# Gfarm distributed file system

Osamu Tatebe

University of Tsukuba

# Gfarm file system

- Award-winning file system since 2000
  - Distributed infrastructure award in SC03
  - Most Innovative Use of Storage In Support of Science Award in SC05
  - Winner – Large Systems in HPC Storage Challenge in SC06
- Open Source distributed file system
  - http://sf.net/projects/gfarm/
- Supported by NPO OSS Tsukuba Support Center
- Features
  - Scaled-out performance in wide area
  - Data access locality, file replica
  - No single point of failure
    - Automatic file replica creation in case of storage failure
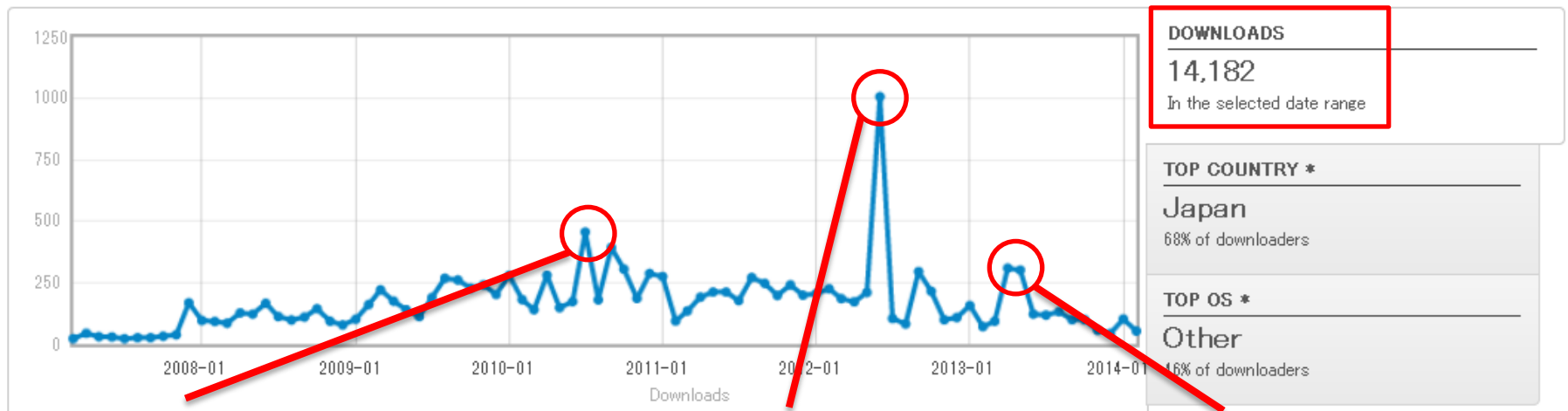    - Hot stand-by MDS

# # downloads

- 14,182 downloads since March, 2007

**Gfarm File System**

Summary | Files | Reviews | Support | Mailing Lists | Trac | News | Code | Tickets | Wiki

🏠 Home (Change File)

Date Range: 2007-03-01 to 2014-02-18

**DOWNLOADS**
14,182
In the selected date range

**TOP COUNTRY \***
Japan
68% of downloaders

**TOP OS \***
Other
16% of downloaders

*(graph: Downloads over time from 2008-01 to 2014-01)*

2010/7
Version 2.3.2, 2.4.0
456 downloads

2012/6
HPCI installation etc
1,007 downloads
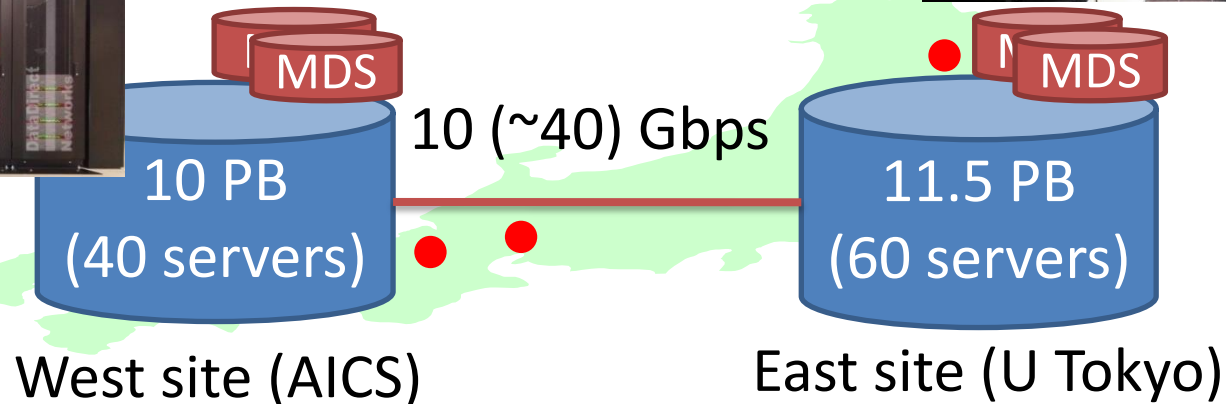
2013/4,5
Version 2.5.8
610 downloads

