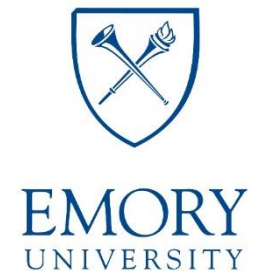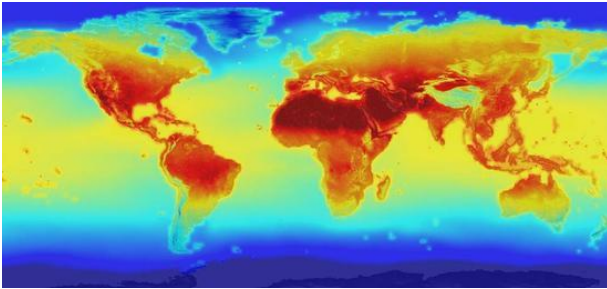# HPC China 2019: 并行存储系统论坛

# Automatic Application-Aware Forwarding Resource Allocation
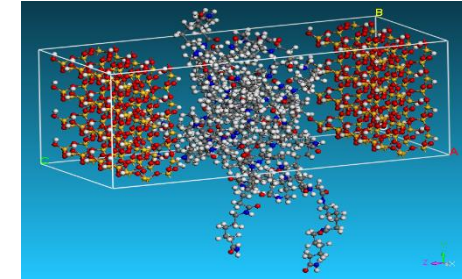
Xu Ji, Bin Yang, Tianyu Zhang, Xiaosong Ma, Xiupeng Zhu, Xiyang Wang, Nosayba EI-sayed, Jidong Zhai, Weiguo Liu, **Wei Xue**
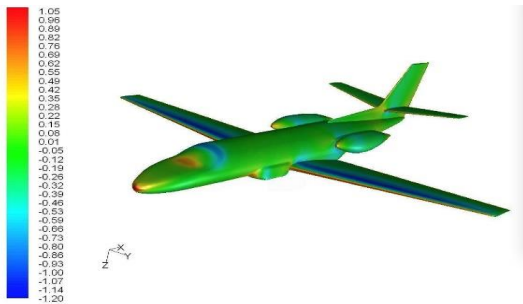
# Storage Crucial for Supercomputing



