# System Sotware Research (2)
## Towards Next-Generation Supercomputing

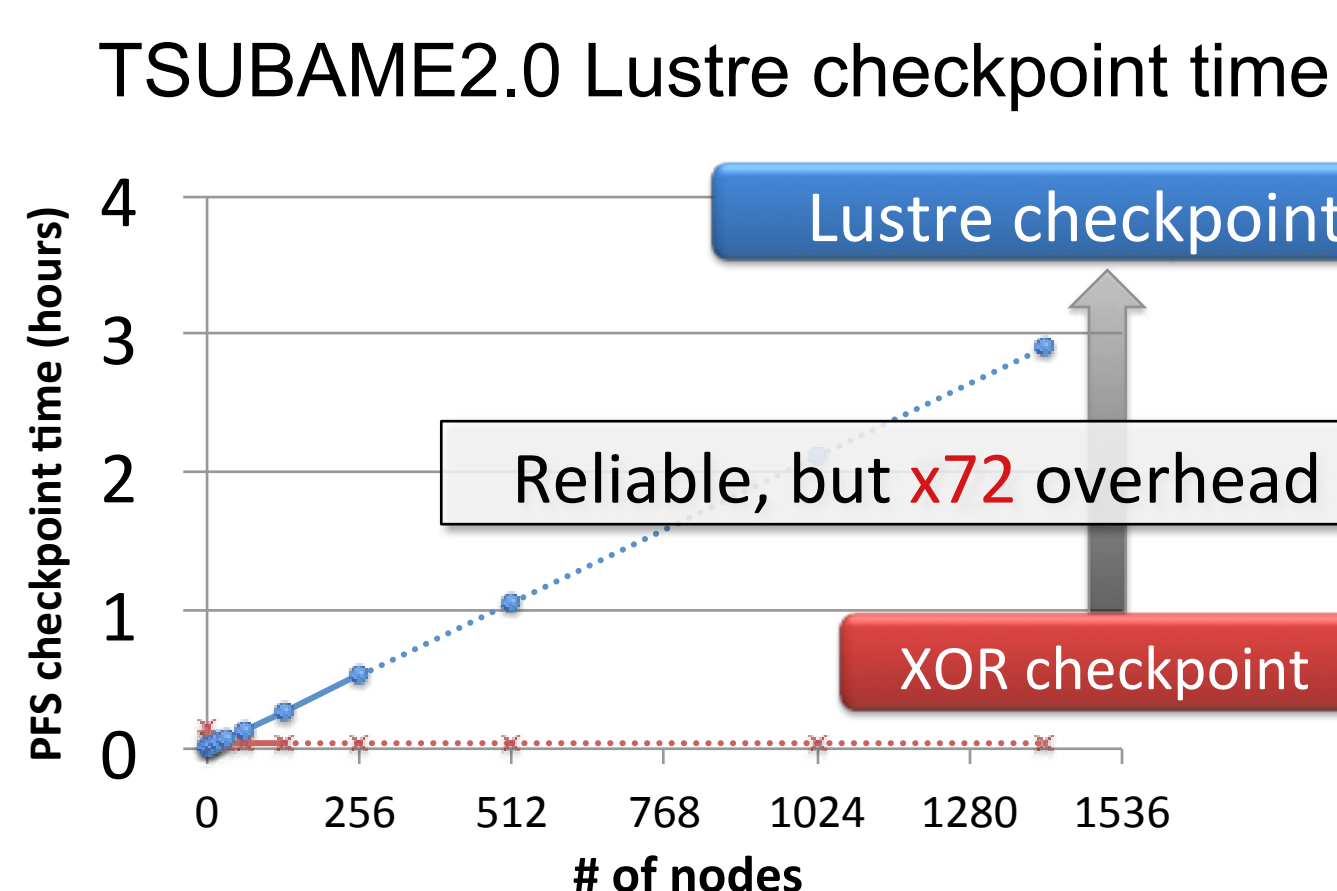# Fault tolerant Infrastructure for Billion-Way Parallelization

