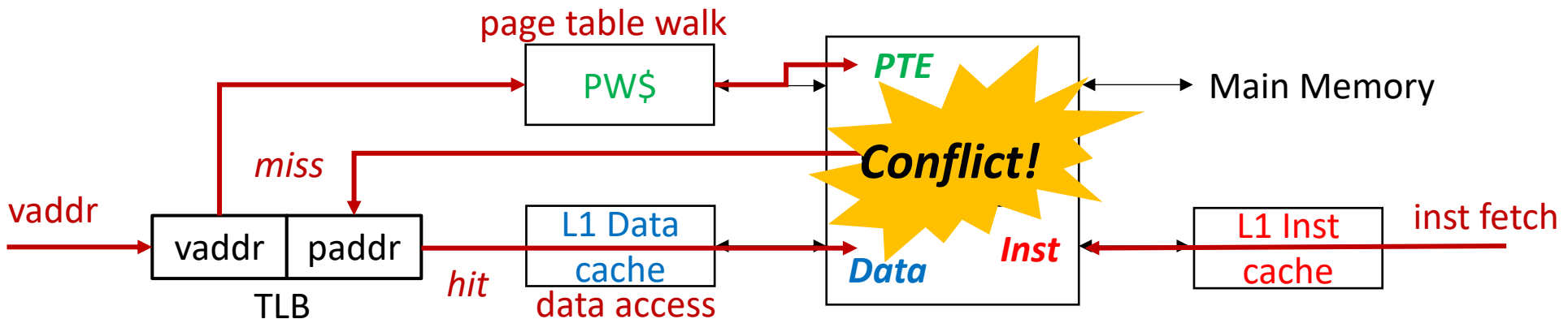
# Advanced HW Cache Management

## [E. Arima DSD'20]

Background: significant **cache conflicts** among different classes of cachelines on lower level caches

- Different classes: data, instructions (codes), and PTEs
- Even worse because: (1) PTE accesses are becoming critical on modern systems esp. for virtual machines; (2) code footprints are increasing esp. for modern server workloads

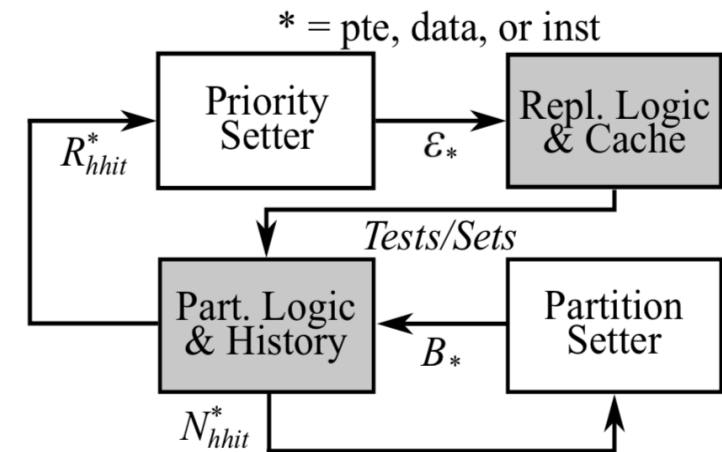
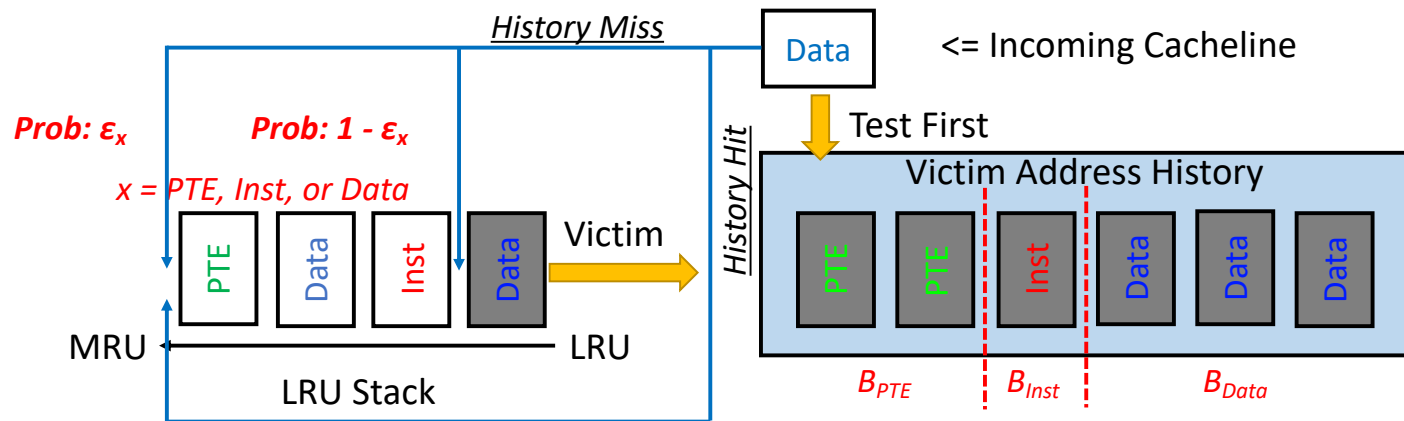
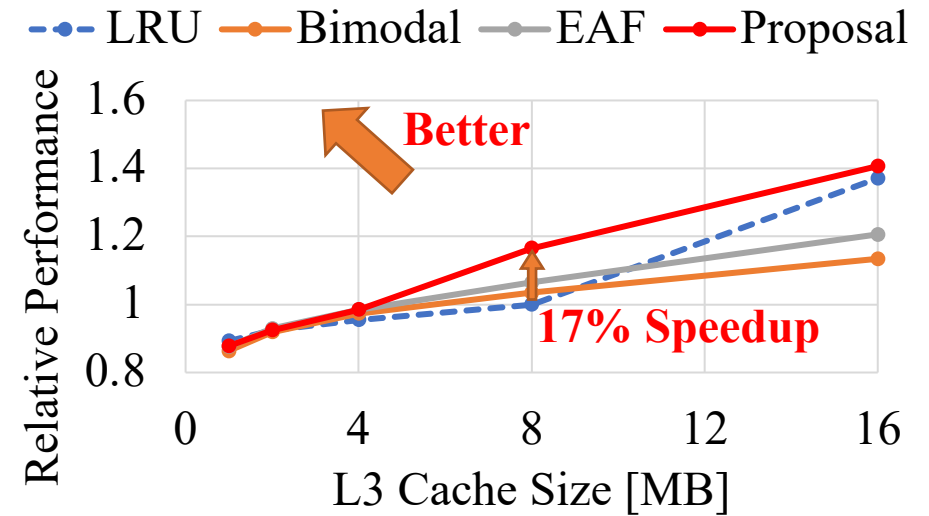
**We should set appropriate allocation priorities to these classes!**