Climate simulation



Molecular dynamics



Fluid dynamics



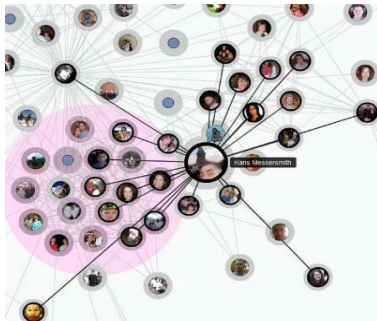Sunway TaihuLight



Summit



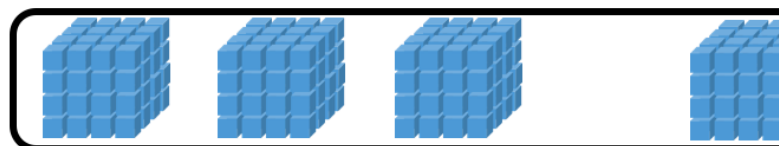Sierra



Bioscience



Graph computing
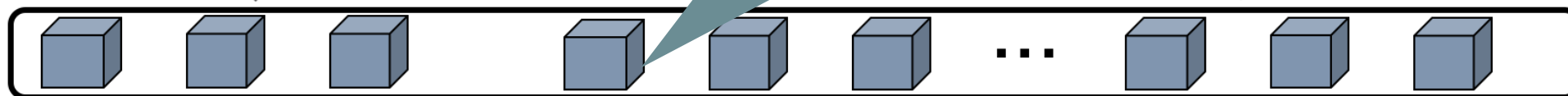


Neural network

# I/O Forwarding Layer



I/O Forwarding architecture
- Enables lean OS on compute nodes
- Reduces #connection
- Connects different network domains
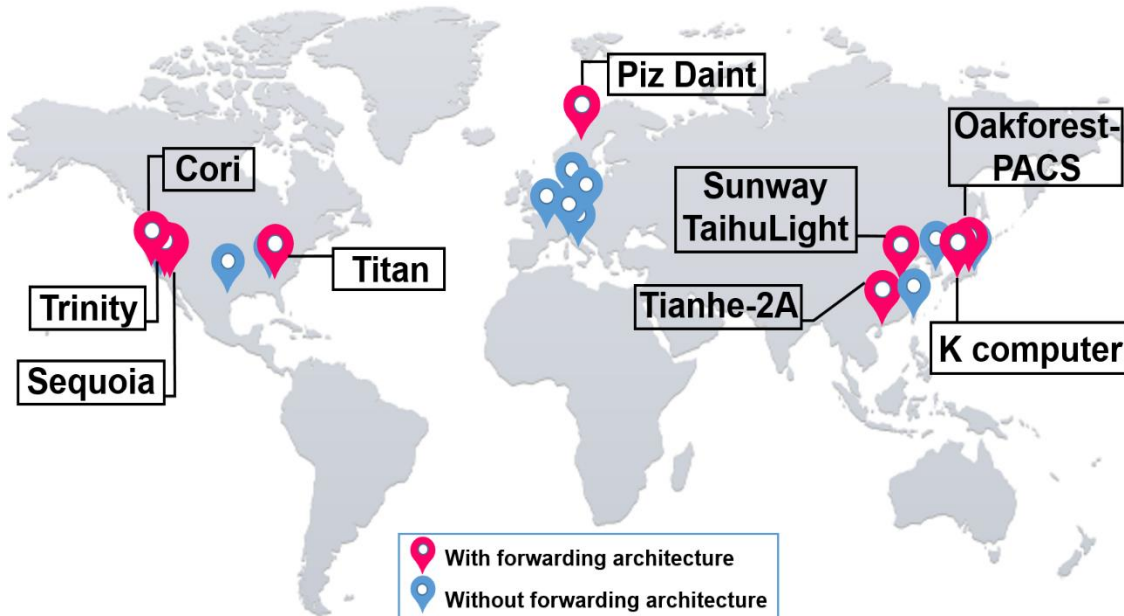- New layer of perfecting/caching

Storage system crucial for speed + scalability
- Performance bottleneck
- Contention point => inter-job performance interference

# I/O Forwarding Widely Adopted

- 9 out of top 20 supercomputers use I/O forwarding (Nov 2018)



| Rank | Machine | Vendor | # Compute nodes | # Forwarding nodes | File system |
|------|---------|--------|-----------------|--------------------|-------------|
| 3 | TaihuLight | NRCPC | 40,960 | 240 | Lustre |
| 4 | Tianhe-2A | NUDT | 16,000 | 256 | Lustre+H2FS |
| 5 | Piz Daint | Cray | 5,320 | 54 | Lustre+GPFS |
| 6 | Trinity | Cray | 14,436 | 576 | Lustre |
| 9 | Titan | Cray | 18,688 | 432 | Lustre |
| 10 | Sequoia | BlueGene | 98,304 | 768 | Lustre |
| 12 | Cori | Cray | 12,076 | 130 | Lustre+GPFS |
| 14 | Oakforest-PACS | Fujitsu | 8,208 | 50 | Lustre |
| 18 | K computer | Fujitsu | 82,944 | 5,184 | FEFS |

The ratio between compute and forwarding nodes range between 170 to 1 and 16 to 1

# Challenges with I/O Forwarding

- Forwarding resource *provisioning*
  - How many forwarding nodes and how much forwarding capacity?
    - Need to consider application I/O demands, machine utilization, budget constraint
    - Design and procurement often finish years before application test runs

- Forwarding resource *management*
  - How many forwarding nodes to allocate to each job?
  - Which forwarding nodes to assign to each job?


- Current common practice: *Fixed Forwarding Mapping (FFM)*
  - Static compute-to-forwarding node mapping

# Sample Forwarding Configuration: TaihuLight



TaihuLight compute nodes (40960)

40 cabinets

Supernode (256 nodes)

...

I/O forwarding nodes (240)

Storage nodes (2*144)

Interconnect network

6 OSTs

Sugon DS800 Disk arrays (144)

...

...