## Asynchronous Checkpointing System

Background: Increasing system failures

- A node failure occurred every  13 hours  on average
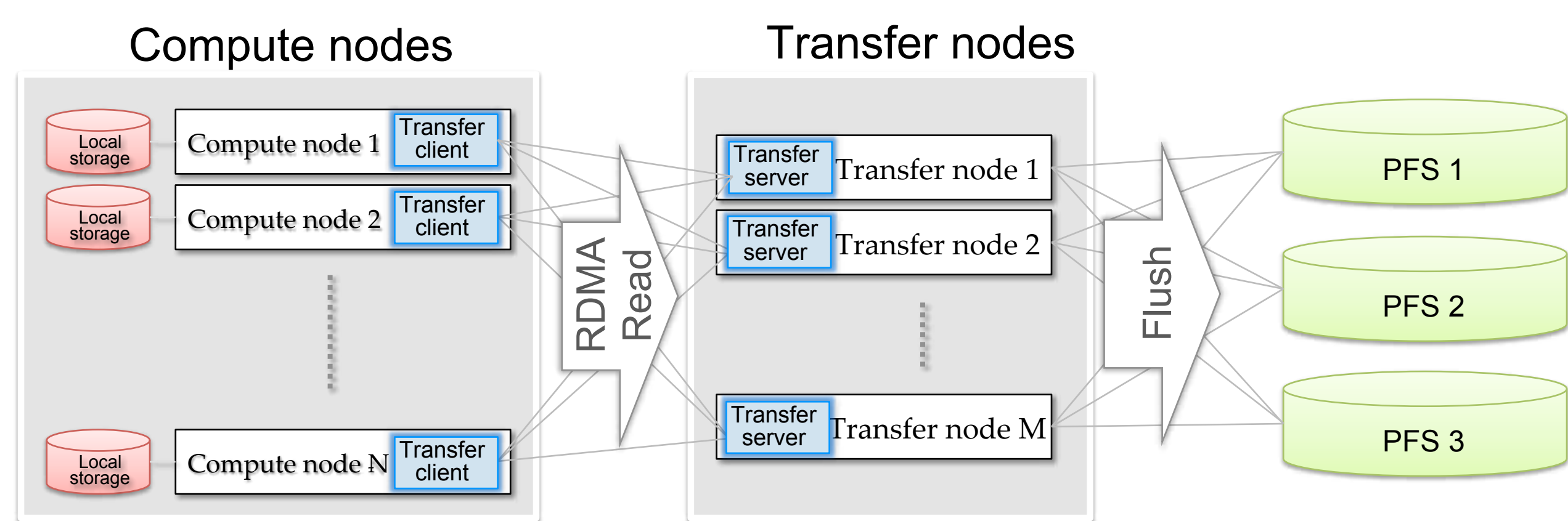