

MATRIX: MAny-Task computing execution fabRiC at eXascale

Ke Wang

Data-Intensive Distributed Systems Laboratory
Computer Science Department
Illinois Institute of Technology

CS554: Data-Intensive Computing, IIT
February 4th, 2015

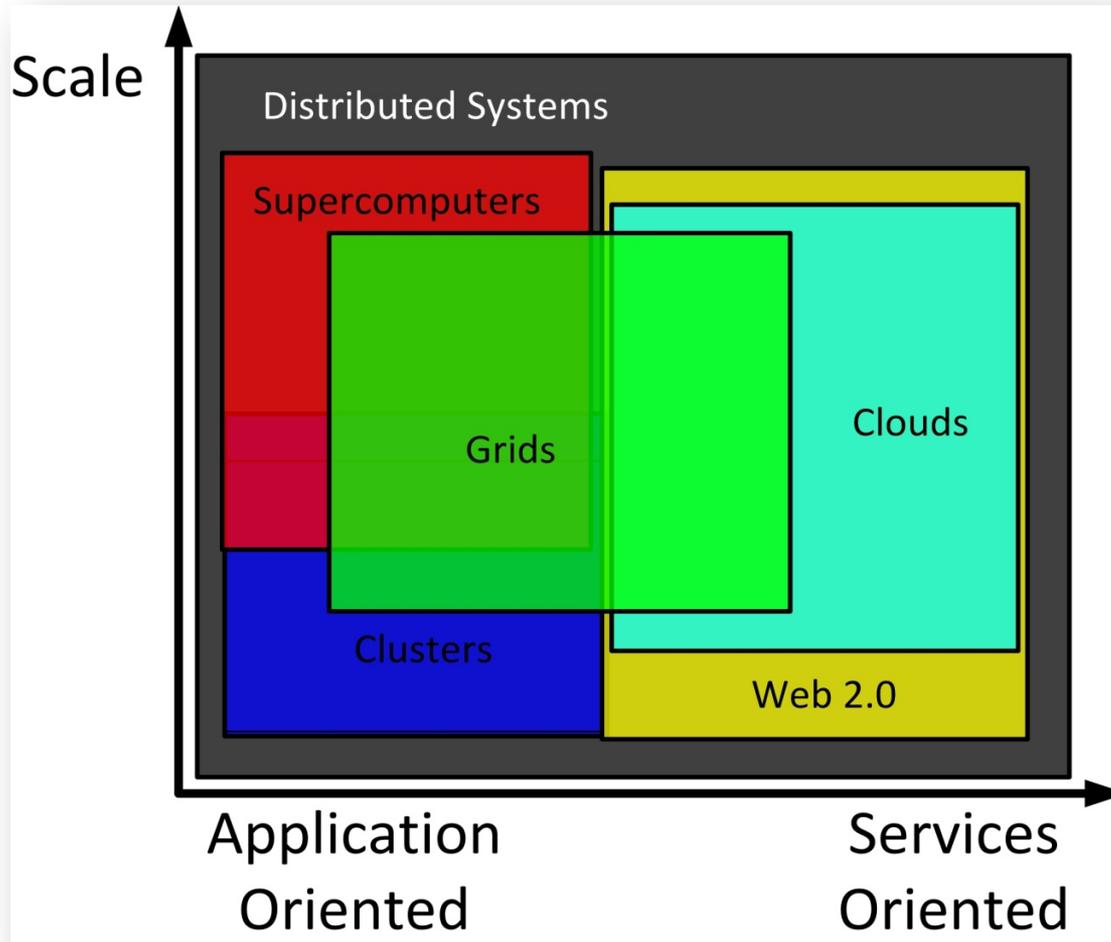
Outline

- **Introduction & Motivation**
- **Problem Statement**
- **Proposed Work**
- **Evaluation**
- **Conclusions**
- **Future Work**

Outline

- **Introduction & Motivation**
- Problem Statement
- Proposed Work
- Evaluation
- Conclusions
- Future Work

Distributed Systems



Exascale Computing

PERFORMANCE DEVELOPMENT



— 1 Billion cores

http://s.top500.org/static/lists/2014/06/TOP500_201406_Poster.pdf

<http://users.ece.gatech.edu/mrichard/ExascaleComputingStudyReports/ECSS%20report%20101909.pdf>

MATRIX: MAny-Task computing execution fabRlc at eXascale

Resource Manager

- Manages resources
 - compute
 - storage
 - network
- Job scheduling/management
 - resource allocation
 - job launch
- Data management
 - data movement
 - caching

Resource Manager Examples

User commands



ActiveBatch® Architecture

The ActiveBatch Enterprise Job Scheduler provides a central point of automation and scheduling that simplifies and manages the integration of applications, platforms and databases into "end to end" workflows across the enterprise.

GUI / Programmatic Interface



ActiveBatch Job Scheduler

Integration

- Amazon EC2
- IBM Cognos BI
- IBM InfoSphere DataStage
- IBM Netezza
- Informatica PowerCenter
- Microsoft Dynamics AX
- Microsoft SharePoint
- Microsoft System Center Suite
 - Operations Manager
 - Orchestrator
 - Service Manager
 - Virtual Machine Manager
- Microsoft Team Foundation Server
- Microsoft Windows Azure
- Oracle E-Business Suite
- Oracle PeopleSoft
- SAP BusinessObjects
- SAP NetWeaver & BW-SCH
- Teradata
- VMware
 - Web Services / REST
 - Other 3rd Party Apps

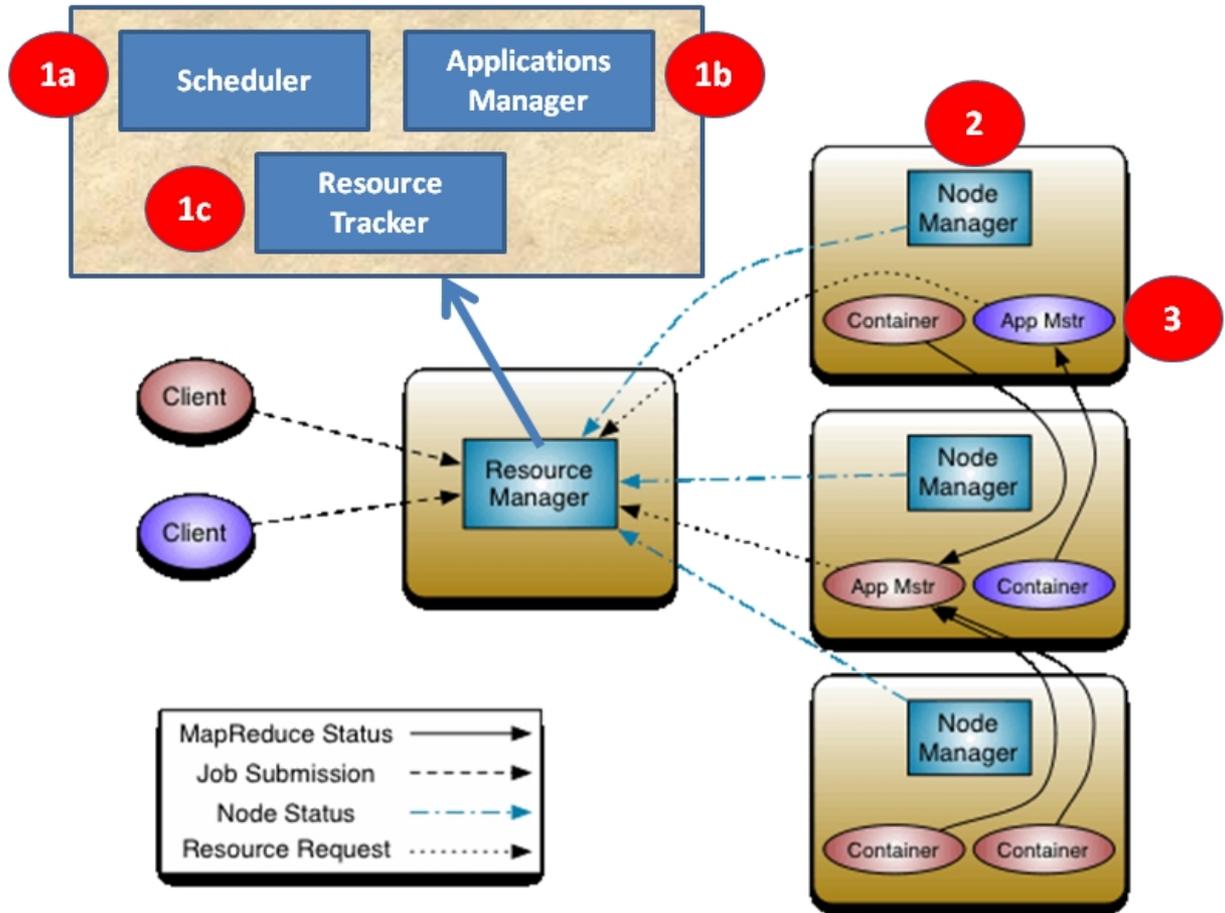


TCP / IP
SSL (TLS)