# Collaboration

- Japan Lattice Data Grid (JLDG)
  - Division of Particle Physics, Division of Computational Informatics
  - KEK, Osaka Univ, Hiroshima Univ, Nagoya Univ, Kyoto Univ, Kanazawa Univ, RIKEN, Univ Tokyo
- GPV/JMA Archive
  - Division of Global Environmental Science, Division of Computational Informatics
- HPCI Shared Storage
  - RIKEN, Nine National Universities, AIST, JAMSTEC, …
- NICT Science Cloud
  - NICT

# HPCI SHARED STORAGE
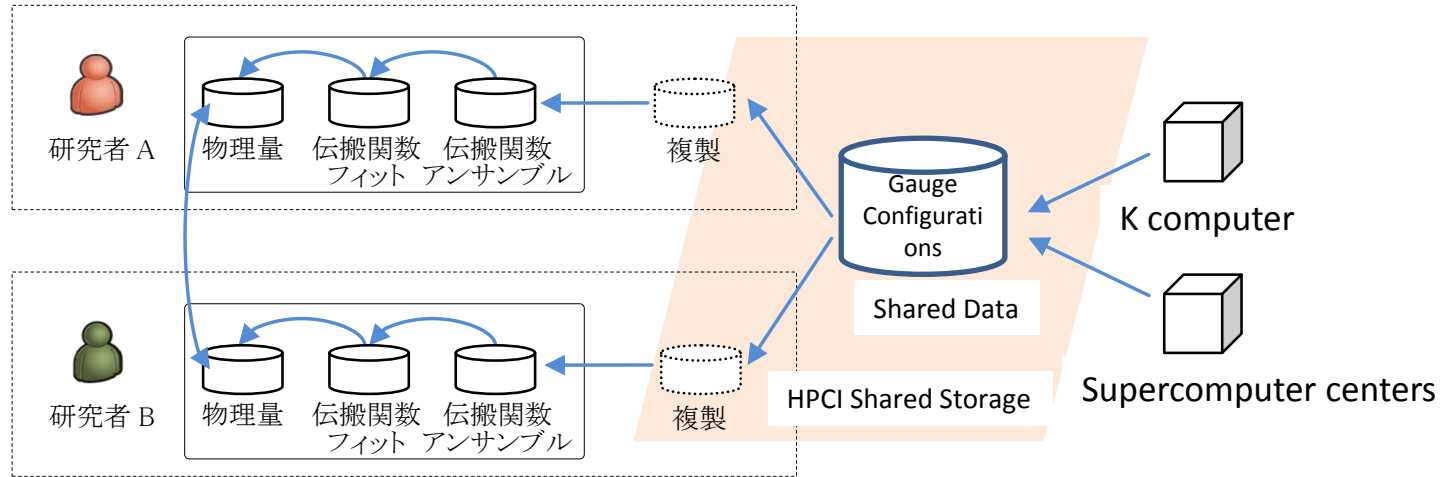
# HPCI Shared Storage

- HPCI – High Performance Computing Infrastructure
  - "K", Hokkaido, Tohoku, Tsukuba, Tokyo, Titech, Nagoya, Kyoto, Osaka, Kyushu, RIKEN, JAMSTEC, AIST
- A 20PB single distributed file system consisting East and West sites
- Grid Security Infrastructure (GSI) for user ID
- Parallel file replication among sites
- Parallel file staging to/from each center



10 (~40) Gbps

MDS

MDS

10 PB
(40 servers)

11.5 PB
(60 servers)

West site (AICS)

East site (U Tokyo)