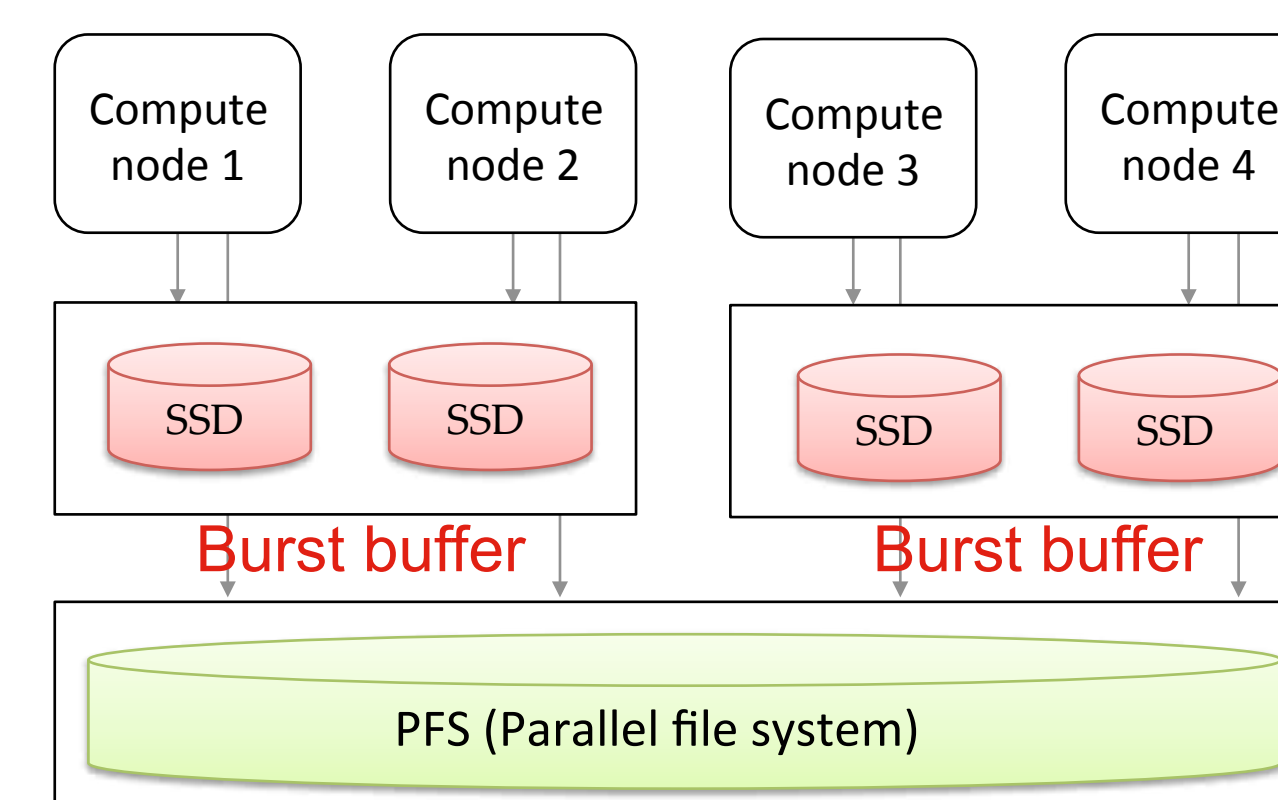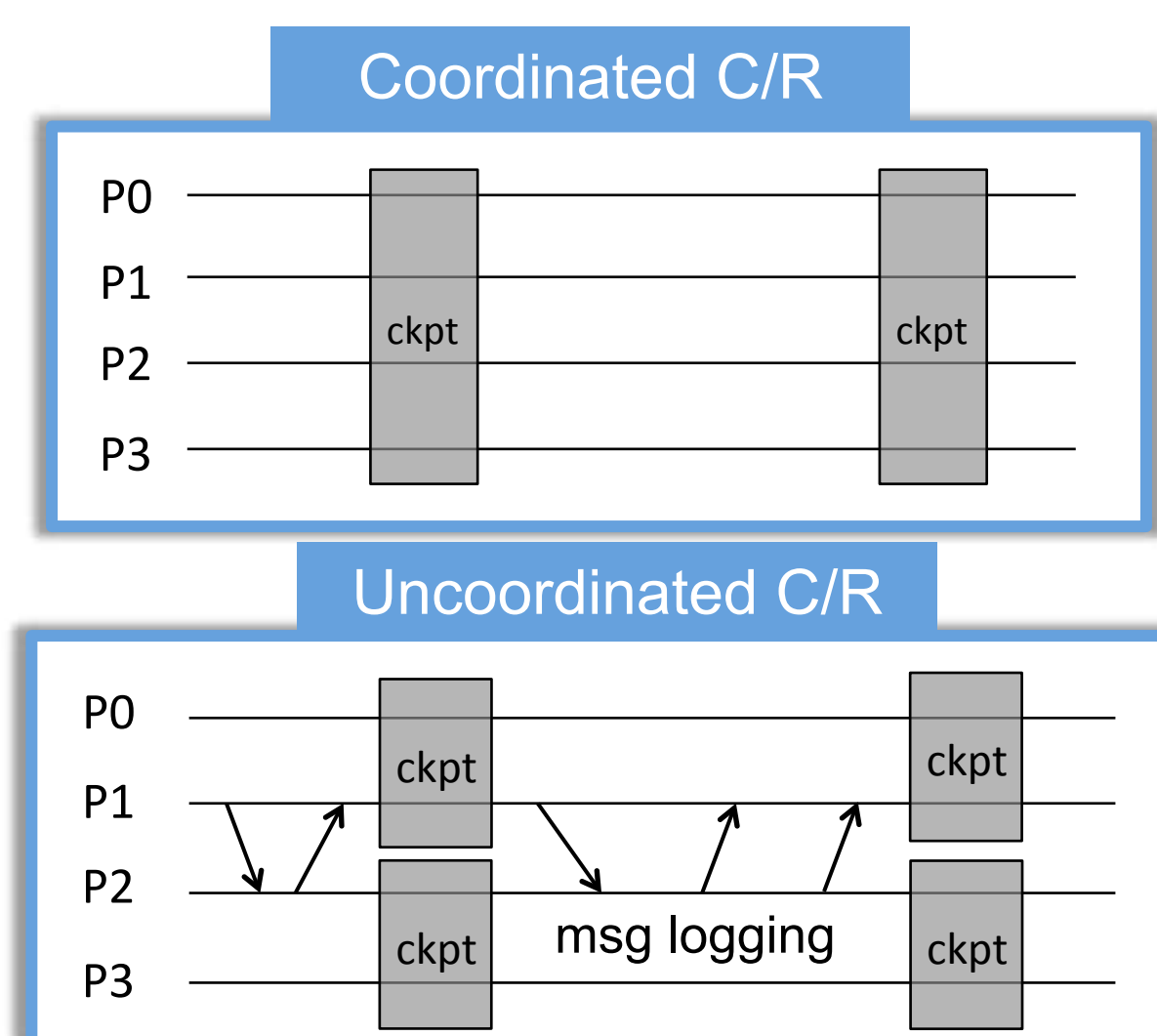- A parallel file system (PFS) checkpointing overhead ( 3 hours )