Metadata nodes (2)

IB FDR          IB FDR

# FFM Problem 1: Resource Misallocation


BW-heavy

**CESM** climate simulator


1-1 I/O mode

**CAM** (2017 GB finalist) atmospheric model


Metadata-heavy

**DNDC** agro-ecosystems


IOPS-heavy

**APT** particle dynamics


BW-heavy

**LAMMPS** molecular dynamics


N-1 I/O mode

**AWP** (2017 GB prize) earthquake simulation


BW-heavy

**XCFD** fluid dynamics


1-1 I/O mode

**WRF** weather forecasting


BW-heavy

**Shentu** (2018 GB finalist) graph engine



**swDNN** neuronal network

## Apps have different I/O behaviors!

# FFM Problem 1: Resource Misallocation



**N-N I/O mode**

$P_1$  $P_2$  ...  $P_N$

File1  File2  FileN

**N-1 I/O mode**

$P_1$  $P_2$  ...  $P_N$

File  ...

**1-1 I/O mode**

$P_1$

File



**WRF** weather forecasting

■ WRF-256(1-1)  ■ XCFD-128(N-N)

I/O speedup vs # Forwarding nodes



**XCFD** fluid dynamics

# FFM Problem 2: I/O Interference

I/O interference brings I/O performance inconsistency and degradation



**AWP** earthquake simulation

**Shentu** graph engine

**LAMMPS** molecular dynamics

compute nodes    fwd nodes

dedicated fwd nodes

compute nodes    fwd nodes

shared fwd nodes

Time in seconds

■ AWP-256    ■ Shentu-1024    ■ LAMMPS-256

dedicated fwd nodes: AWP-256 = 1692, Shentu-1024 = 966, LAMMPS-256 = 306

shared fwd nodes: AWP-256 = 1758, Shentu-1024 = 2832, LAMMPS-256 = 378

Applications may suffer from I/O interference at I/O forwarding layer!

# FFM Problem 3: Forwarding Node Anomaly

Fail-slow is also a significant problem to FFM



App performance severely hurt when assigned fail-slow forwarding nodes

Fail-Slow at Scale: Evidence of Hardware Performance Faults in Large Production Systems, Gunawi, FAST'18

# Our Solution: *DFRA (Dynamic Forwarding Resource Allocation)*

Main idea: to ***upgrade*** forwarding node allocation on demand



Normal - default policy

Disruptive – private allocation

Demanding - more allocation

# Per-App I/O Profiling : *Beacon End-to-end I/O Monitor*

- Automatic, transparent multi-level I/O monitoring
  - Allows DFRA to look up node status and per-application I/O profiles
  - Intuition: HPC applications have **consistent I/O behavior** across runs (jobs)



- Deployed at TaihuLight since April, 2017
- Paper at NSDI '19
  - "End-to-end I/O Monitoring on a Leading Supercomputer", by Yang et al.
- Code and I/O monitoring data released
  - https://github.com/Beaconsys/Beacon

# Automatic Forwarding Node Scaling

- Identifying jobs **needing more than default FFM allocation**
- Target job eligible for consideration if application
  - has enough I/O volume ($> V_{min}$), and
  - has enough nodes performing I/O ($> N_{min}$), and
  - is **not** metadata-bound (*avarage metadata queue length* $< W_{metadata}$)

- Estimating adjusted forwarding node number $S$
  - $S = N_{I/O} * B_c / B_f$
    $B_f$ : I/O bandwidth per forwarding node
    $B_c$ : I/O bandwidth per compute node
    $N_{I/O}$ : #Nodes performance I/O
  - Upgrade if $S > N_f$
    - Do nothing otherwise

# I/O Interference Avoidance

- Identifying jobs **requiring dedicated forwarding nodes**
  - Based on comprehensive inter-application I/O interference study
  - Pair-wise interference measurement on 8 apps with representative behaviors
    - 256-node runs, two apps sharing one forwarding node

| Apps | MPI-IO$_N$ | APT | DNDC | WRF$_1$ | WRF$_N$ | Shentu | CAM | AWP |
|---|---|---|---|---|---|---|---|---|
| MPI-IO$_N$ | *(2.1,2.1) | **(1.1,9.3)** | **(4.8,1.1)** | (1.0,1.0) | *(2.1,2.0) | **(1.3,4.5)** | (1.0,1.0) | **(3.3,1.1)** |
| APT | - | *(2.0,2.1) | **(33.3,1.0)** | (1.0,1.0) | **(4.3,1.4)** | **(6.3,1.3)** | (1.0,1.0) | **(50.0,1.1)** |
| DNDC | - | - | *(2.0,2.0) | **(1.0,25.0)** | **(1.0,11.1)** | (1.1,16.7) | **(1.0,33.3)** | *(2.2,2.4) |
| WRF$_1$ | - | - | - | (1.0,1.0) | (1.0,1.0) | (1.0,1.0) | (1.0,1.0) | **(50.0,1.0)** |
| WRF$_N$ | - | - | - | - | *(2.1,2.1) | *(2.0,2.3) | (1.0,1.0) | **(12.5,1.3)** |
| Shentu | - | - | - | - | - | *(2.0,2.0) | (1.0,1.0) | **(12.5,1.1)** |
| CAM | - | - | - | - | - | - | (1.0,1.0) | **(100.0,1.0)** |
| AWP | - | - | - | - | - | - | - | *(2.0,2.0) |