# Advanced HW Cache Management [E. Arima DSD'20]

**Approach:** (1) augment a functionality to set a unique priority to each class;  
(2) partition a module used for estimating reuseness behavior (history of victims);  
(3) provide a control system to optimize the **priorities/partitions at the same time**

**Evaluation result:** considerable speed-up over existing approaches





# Conclusion

**Summary:** We are aiming to provide various **insights** that can potentially lead to significant performance, energy efficiency, or any other sorts of improvement for current/future supercomputing systems through both SW & HW based approaches.

## Related Publications:

- E. Arima, T. Hanawa, C. Trinitis, M. Schulz "Footprint-Aware Power Capping for Hybrid Memory Based Systems" In *Proc. of ISC High Performance*, pp.347–369, (2020) [Youtube](#)
- Y. Kodama, T. Odajima, E. Arima, and M. Sato "Evaluation of Power Controls on Supercomputer Fugaku" In *Proc. of CLUSTER (EEHPC volume)*, pp.xx–xx, (2020) [Youtube](#)
- E. Arima, M. Schulz "Pattern-Aware Staging for Hybrid Memory Systems" In *Proc. of ISC High Performance*, pp.474–495, (2020) [Youtube](#)
- E. Arima "Classification-Based Unified Cache Replacement via Partitioned Victim Address History" In *Proc. of DSD*, pp.101–108, (2020)

🌐 **Website:** <https://www.cspp.cc.u-tokyo.ac.jp/arima/index-e.html>

✉ **E-mail:** [arima@cc.u-tokyo.ac.jp](mailto:arima@cc.u-tokyo.ac.jp)



*That's it! Thank you for watching!*

## Staff

Director: Eishi Arima

Speaker: Eishi Arima

Investigator: Eishi Arima

Editor: Eishi Arima

## Grants

JSPS Grant-in-Aid for Early-Career Scientists (JP20K19766) – PI: EA

JSPS Grant-in-Aid for Early-Career Scientists (JP18K18021) – PI: EA

JSPS Grant-in-Aid for Research Activity Start-up (JP16H06677) – PI: EA

Research on Processor Architecture, Power Management, System Software  
and Numerical Libraries for the Post K Computer System of RIKEN

## Special Thanks to

Folks at ITC, U. Tokyo; CAPS, TUM; R-CCS, RIKEN; and CASC, LLNL



*The End*