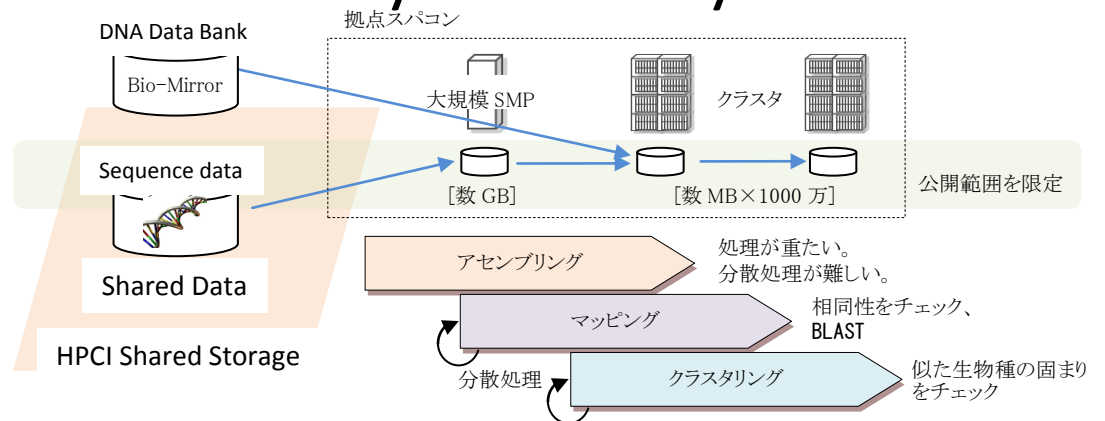Picture courtesy by Hiroshi Harada (U Tokyo)
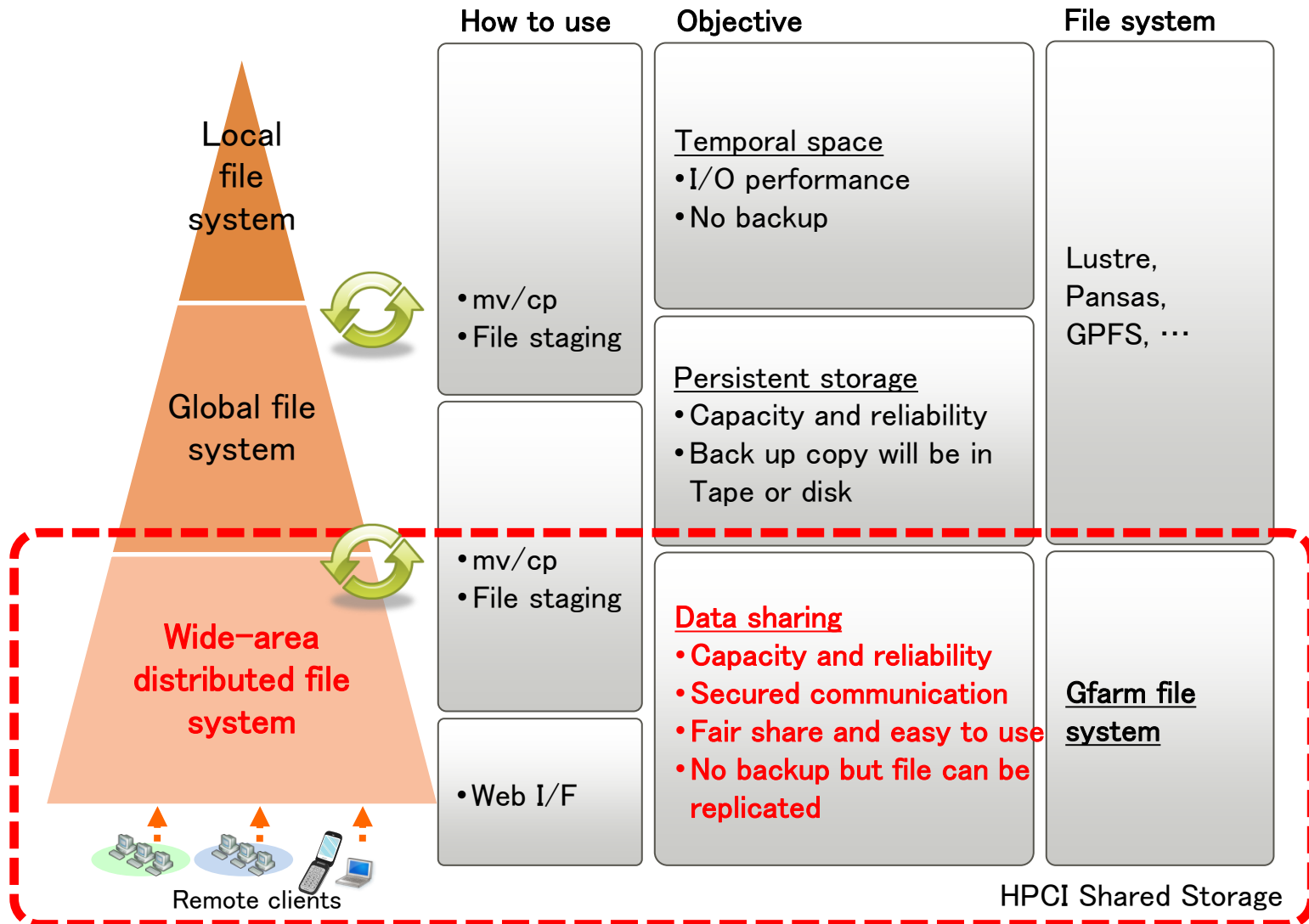
# Usage Scenario