I/O slowdown factor pairs

# I/O Interference Study: Summary

- Identified I/O interference **root causes**:
  - N-1 I/O mode, metadata-heavy, and high-bandwidth workloads;
- High-IOPS workloads more vulnerable

- Detect I/O interference by checking both **_target job_** and **_existing neighbors_** on shared forwarding nodes to be allocated via FFM
  - Assign target job dedicated forwarding nodes if either party belongs to above categories

# DFRA Workflow

# Implementation and Deployment Status

- Implemented to be used with SLURM-based scheduler, in **C and Python**

- **Remapping to more/dedicated forwarding** nodes when job approved for upgrade
  - By relinking RDMA connection
  - **Per-job** basis, new mapping removed at end of job run
  - Currently no "downgrading"

- Partially deployed on TaihuLight production system since *Feb 2018*
  - Users "*opt in*" with job submission command
  - Intend to switch to "*opt out*" in the future

# Evaluation 1: Upgrade Eligibility of Historical Jobs

- 18 months I/O history analysis (Apr 2017 – Sep 2018)

```
[2016-08-27 10:14:44] T lwfs-fuse: 69742: OPEN() btio.full.out => 0x200095032d0
[2016-08-27 10:14:44] T lwfs-fuse: 69745: WRITE (0x200095032d0, size=161320, offset=0, count=2)
[2016-08-27 10:14:44] T lwfs-fuse: 69740: RELEASE 0x20009500e40 (FLUSH implied)
[2016-08-27 10:14:49] T lwfs-fuse: 69805: OPEN() btio.full.out => 0x20009504160
[2016-08-27 10:14:49] T lwfs-fuse: 69811: READ (0x20009504160, size=262144, offset=262144, count=4)
```