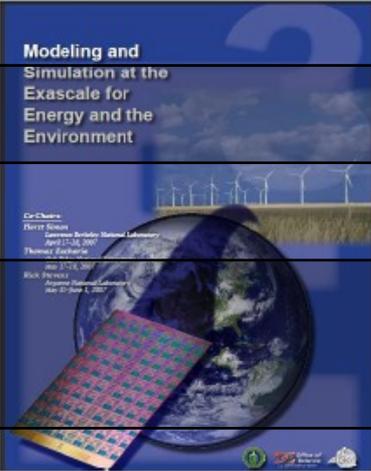
Execution Agents



Resource Manager: Sub-components



Motivation

1970s	<p>Medical Image Processing: Functional Magnetic Resonance</p> <p>Ensemble simulations are important for future exascale platforms</p>	
1980s	<p>Chemistry Domain: MolDyn</p> <p>Molecular Dynamics: DOCK</p> <p><i>One approach to dealing with uncertainty is to perform multiple ensemble runs (parameter sweeps) with various combinations of the Run in parallel. In the space of parameters will be of high dimension, we will have to address the challenges of designing efficient parameter sweep methods for high-dimensional spaces. Recent advances in approximation theory and data mining methods, such as sparse grids, offer new approaches to this problem. Furthermore, recent results in approximation theory can be used to guide us in using exascale computing power to search for efficient methods.</i></p>	
1990s	<p>Production Runs in Drug Design</p> <p>Economic Modeling: MARS</p>	
2000s	<p>Large-scale Astronomy Application Evaluation</p>	
2010s	<p>Astronomy Domain: Montage</p>	
	<p>Data Analytics: Sort and WordCount</p> <p><i>(Source: Horst Simon et al. Modeling and Simulation at the Exascale for Energy and the Environment)</i></p>	

Outline

- Introduction & Motivation
- **Problem Statement**
- Proposed Work
- Evaluation
- Conclusions
- Future Work

Problem Statement

Job Scheduling System Challenges

- **Scalability**
 - System scale is increasing
 - Workload size is increasing
 - Processing capacity needs to increase
- **Efficiency**
 - Allocating resources fast
 - Making fast enough scheduling decisions
 - Maintaining a high system utilization
- **Reliability**
 - Still functioning well under failures

Outline

- Introduction & Motivation
- Problem Statement
- **Proposed Work**
- Evaluation
- Conclusions
- Future Work

Job Scheduling Systems

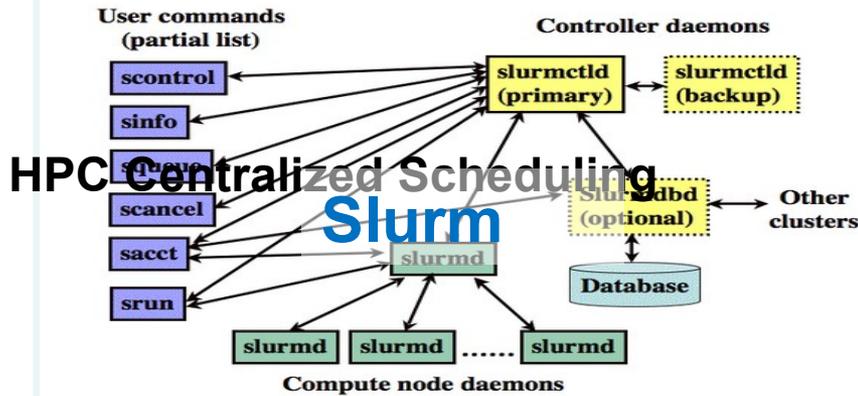
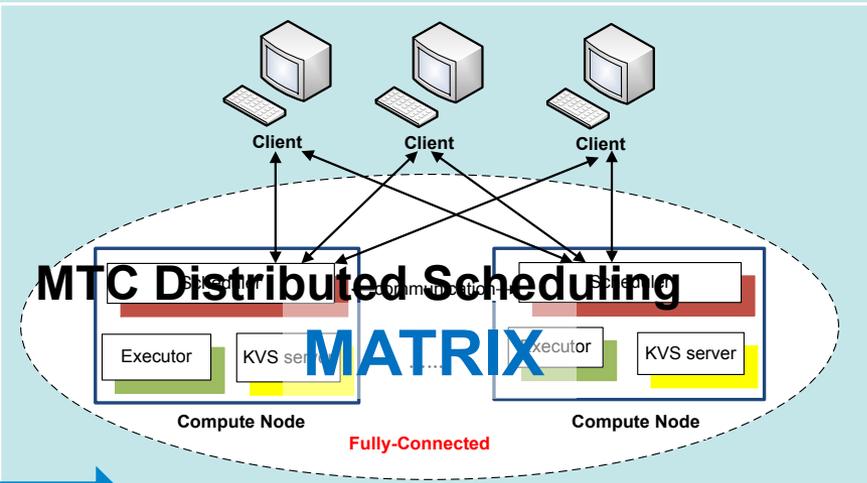
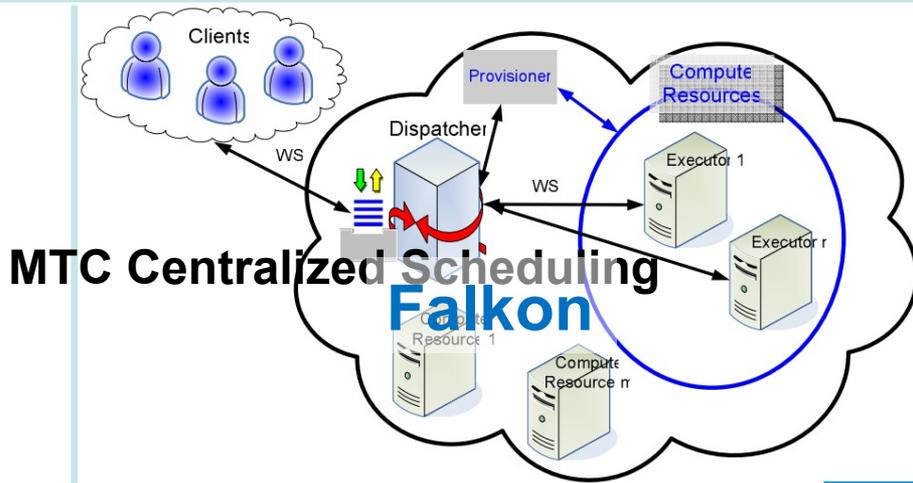
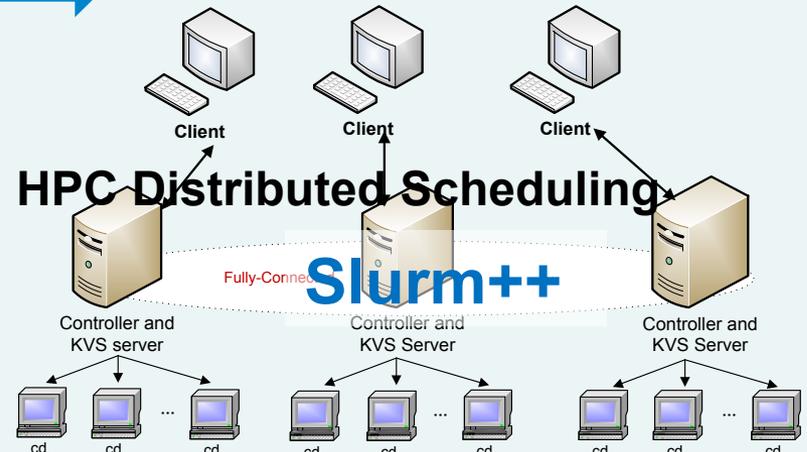
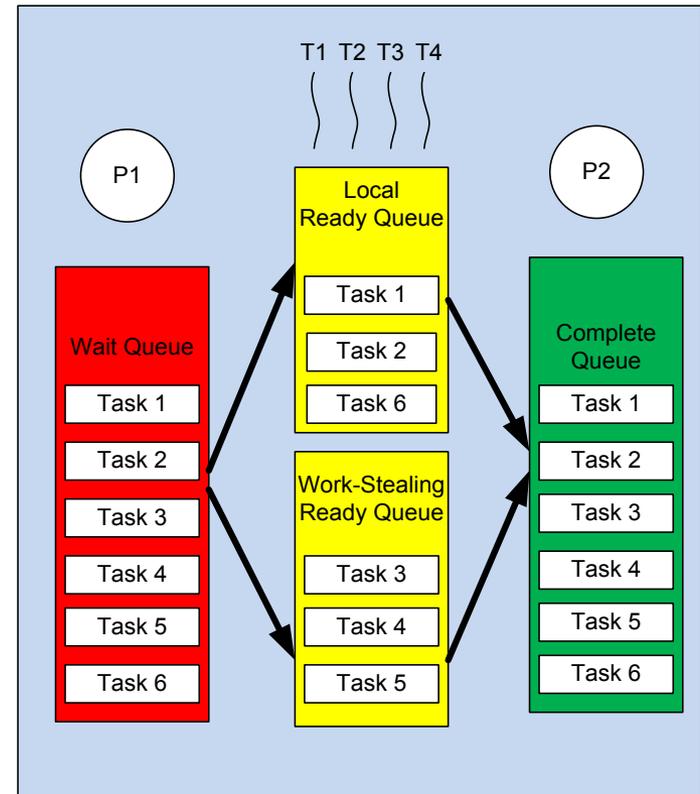
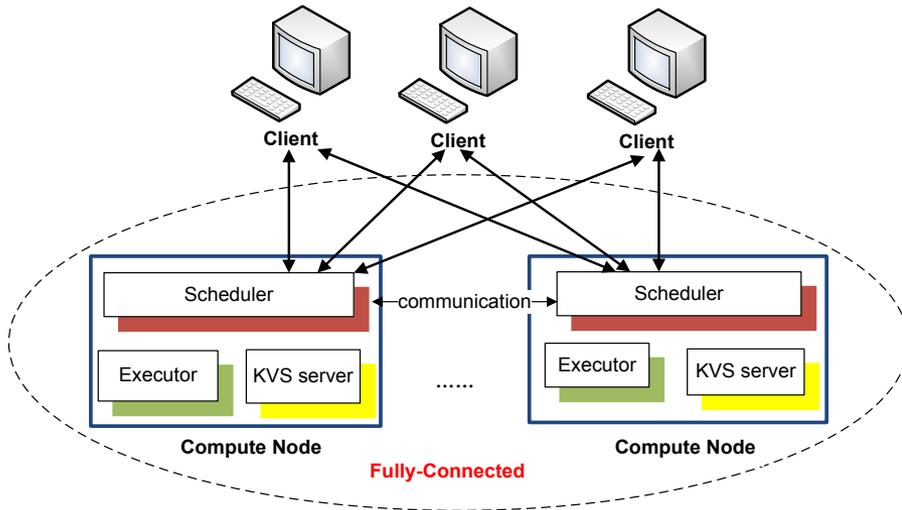