- QCD



- Large scale simulation (ex. cosmic simulation)
  - Simulation data obtained by the K computer, will be analyzed and visualized by University's supercomputer
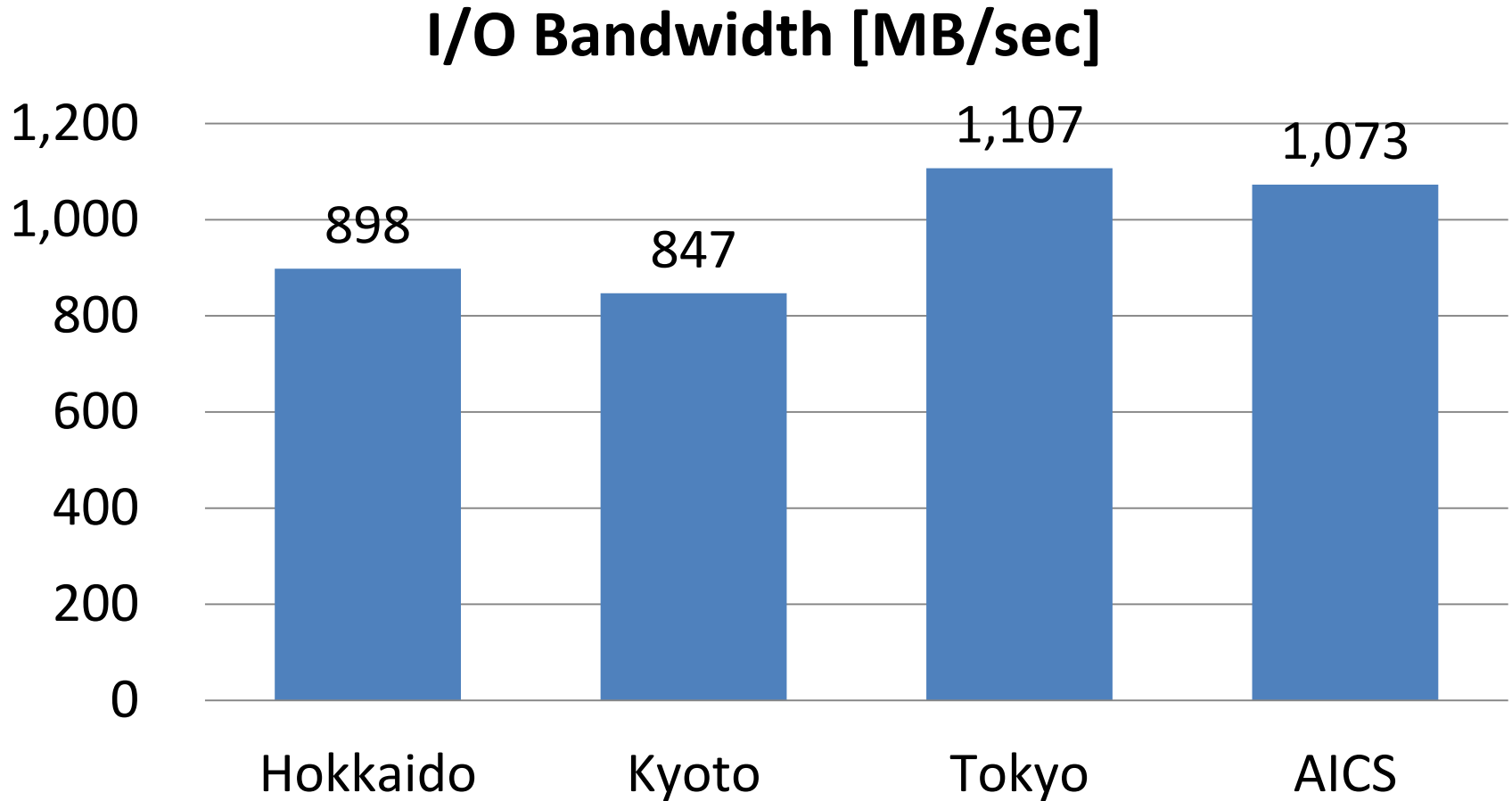
- Life Science