Objective : Reduce PFS checkpoint overhead

Proposed method : Implementation and modeling of an asynchronous checkpointing

TSUBAME2.0 Lustre checkpoint time



Reliable, but x72 overhead

Lustre checkpoint
XOR checkpoint

PFS checkpoint time (hours) / # of nodes

Compute nodes          Transfer nodes



Local storage — Compute node 1 — Transfer client
Local storage — Compute node 2 — Transfer client
Local storage — Compute node N — Transfer client
RDMA Read
Transfer server — Transfer node 1
Transfer server — Transfer node 2
Transfer server — Transfer node M
Flush
PFS 1
PFS 2
PFS 3

x1.1 ~ x1.8 System efficiency improvement

$$Efficiency = \frac{ideal\_runtime}{expected\_runtime}$$



Efficiency / Failure rate x1 / Failure rate x2 / Failure rate x10

- PFS cost x1 / Asynchronous
- PFS cost x1 / Synchronous
- PFS cost x2 / Asynchronous
- PFS cost x2 / Synchronous
- PFS cost x10 / Asynchronous
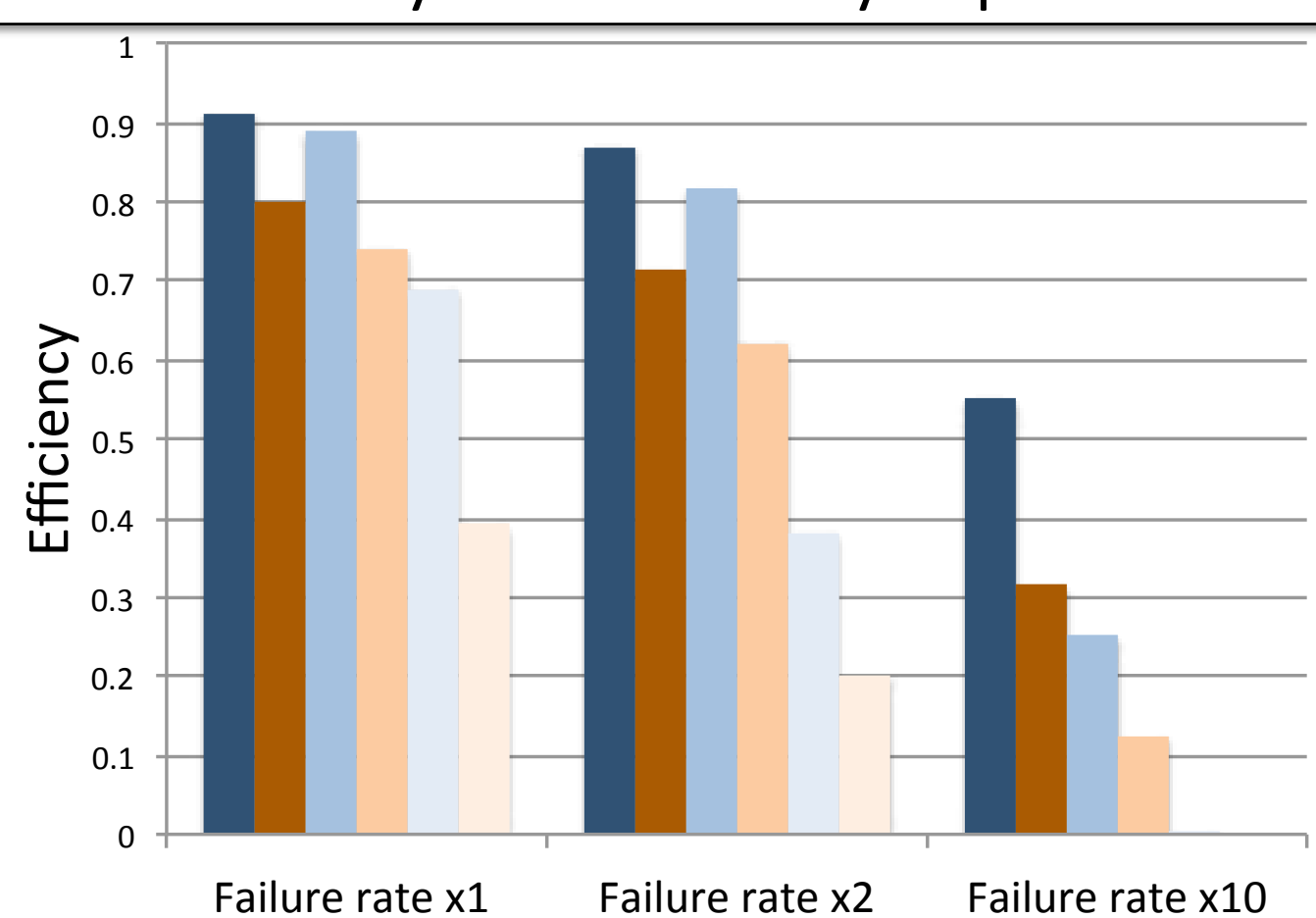- PFS cost x10 / Synchronous