Figure 1. SLURM components

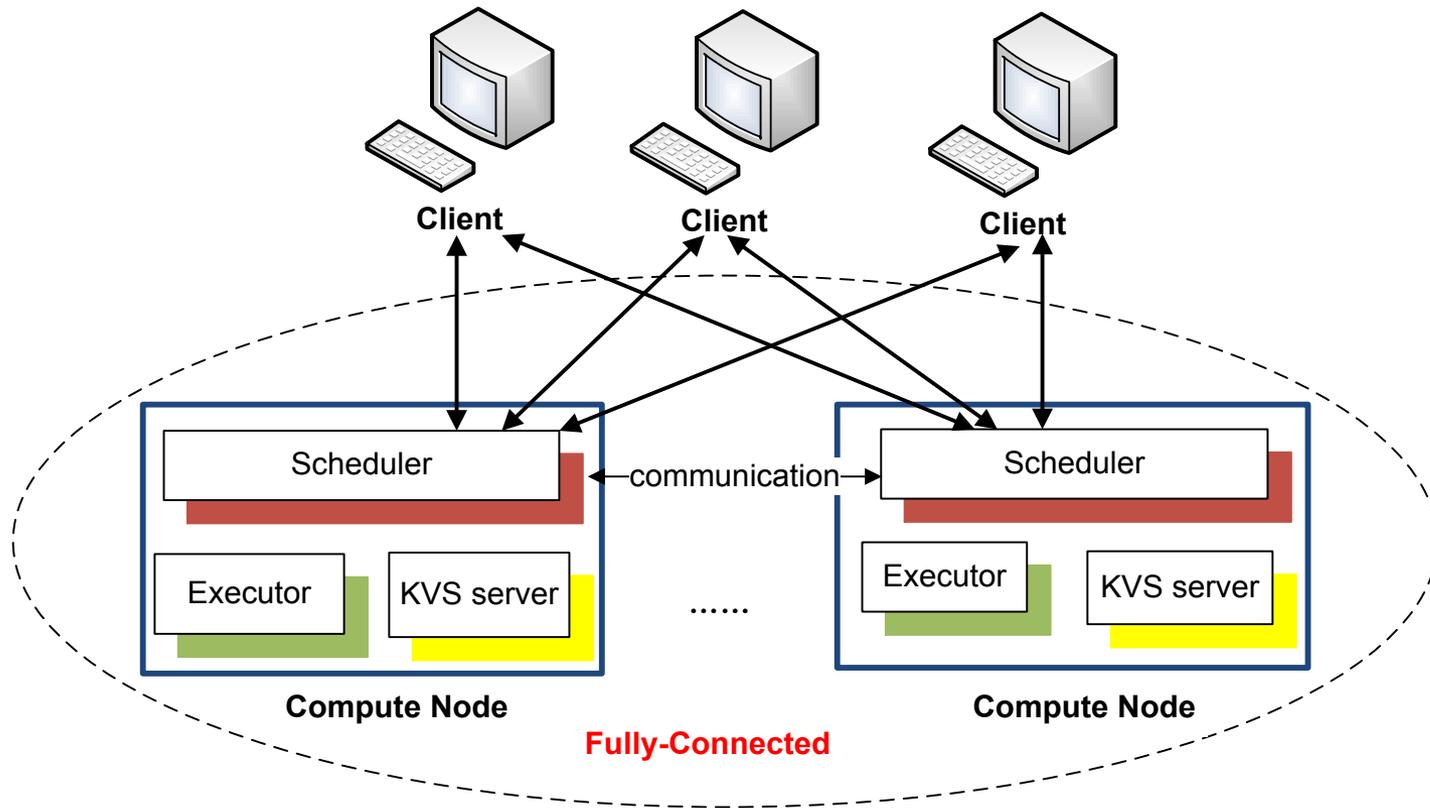


MTC Distributed Scheduling



Scheduler Specification

MTC Scheduling Architecture

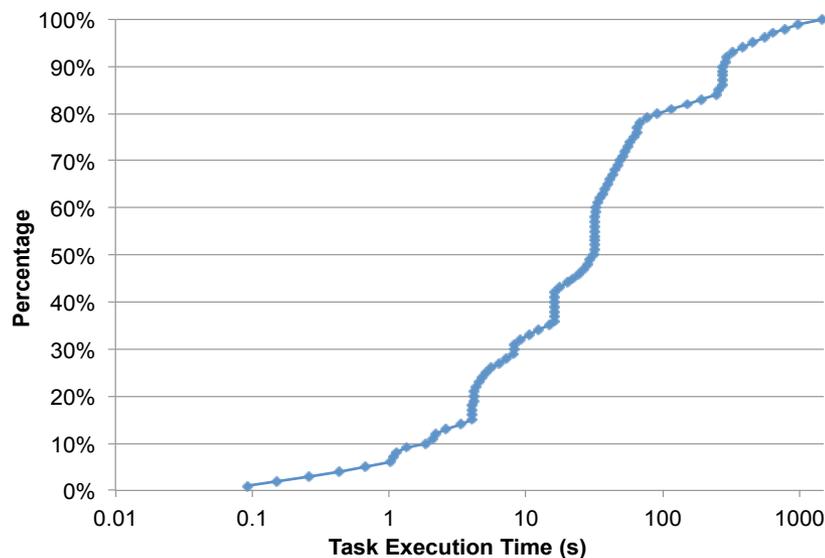


MTC Distributed Scheduling Challenges

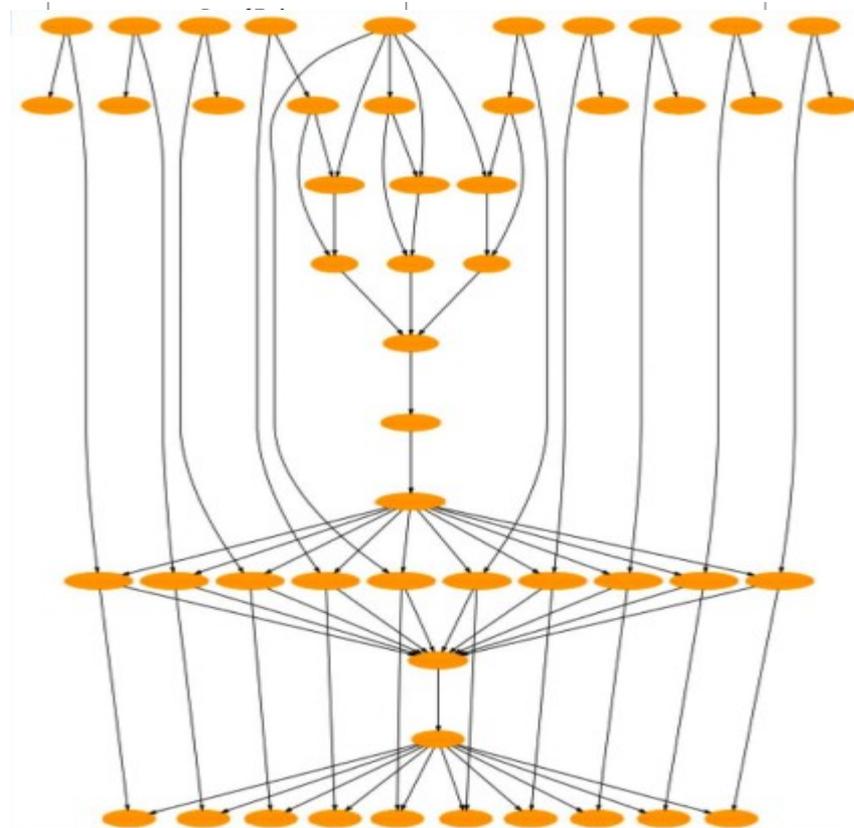
- **Fine-grained Workloads**
- **Load Balancing**
- **Data-Aware Scheduling**

Fine-grained Workloads

MTC application trace



MTC application DAGs



- **17-month period on IBM BG/P machine**
 - Minimum length: 0 sec
 - Maximum length: 1470 sec
 - Medium length: 30 sec
 - Average length: 95 sec