# Storage structure of HPCI Shared Storage

|  | How to use | Objective | File system |
|---|---|---|---|
| Local file system | • mv/cp<br>• File staging | **Temporal space**<br>• I/O performance<br>• No backup | Lustre,<br>Pansas,<br>GPFS, … |
| Global file system | | **Persistent storage**<br>• Capacity and reliability<br>• Back up copy will be in Tape or disk | |
| Wide-area distributed file system | • mv/cp<br>• File staging<br><br>• Web I/F | **Data sharing**<br>• Capacity and reliability<br>• Secured communication<br>• Fair share and easy to use<br>• No backup but file can be replicated | **Gfarm file system** |

Remote clients

HPCI Shared Storage

# How to use HPCI shared storage

```
% mount.hpci                          # mount command
Update proxy certificate for gfarm2fs
timeleft : 167:50:40  (7.0 days)
Mount GfarmFS on /gfarm/hp120273/tatebe
% df -H /gfarm/hp120273/tatebe
Filesystem        Size   Used  Avail Use% Mounted on
fuse              23P   2.9P   20P  14% /gfarm/hp120273/tatebe
% cd /gfarm/hp120273/tatebe
% gfpcopy –P /work/CSI/tatebe/data .  #parallel copy command
….
copied_file_num: 10
copied_file_size: 6553600000
total_throughput: 70.233735 MB/s
total_time: 93.311284 sec.
% gfncopy –s 2 data                   #specify # replicas
(file repilcas are automatically created on background）
```

# Initial Performance Result

**I/O Bandwidth [MB/sec]**



File copy performance of 300 1GB files

# Related Work

- XSEDE-Wide File System (GPFS)
  - Planned, but not in operation yet
- DEISA Global File System
  - Multicluster GPFS
    - RZG, LRZ, BSC, JSC, EPSS, HLRS, …
    - Site name included in the path name – no distribution transparency
      - files cannot be replicated across sites
  - PRACE does not provide global file system
    - Limitation of operation systems that can mount
    - PRACE does not assume to use multiple sites
- Distant access to Lustre File System
  - Many researches in TeraGrid
    - Increase the number of pending requests
  - SC11 paper showed the performance could not be simply improved in private 100Gbps network

# PARALLEL AND DISTRIBUTED DATA ANALYSIS

# Hadoop Gfarm plugin [Mikami, Ohta, Tatebe, IEEE/ACM Grid 2011]

- Design and Implement Gfarm-Hadoop plugin to access POSIX compatible Gfarm file system from Hadoop apps

- Compare with HDFS, PVFS and GlusterFS



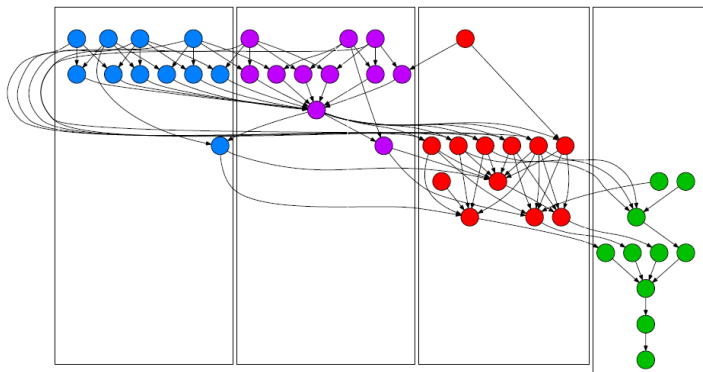**Concurrent write performance**

# Pwrake workflow engine

- Rake extension – parallel and distributed workflow language and execution engine
- http://github.com/masa16/Pwrake/
- Gfarm file system support
  - Automatic mount/umount of Gfarm file system
  - Data aware job scheduling
- Masahiro Tanaka, Osamu Tatebe, "**Pwrake: A parallel and distributed flexible workflow management tool for wide-area data intensive computing**", Proceedings of ACM International Symposium on High Performance Distributed Computing (HPDC), pp.356-359, 2010
- Masahiro Tanaka and Osamu Tatebe , "**Workflow Scheduling to Minimize Data Movement using Multi-constraint Graph Partitioning**", Proceedings of IEEE/ACM International Symposium on Cluster, Cloud and Grid Computing (CCGrid), 2012 (to appear)