## Multi-tier Resilient Storage Design

- A burst buffer is a storage space to bridge the gap in latency and bandwidth between node-local storage and the PFS
  - Shared by a subset of compute nodes
- Although additional nodes are required, several advantages
  - More Reliable because burst buffers are located on a smaller # of nodes
  - Efficient utilization of storage resources with uncoordinated checkpointing
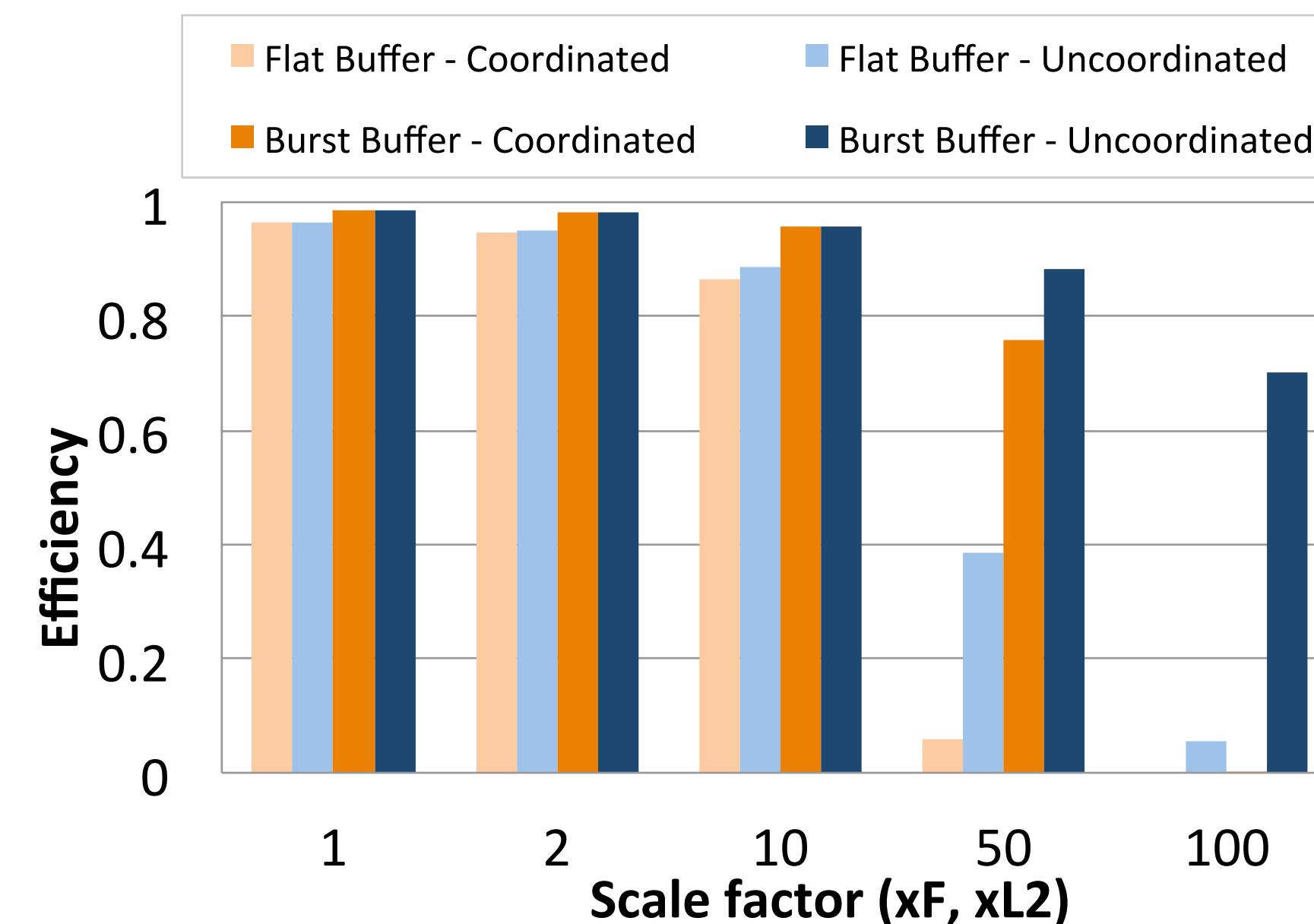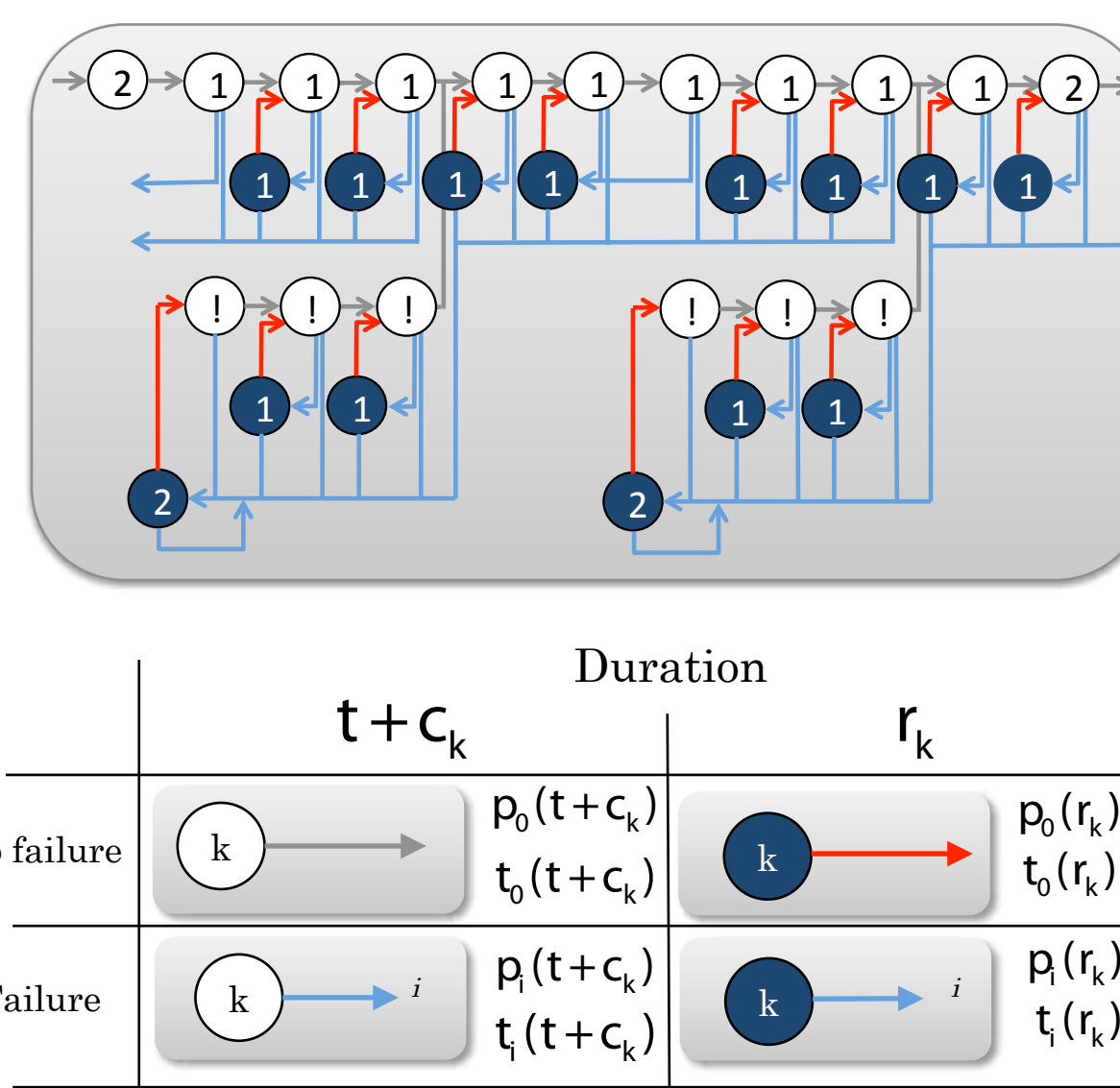
Coordinated C/R