Load Balancing

Steal tasks

- **Load balancing**

- Distribute workloads as evenly as possible

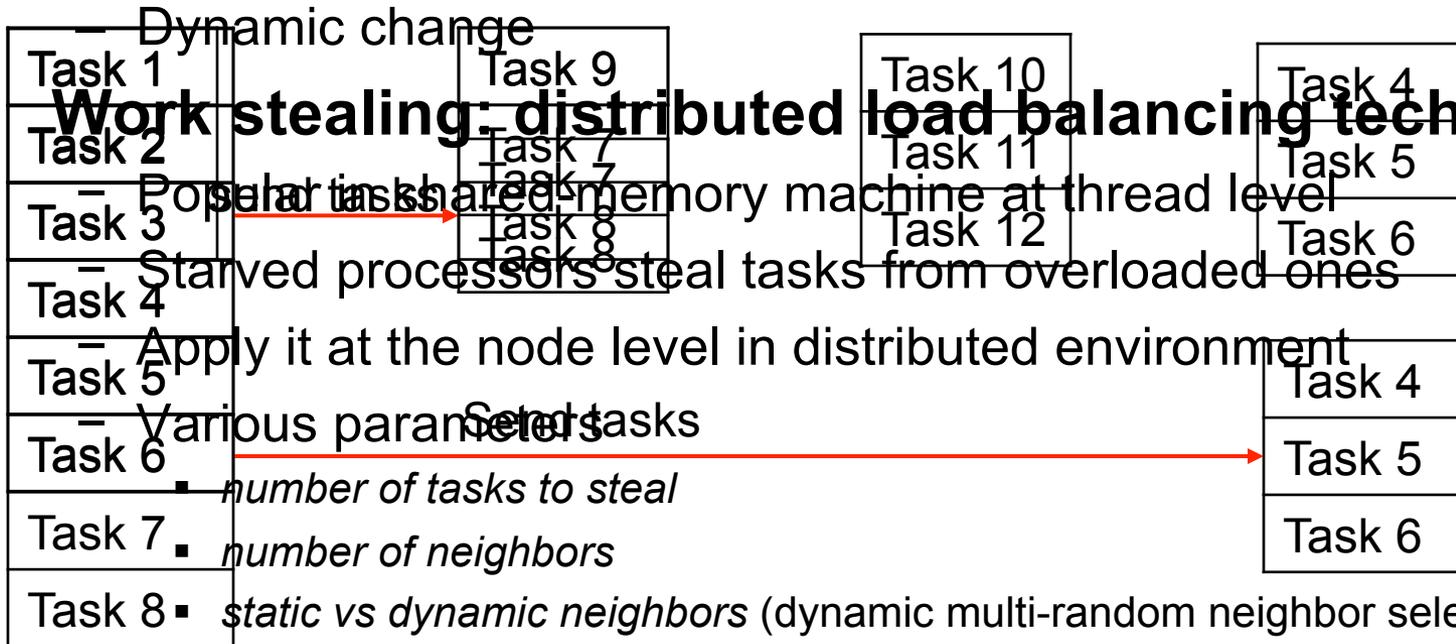
Scheduler 1

Scheduler 2

Scheduler 3

Scheduler 4

- **Work stealing: distributed load balancing technique**



- *number of tasks to steal*

- *number of neighbors*

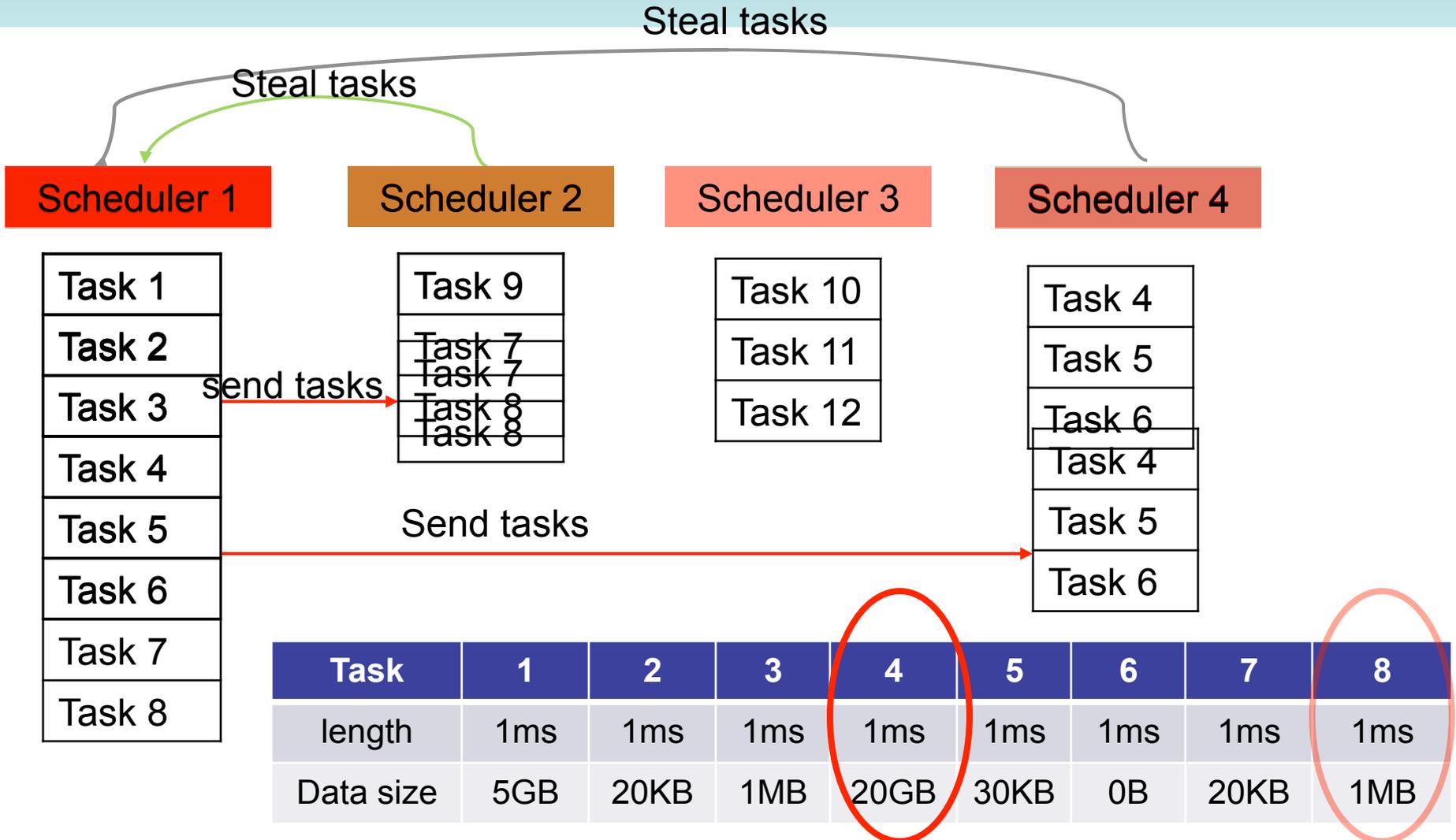
- *static vs dynamic neighbors (dynamic multi-random neighbor selection)*

- *polling interval (exponential back-off)*

Data-Aware Scheduling

- **Big-data era**
 - data-intensive applications
 - Data-flow driven programming models (MTC)
 - Workflow, task execution framework
- **Work Stealing**
 - Original work stealing is data locality-oblivious
 - Migrating tasks randomly comprise data-locality
 - Propose a Data-aware Work Stealing technique

Locality-Oblivious Work Stealing

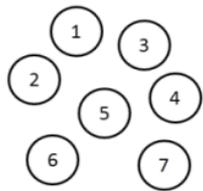


Data-Aware Work Stealing

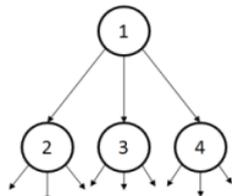
```

typedef TaskMetaData
{
    int num_wait_parent; // nu
    vector<string> parent_list; //
    vector<string> data_object;
    vector<long> data_size; // d
    long all_data_size; // all data