# Data aware workflow scheduling
## [Tanaka, Tatebe, CCGrid 2012]

job scheduling by **multi-constraint graph partitioning** to minimize data transfer and maximize parallel job executions

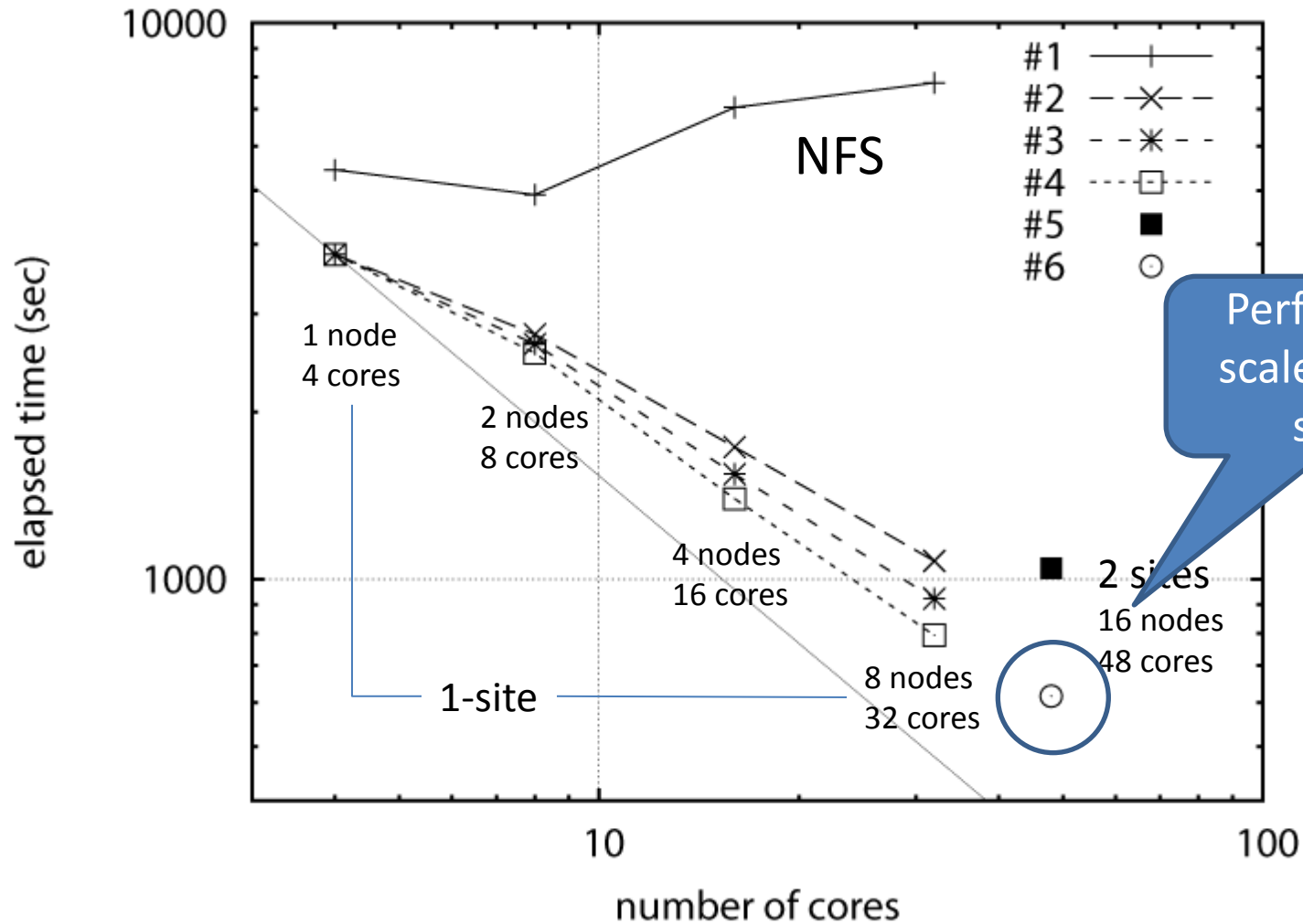Simple graph partitioning | **Multi-constraint graph partitioning**



Load imbalance happens in each parallel step

Reduce 14% of data transfer
Improve 31% of performance

# Performance result of Montage astronomical data analysis workflow

# SELECTED FEATURES

# Consistency check and repair

- Consistency check and repair at MDS startup
- Consistency check and repair at file server startup in parallel
- # replicas is automatically maintained in case of file creation, file server failure, and changing # replicas
  - # replicas can be specified in each directory
    % **gfncopy** –s 3 /home/tatebe

