

# System Software for Post Petascale **Data Intensive Science**

http://www.hpcs.cs.tsukuba.ac.jp/project/crest-ppfs/en/



## **Runtime System**

- **Pwrake: Scalable Workflow System** Workflow DAG Partitioning **Pwrake Master** *Pwrake* is a workflow engine for data-intensive sciences with the following powerful features. Sub-Master Sub-Master • Rake as a Workflow Language [HPDC 2010]
  - Rake is a DSL for a powerful build tool in Rake, which enables to describe complex scientific workflows
  - Scalable I/O Performance based on File Locality.
    - Workflow scheduling to minimize data transfer based on Multi-Constraint Graph Partitioning [CCGrid 2012].

1 serv

2 servs

★ 3 servs

-4 servs

8 servs

14 servs

- Post Petascale system is the next target of Pwrake.
  - Hierarchical Pwrake System manages 100M tasks executed on one million cores.



## **Distributed File System**

### **PPMDS:** A Distributed Metadata Server



- Fine-grained parallelization, the system manages directory namespace in parallel by ordered key-value
- stores. The key consists of a pair of a parent inode number and a file name. The value is the metadata.
- The largest granularity of locking is a key-value pair.
- Nonblocking transactions across multiple key-value servers based
- on Dynamic STM to update multi key-value

pairs transactionaly.





### **File System for next-generation storage**



#### **Redundant data store across nodes**



# Read performance of 64MiB data from 2 server



## **File Access with Infiniband RDMA**

Infiniband is a high throughput and low latency network. The objective of this research is to improve the performance of distributed file systems by using **RDMA** (Remote Direct Memory Access). This function has a

significant effect to reduce the





overhead of communication.

Read size [byte]

#### **Gfarm: Wide-area Distributed File System**





http://www.ccs.tsukuba.ac.jp/