    vector<string> children; // c
}
    
```

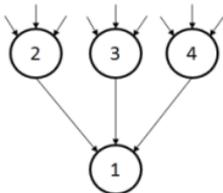
Bag of Tasks



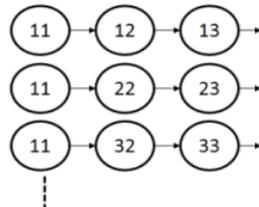
Fan-Out DAG



Fan-In DAG



Pipeline DAG



Wait Queue (WaitQ): holds tasks that are waiting for parents to complete

Dedicated Local Ready Queue (LReadyQ): holds ready tasks that can only be executed on local node

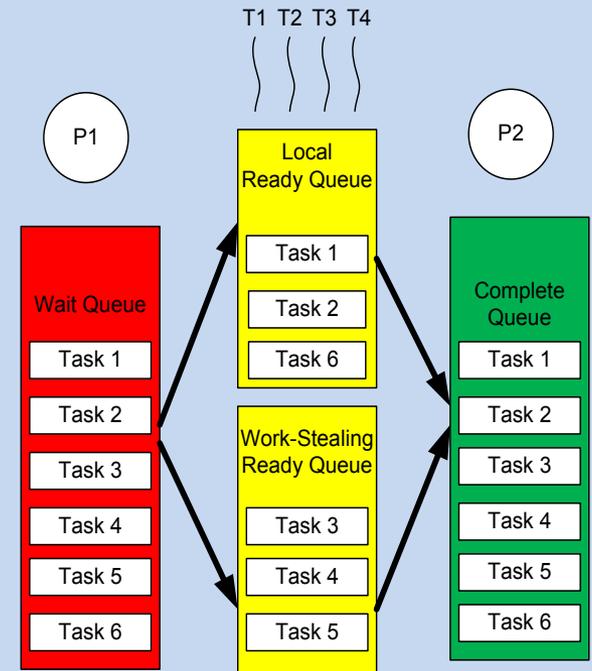
Shared Work-Stealing Ready Queue (SReadyQ): holds ready tasks that can be shared through work stealing

Complete Queue (CompleteQ): holds tasks that are completed

P1: a program that checks if a task is ready to run, and moves ready tasks to either ready queue according to the decision making algorithm

P2: a program that updates the task metadata for each child of a completed task

T1 to T4: executor has 4 (configurable) executing threads that executes tasks in the ready queues and move a task to complete queue when it is done



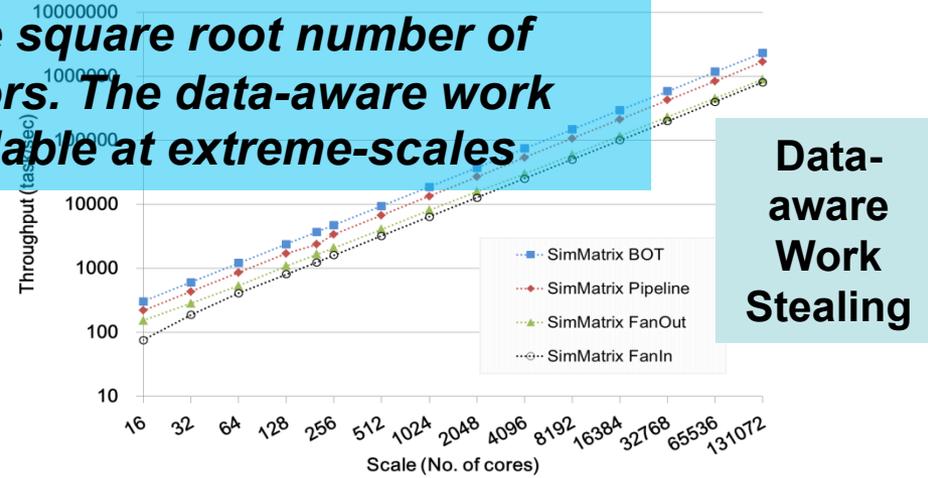
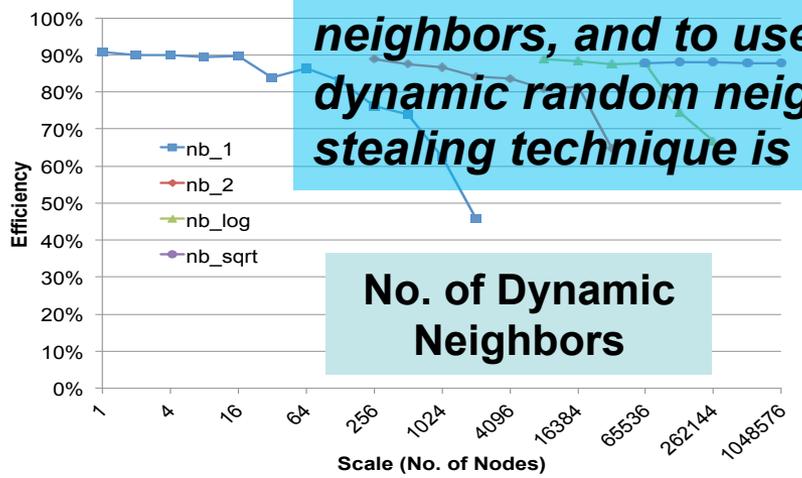
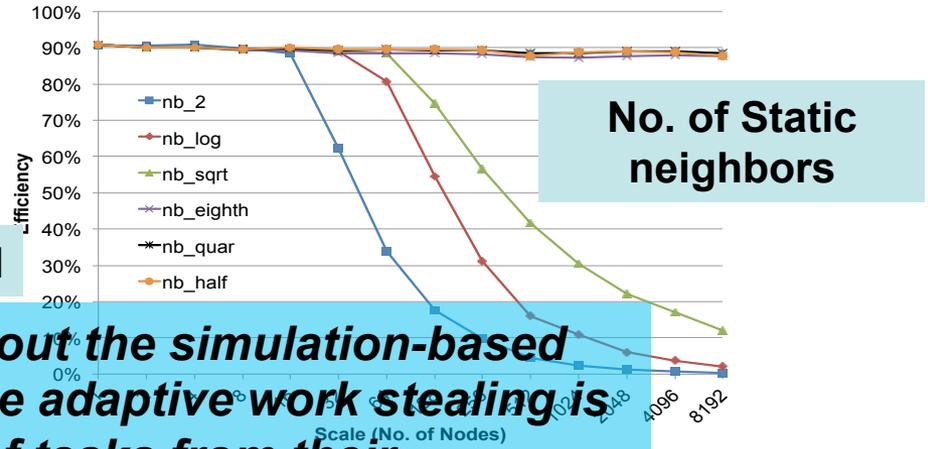
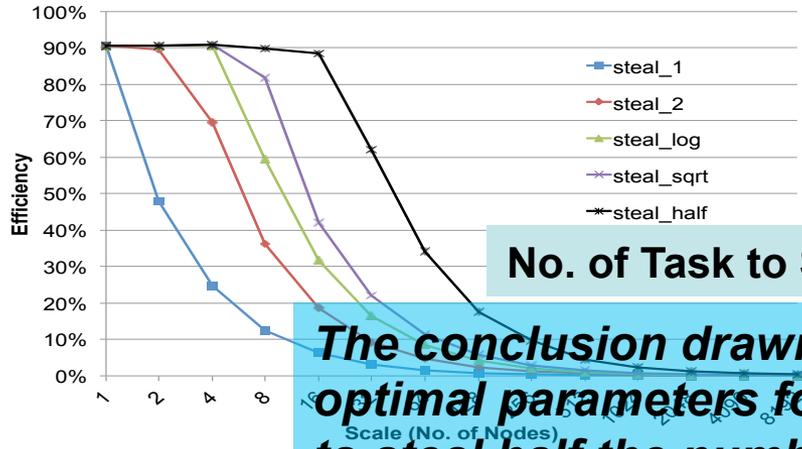
Scheduling Policies