Inefficient I/O processes , 8%

Benefit from DFRA, 14%

Low I/O volume, 78%

**By job count**

Low I/O volume, 18%

Inefficient I/O processes , 2%

Benefit from DFRA, 80%

**By core-hours consumed**

Few jobs but consuming considerable core-hours are benefiting from DFRA

# Evaluation 2: DFRA Effectiveness

Overall performance from 10 runs (each box) in production environment
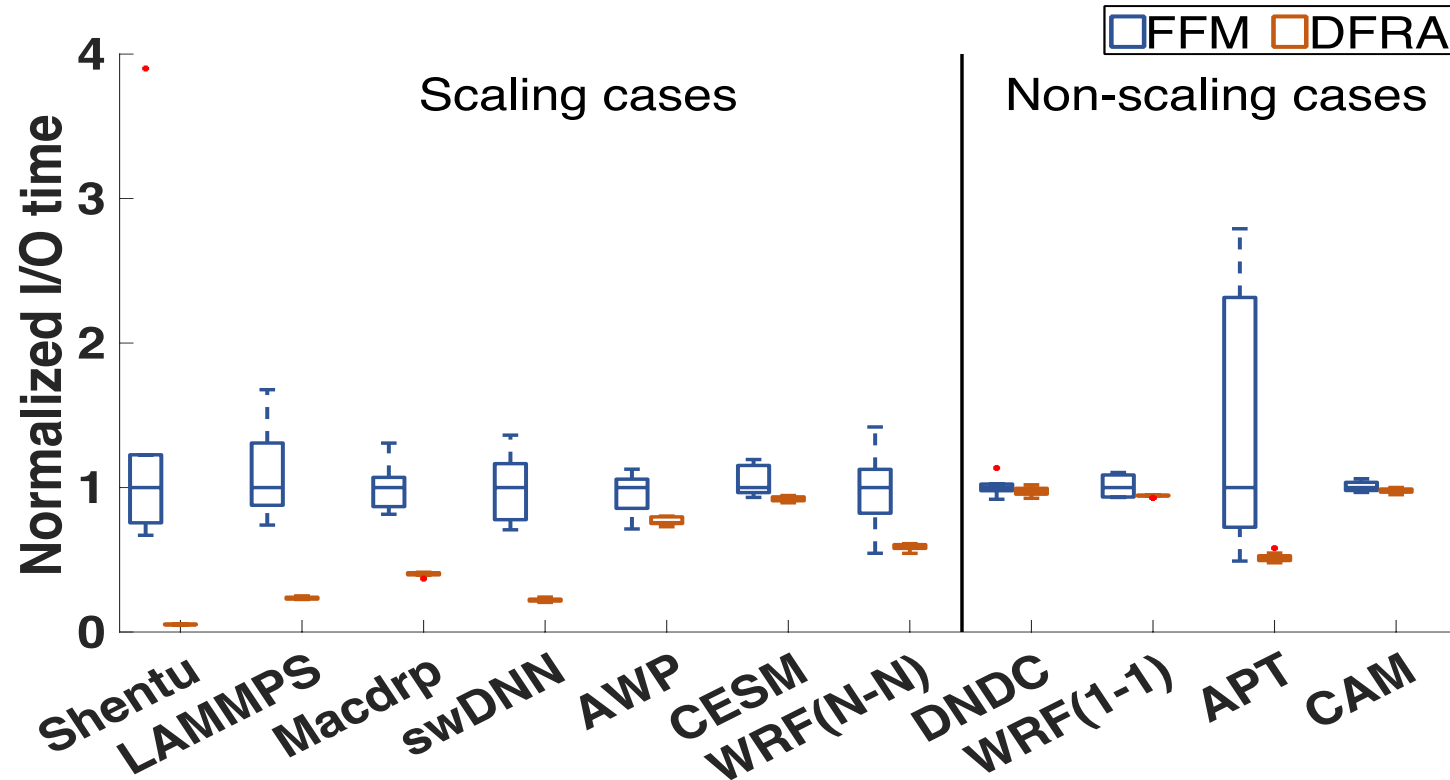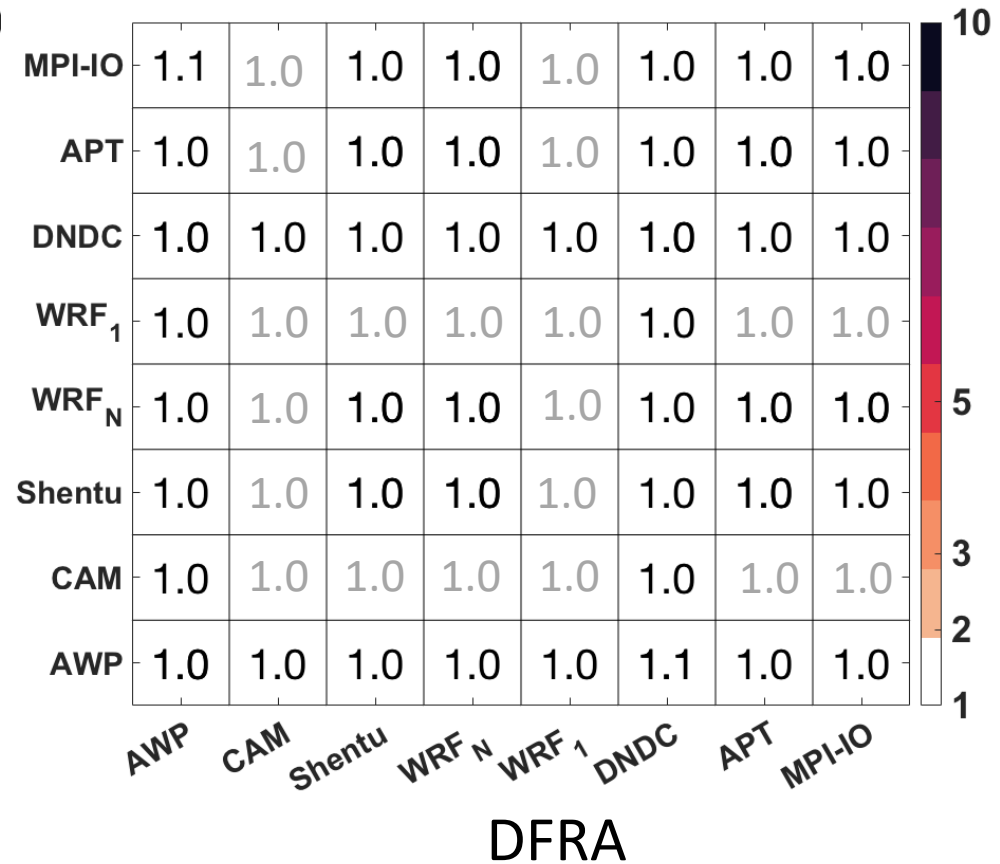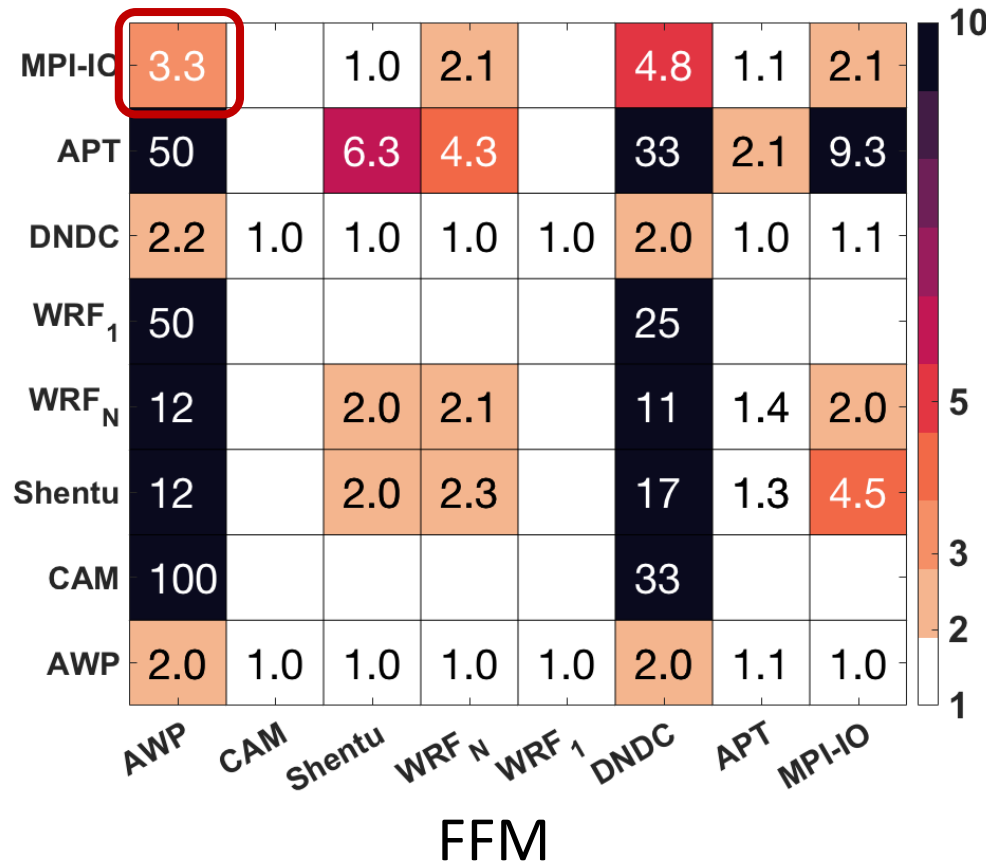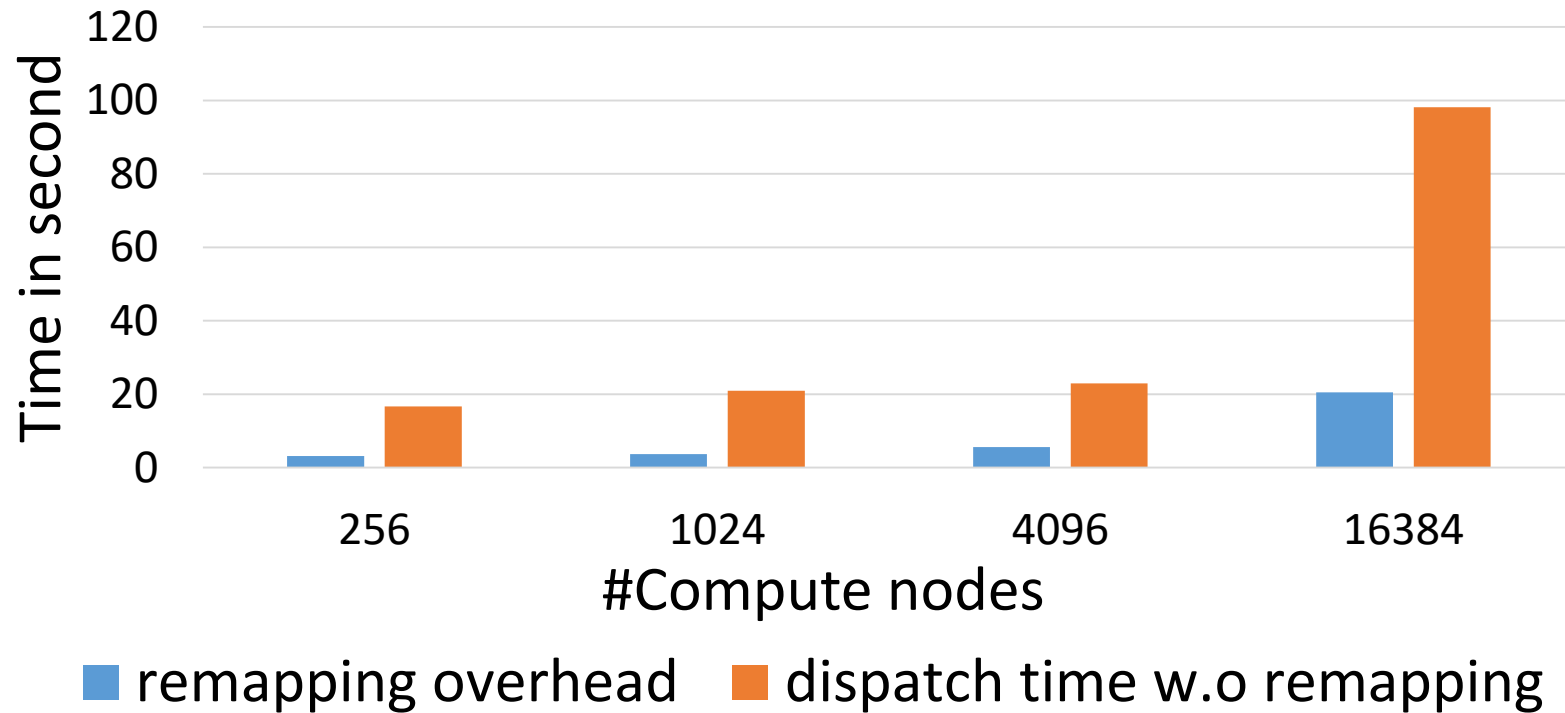