P0 P1 P2 P3 ckpt ckpt

Uncoordinated C/R

P0 P1 P2 P3 ckpt msg logging ckpt

Compute node 1 / Compute node 2 / Compute node 3 / Compute node 4
SSD SSD SSD SSD
Burst buffer        Burst buffer
PFS (Parallel file system)



Duration

| | $t + c_k$ | $r_k$ |
|---|---|---|
| No failure | k — $p_0(t+c_k)$ $t_0(t+c_k)$ | k — $p_0(r_k)$ $t_0(r_k)$ |
| Failure | k — i $p_i(t+c_k)$ $t_i(t+c_k)$ | k — i $p_i(r_k)$ $t_i(r_k)$ |



- Flat Buffer - Coordinated
- Flat Buffer - Uncoordinated
- Burst Buffer - Coordinated
- Burst Buffer - Uncoordinated

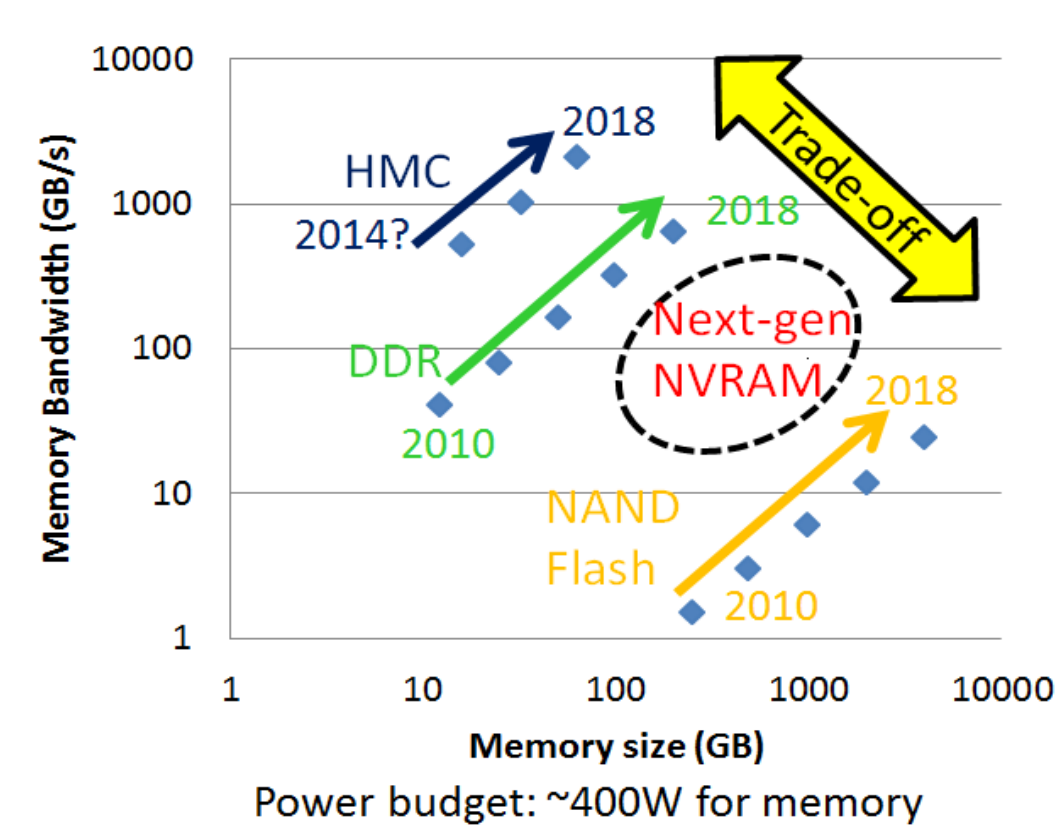Efficiency / Scale factor (xF, xL2)

# Dealing with Deeper Memory Hierarchy

In Exa-scale supercomputing systems, the "**memory wall**" problem will become even higher, which prevents the realization of exa-scale real world simulations.