- **MLB**
 - Maximized Load Balancing
- **MDL**
 - Maximized Data-Locality
- **RLDS**
 - Rigid Load balancing and Data-locality Segregation
- **FLDS**
 - Flexible Load balancing and Data-locality Segregation

SimMatrix: SIMulator for MAny-Task computing execution fabRlc at eXascale

- **A discrete event simulator**
 - Simulates MTC task execution framework
 - 1M nodes, 1B cores, 100B tasks
- **Explore the scalability**
 - fully distributed MTC-architecture
 - work stealing technique
 - data-aware work stealing technique

Exploration of Work Stealing



The conclusion drawn about the simulation-based optimal parameters for the adaptive work stealing is to steal half the number of tasks from their neighbors, and to use the square root number of dynamic random neighbors. The data-aware work stealing technique is scalable at extreme-scales

Simulations have shown that the distributed architecture and the work stealing technique are scalable, now we can do real implementations

MATRIX: TASK EXECUTION FRAMEWORK FOR MTC

MATRIX: MAny-Task computing execution fabRlc at eXascale

- **MATRIX components**

- Client (submit tasks, monitoring execution progress)
- Scheduler (scheduling tasks for execution)
- Executor (execute tasks)

- **MATRIX code (<https://github.com>)**

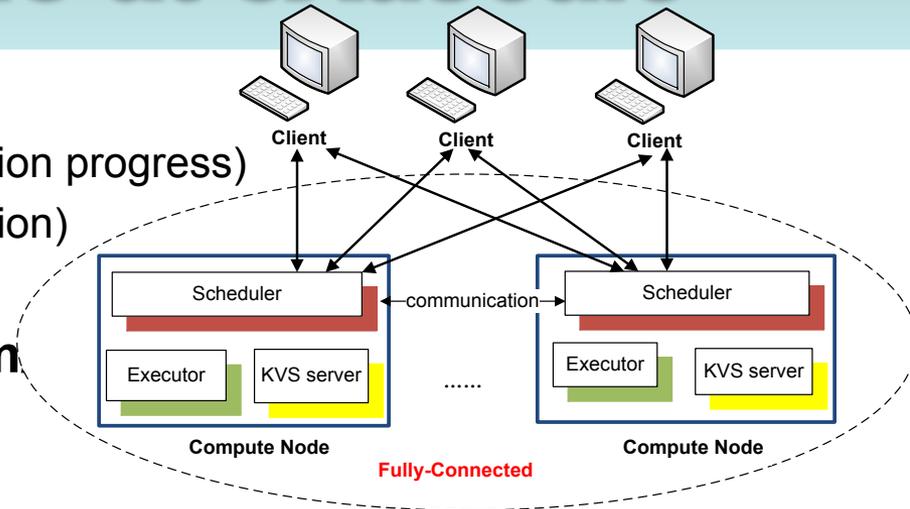
- 5K lines of C++ code
- 8K lines of ZHT code
- 1K lines of auto-generated code from Google protocol buffer

- **MATRIX goal**

- Scalable distributed scheduling system for MTC applications
- The prototype has been finished
- working on running real applications

- **Testbeds:**

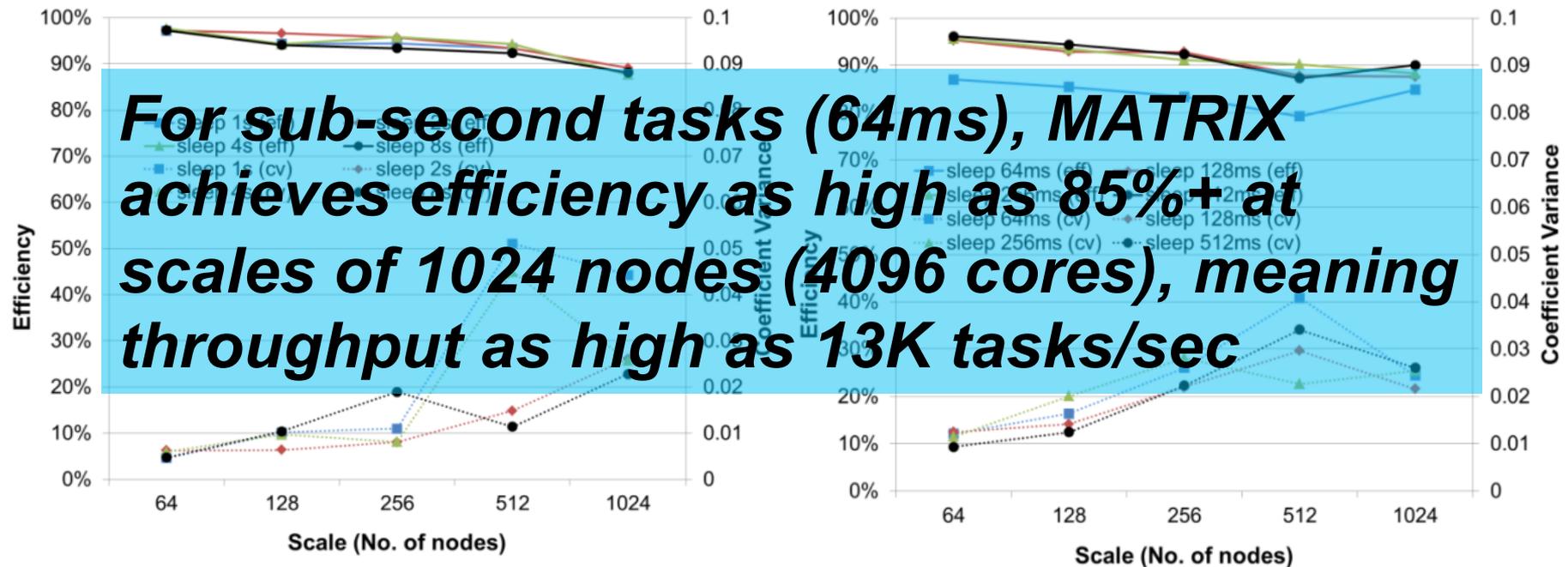
- IBM Blue Gene/P/Q supercomputers (4K cores)
- Probe Kodiak cluster (200 cores)
- **Amazon EC2 (256 cores)**