DFRA reduces I/O variation and improves I/O performance at I/O forwarding layer

# Evaluation 3: I/O Interference Reduction



FFM

| | AWP | CAM | Shentu | WRF$_N$ | WRF$_1$ | DNDC | APT | MPI-IO |
|---|---|---|---|---|---|---|---|---|
| MPI-IO | 3.3 | | 1.0 | 2.1 | | 4.8 | 1.1 | 2.1 |
| APT | 50 | | 6.3 | 4.3 | | 33 | 2.1 | 9.3 |
| DNDC | 2.2 | 1.0 | 1.0 | 1.0 | 1.0 | 2.0 | 1.0 | 1.1 |
| WRF$_1$ | 50 | | | | | 25 | | |
| WRF$_N$ | 12 | | 2.0 | 2.1 | | 11 | 1.4 | 2.0 |
| Shentu | 12 | | 2.0 | 2.3 | | 17 | 1.3 | 4.5 |
| CAM | 100 | | | | | 33 | | |
| AWP | 2.0 | 1.0 | 1.0 | 1.0 | 1.0 | 2.0 | 1.1 | 1.0 |

DFRA

| | AWP | CAM | Shentu | WRF$_N$ | WRF$_1$ | DNDC | APT | MPI-IO |
|---|---|---|---|---|---|---|---|---|
| MPI-IO | 1.1 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 |
| APT | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 |
| DNDC | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 |
| WRF$_1$ | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 |
| WRF$_N$ | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 |
| Shentu | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 |
| CAM | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 |
| AWP | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.1 | 1.0 | 1.0 |