In our project, "Software Technology that Deals with Deeper Memory Hierarchy in Post-petascale Era", we promote research in aspect of "Architecture", "Algorithm" and "System software".
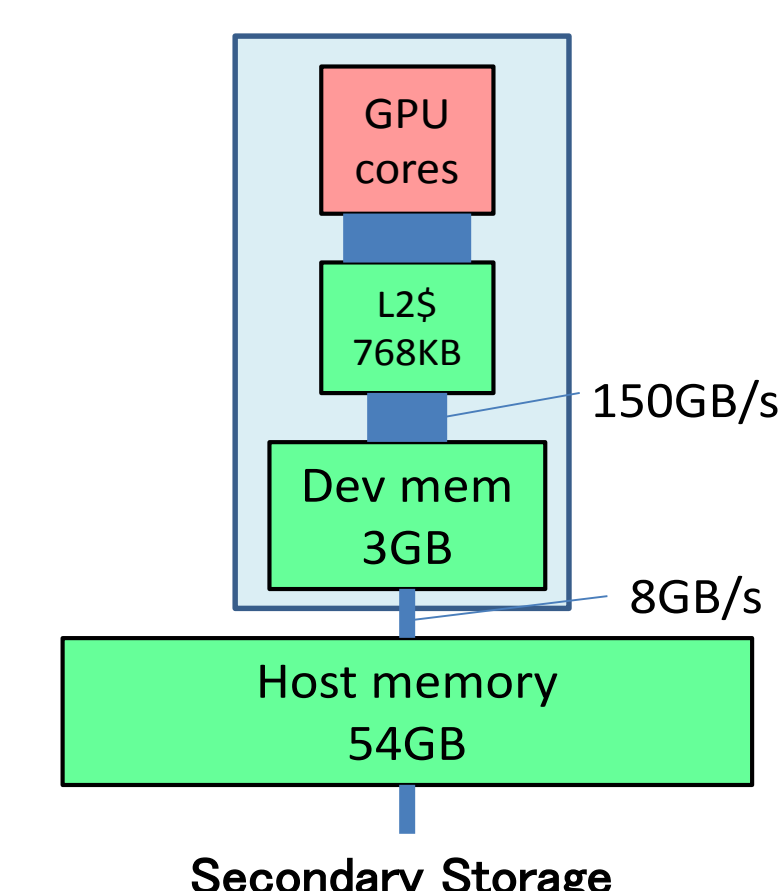
## [Architecture]

To suppose supercomputerm architecture with **deeper memory hierarchy** including hybrid memory devices, including non-volatile RAM (NVRAM).



Memory Bandwidth (GB/s) / Memory size (GB)
Power budget: ~400W for memory
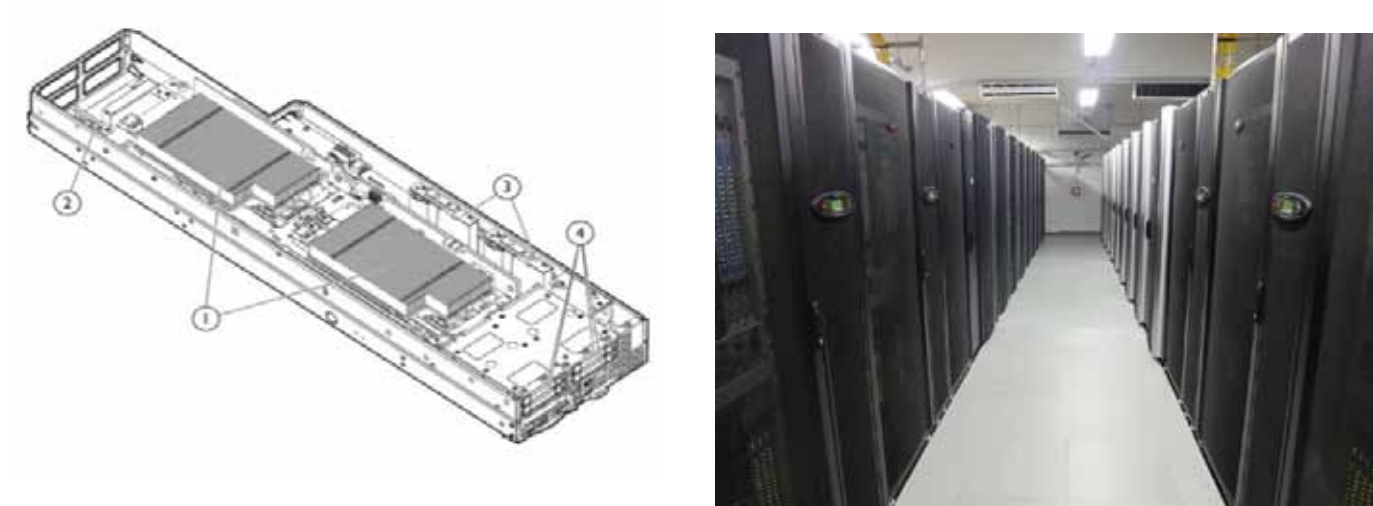HMC 2018 / 2014? / DDR 2010 / Next-gen NVRAM 2018 / NAND Flash 2010 / Trade-off

Hybrid Memory Cube (HMC):
DRAM chips are stacked with TSV technology. It will have advantage in bandwidth over DDR, but capacity will be smaller.