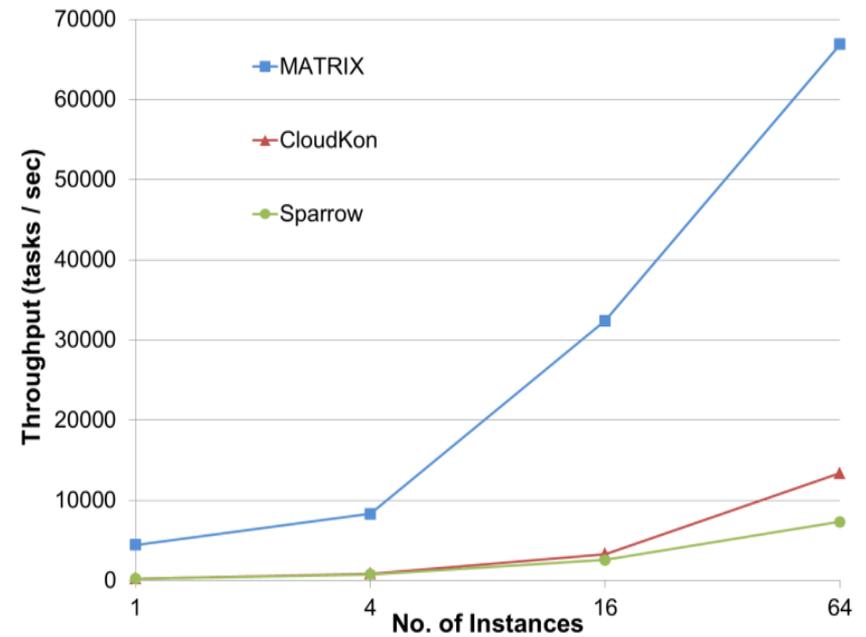
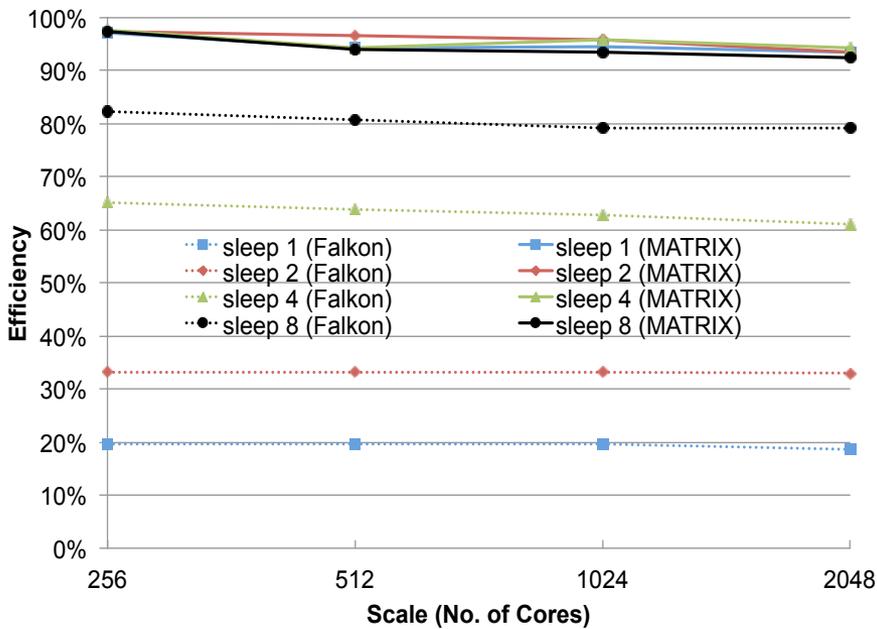
Outline

- Introduction & Motivation
- Problem Statement
- Proposed Work
- **Evaluation**
- Conclusions
- Future Work