# Gfarm zabbix plugin



Realtime server monitoring and automatic ticket issue

# Gfarm ganglia plugin



# Realtime IOPS and bandwidth monitoring

# NPO Tsukuba OSS Support Center

- http://oss-tsukuba.org/
- Established in Apr, 2013
- Gfarm software support
- Inaugural symposium on Aug 28, 2013
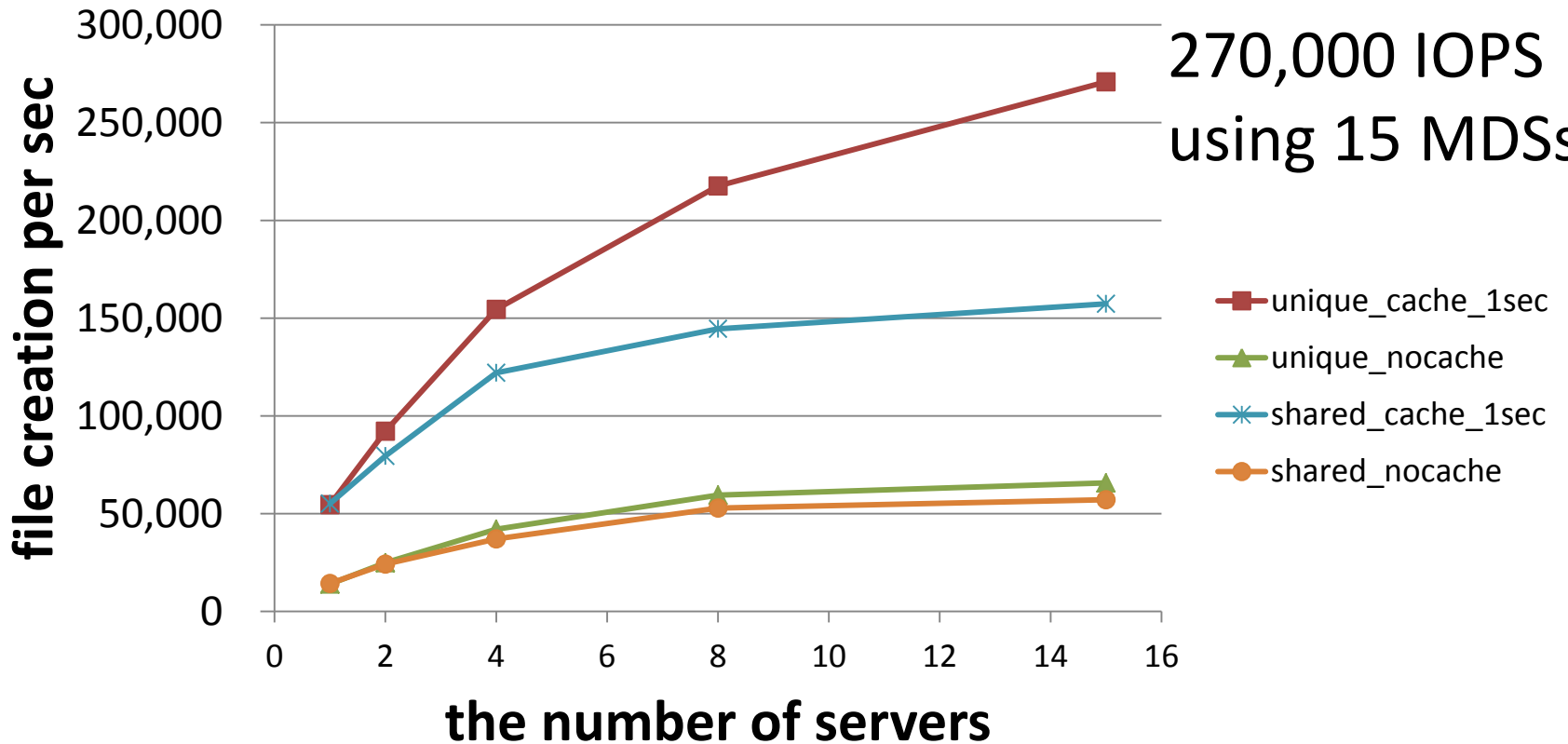
# FUTURE EVOLUTION