DFRA reduce variation and improve I/O performance at I/O forwarding layer

# Evaluation 4: Remapping Overhead vs. Saving



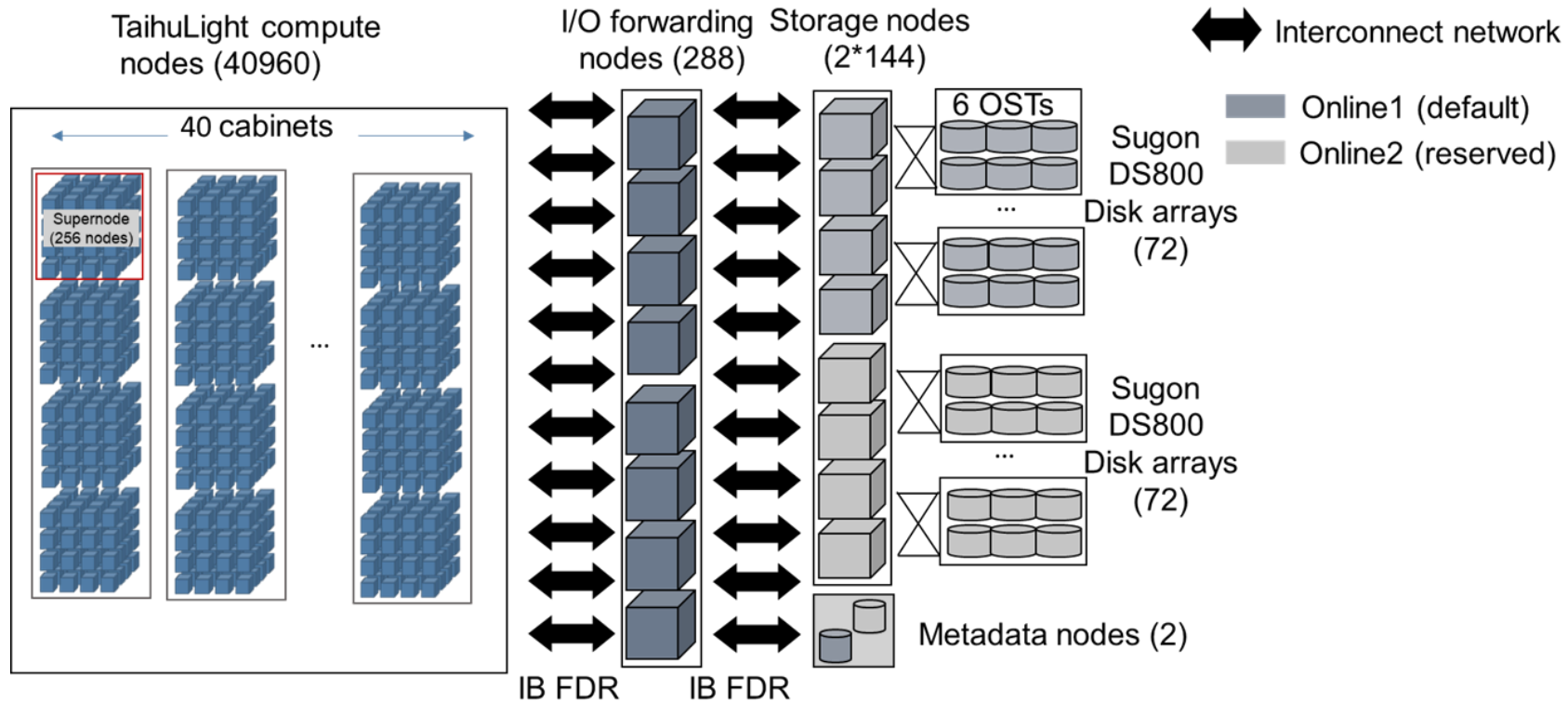Chart — Time in second vs. #Compute nodes

Legend: ■ remapping overhead   ■ dispatch time w.o remapping

Overhead acceptable considering benefit
- Average I/O time saving of 6 minutes for I/O-intensive jobs
- Estimated saving of 200 million core-hours in past 8 months

# Conclusion



**Take away points:**

- Don't guess future user I/O demands
    - Over-provision, give low "basic plan", then upgrade when needed
- DFRA applicable to shared burst buffer management too

# Q&A

Thank you!

Partial I/O monitoring data released at
https://github.com/Beaconsys/Beacon