NAND Flash:
SSDs are already commodity. Newer products, such as IO-drive have O(GB/s) bandwidth.

Next-gen non-volatile RAM (NVRAM):
Several kinds of NVRAM such as STT-MRAM, ReRAM, FeRAM, etc, will be available in a few years.



GPU cores
L2$ 768KB
Dev mem 3GB — 150GB/s
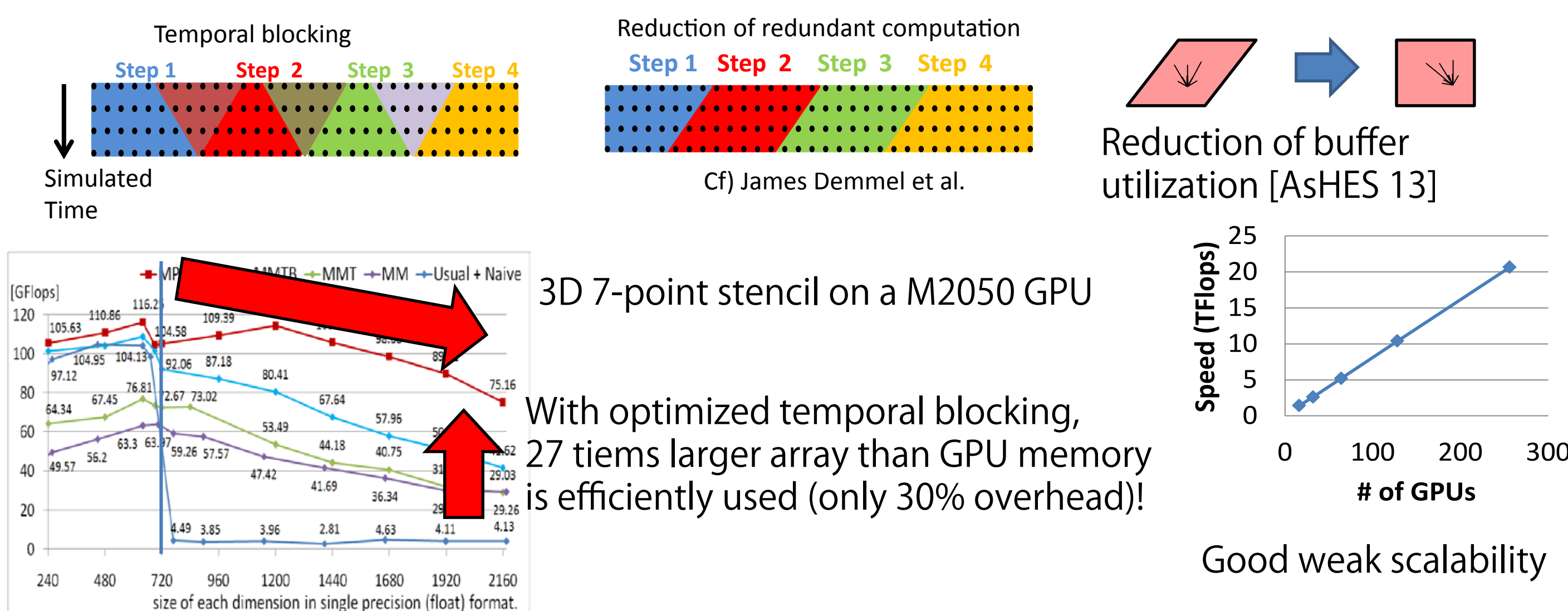Host memory 54GB — 8GB/s
Secondary Storage

Currently, we use TSUBAME2, CPU-GPU hybrid supercomputer as research environment. Here we have memory hierarchy of GPU device memory and Host memory.

## [Algorithm]

To harness hierarchical memory efficiently, we are investigating **locality improvement** of application algorithms. In stencil applications, **temporal blocking** is the key.



Temporal blocking
Step 1 Step 2 Step 3 Step 4
Simulated Time

Reduction of redundant computation
Step 1 Step 2 Step 3 Step 4
Cf James Demmel et al.

Reduction of buffer utilization [AsHES 13]

3D 7-point stencil on a M2050 GPU



[GFlops] / size of each dimension in single precision (float) format.

With optimized temporal blocking, 27 tiems larger array than GPU memory is efficiently used (only 30% overhead)!



Speed (TFlops) / # of GPUs
Good weak scalability

## [System Software]

To support real applications to harness hierarchical memory with lower development efforts, system software support is necessary. Our target includes **locality aware compiler** and **scalable memory management runtime**.

PI: Toshio Endo. Supported by JST-CREST