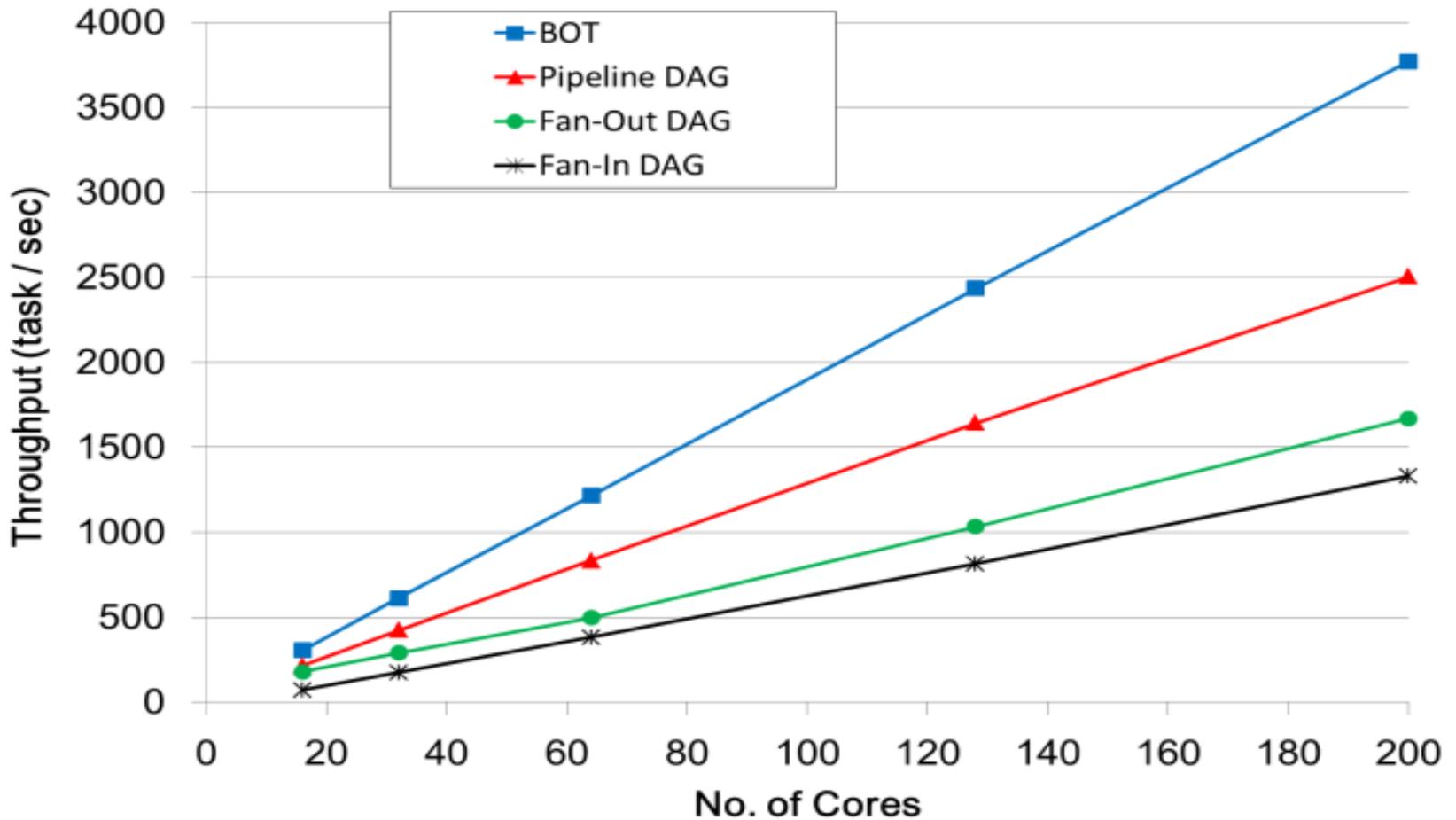
MATRIX running fine-grained tasks



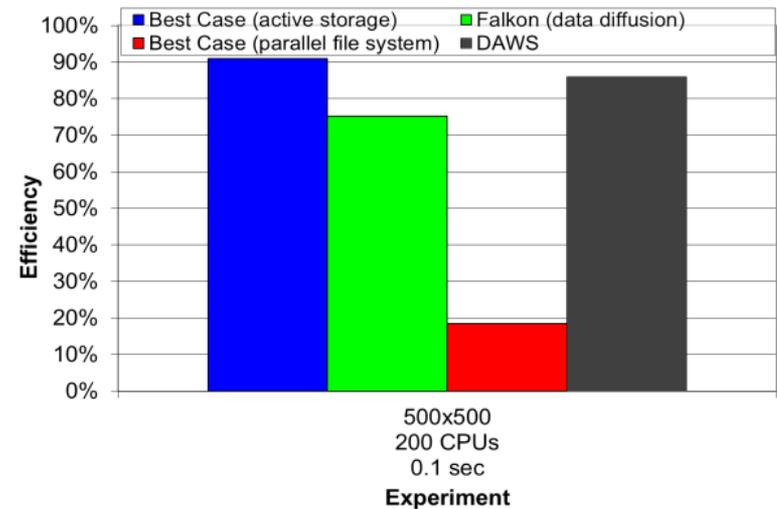
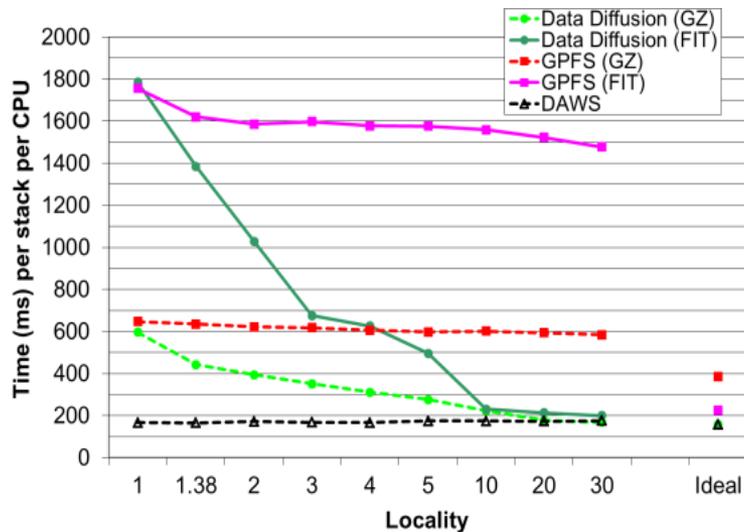
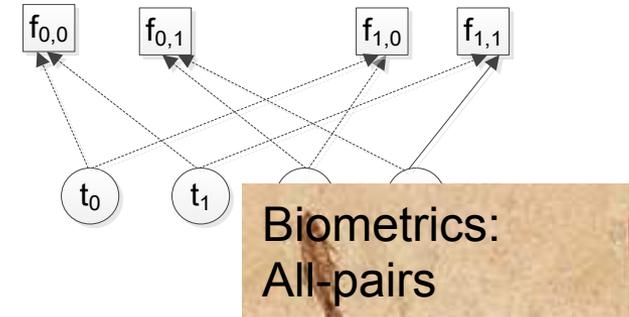
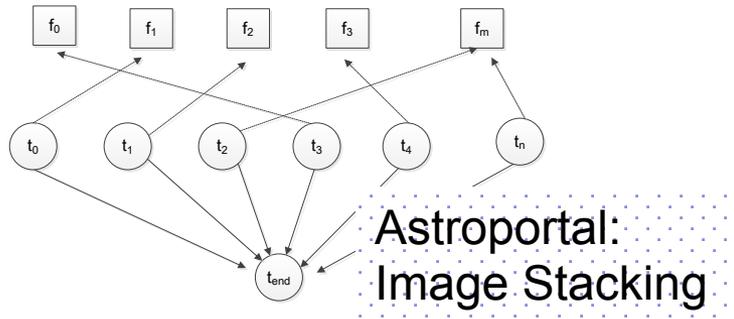
MATRIX comparisons



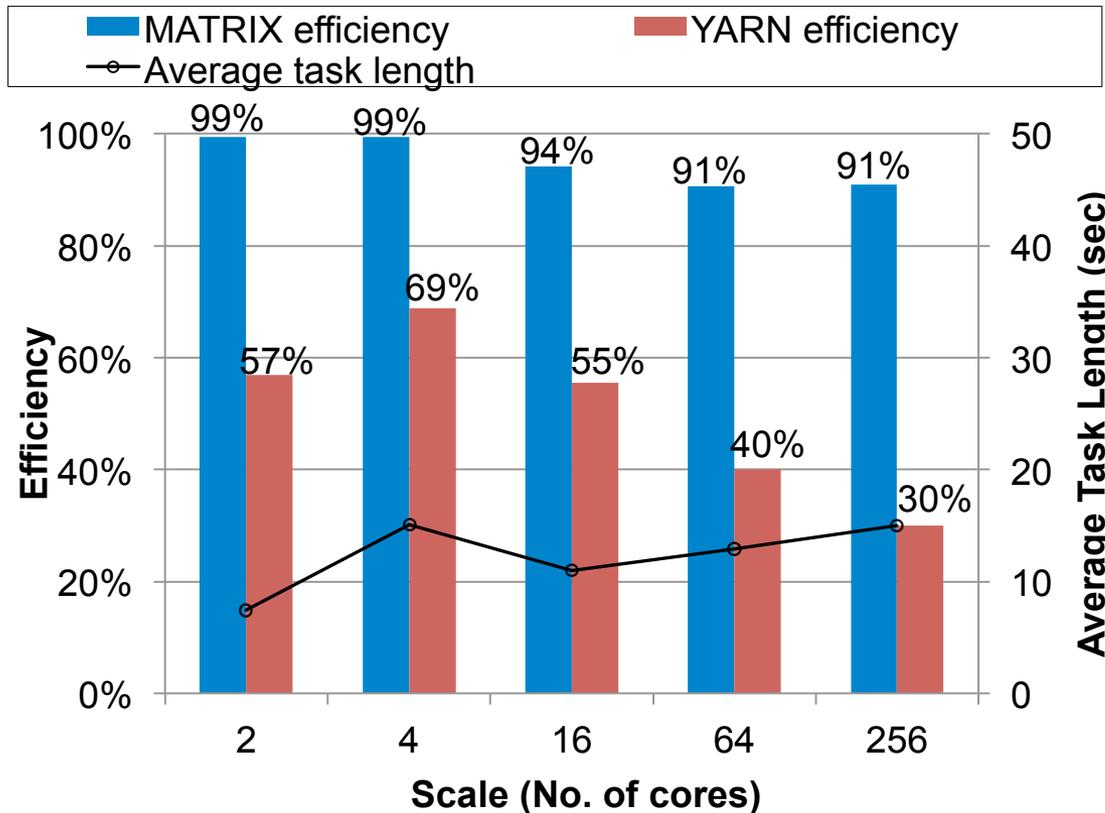
MATRIX running workload DAGs



MATRIX running data-intensive applications



MATRIX vs. Hadoop YARN



- **Bio-Informatics Application**
- **Protein-ligand Clustering**
- **5 Phases of MapReduce Jobs**
- **256MB data per node**
- **First phase has the majority of tasks**

Outline

- Introduction & Motivation
- Problem Statement
- Proposed Work
- Evaluation
- **Conclusions**
- Future Work

Conclusions

- **System scale is approaching exascale**
- **Applications are becoming fine-grained**
- **Next-generation resource management systems need to be highly scalable, efficient and available**
- **Fully distributed architectures are scalable**
- **Data-aware work stealing technique is scalable**

Outline

- Introduction & Motivation
- Problem Statement
- Proposed Work
- Evaluation
- Conclusions
- **Future Work**

Future Work

- **Hadoop Integration of MATRIX**
 - Hadoop scheduler is centralized
 - Replace the Hadoop scheduler with MATRIX
- **Workflow integration of MATRIX**
 - Swift workflow system
 - Will enable the execution of real scientific applications
- **File System integration of MATRIX**
 - FusionFS distributed file system
 - Take care of the data storage and management
 - Expose the data to MATRIX

More Information

- More information:
 - <http://datasys.cs.iit.edu/~kewang/>
- Contact:
 - kwang22@hawk.iit.edu
- Questions?