# Gfarm 2.6

- Will be released Q1, 2014
- Functionality to specify replica location to be created
- Transparent MDS failover
- Performance improvement of gfpcopy
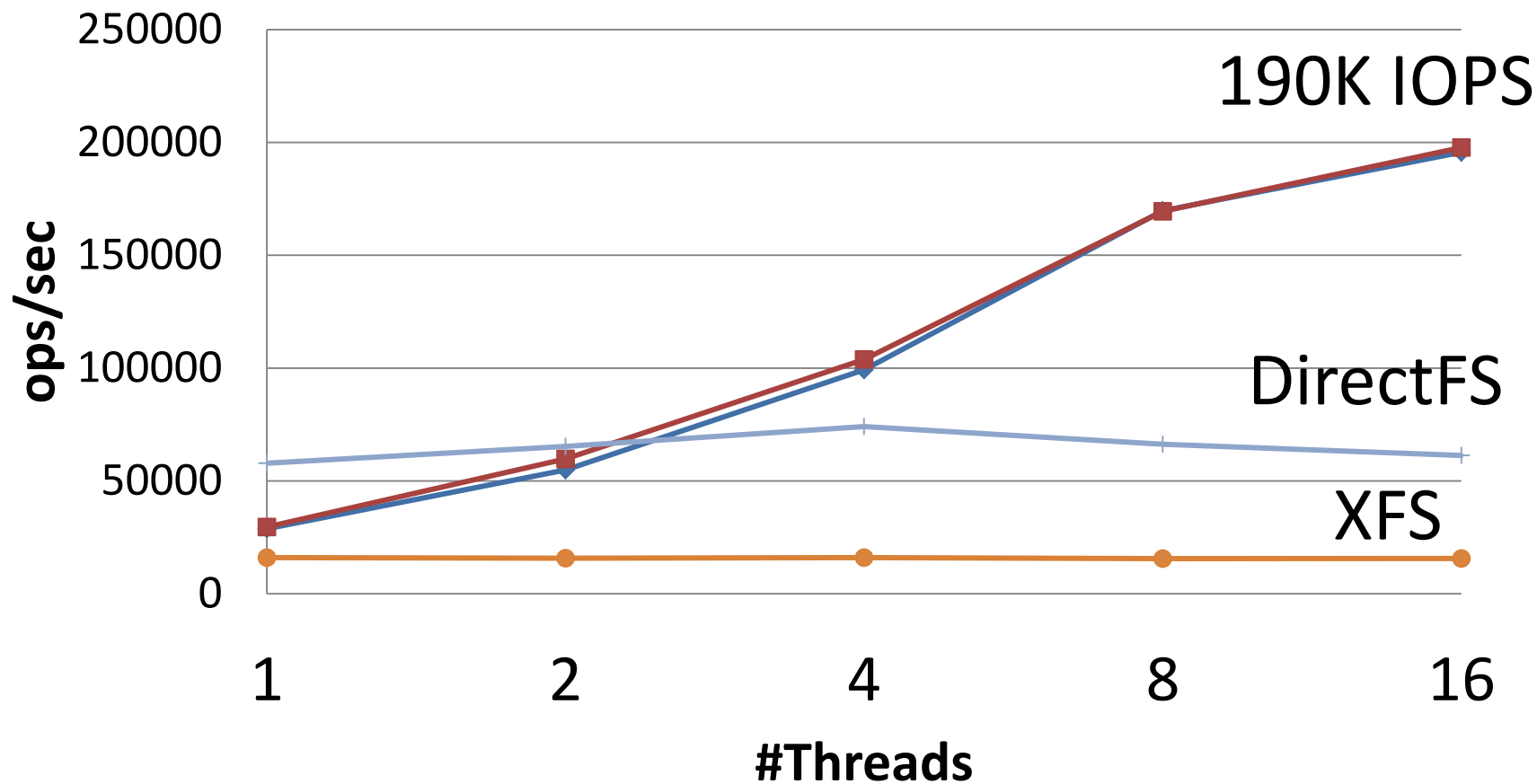
# R&D of distributed MDS [JST/CREST]



distributed dir attr mdtest file creation 128 clients

270,000 IOPS using 15 MDSs

file creation per sec

the number of servers

- unique_cache_1sec
- unique_nocache
- shared_cache_1sec
- shared_nocache

# R&D of object storage for ioDrive [JST/CREST]

**File creation performance [ops/sec]**

# Summary

- Gfarm file system
  - Developed since 2000, O(14,000) downloads
  - HPCI Shared Storage, NICT Science Cloud, Japan Lattice Data Grid (JLDG), companies
- > 1 GB/s parallel copy performance
- Hadoop MapReduce, Workflow, MPI-IO
- Management tools
- NPO OSS Tsukuba Support Center
- R&D for distributed MDS and object store