



THE UNIVERSITY OF  
**CHICAGO**



# Managing and Executing Loosely-Coupled Large-Scale Applications on Clusters, Grids, and Supercomputers

**Ioan Raicu**

Distributed Systems Laboratory  
Computer Science Department  
University of Chicago

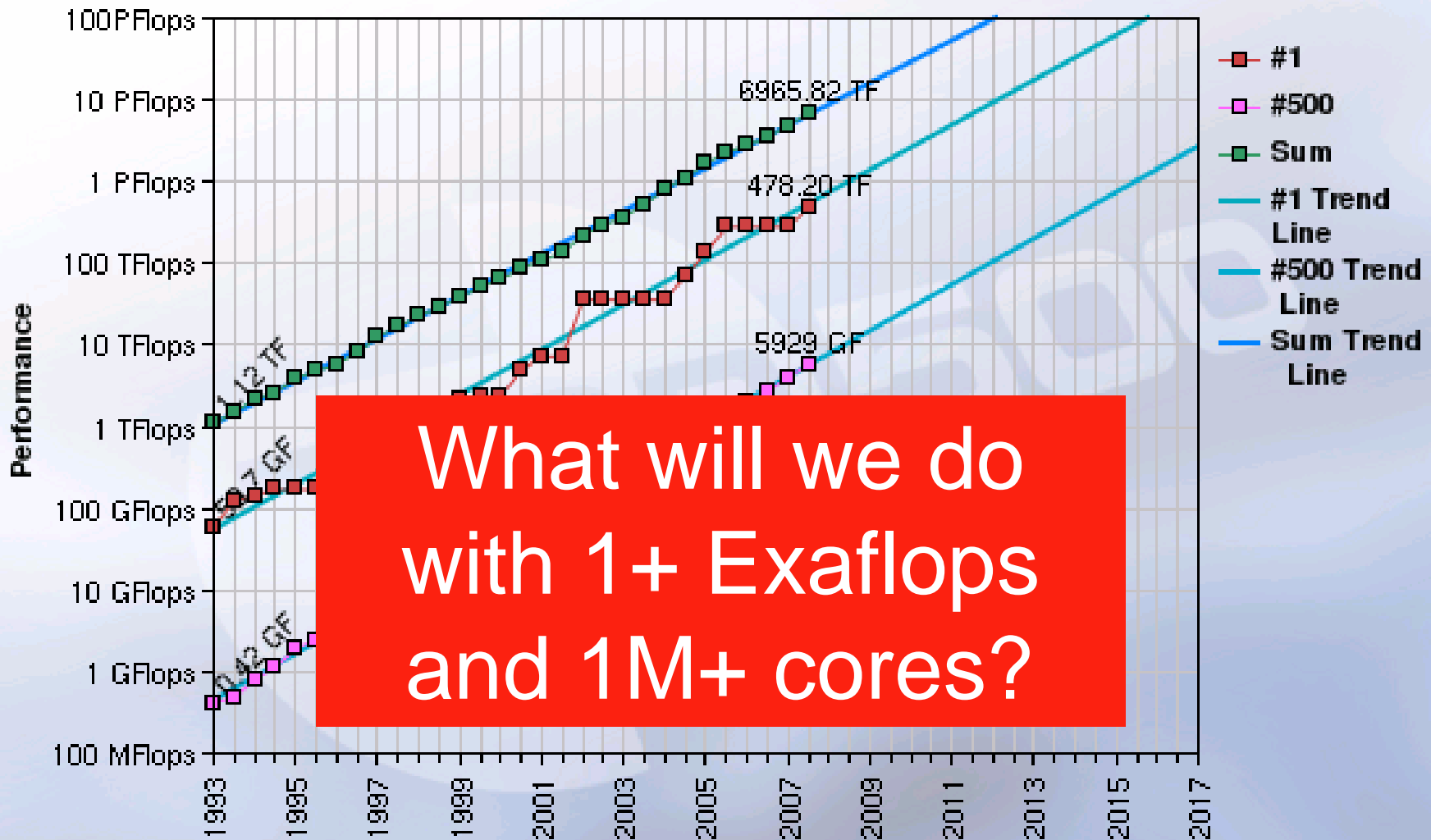
**Collaborators:**

Ian Foster (UC/CI/ANL), Yong Zhao (MS), Mike Wilde (CI/ANL),  
Zhao Zhang (CI), Alex Szalay (JHU), Jerry Yan (NASA/ARC), Catalin  
Dumitrescu (FANL), many others from Swift and Falcon teams



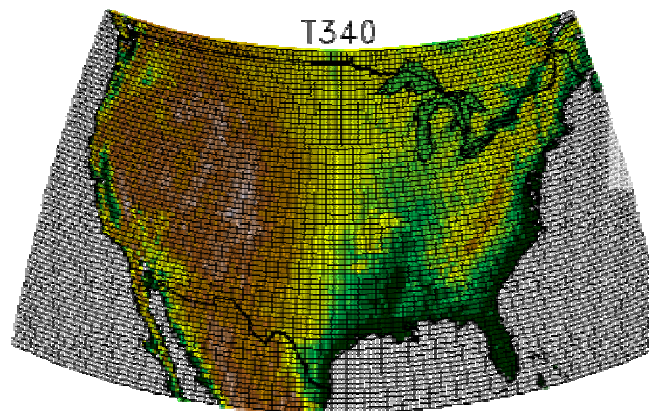
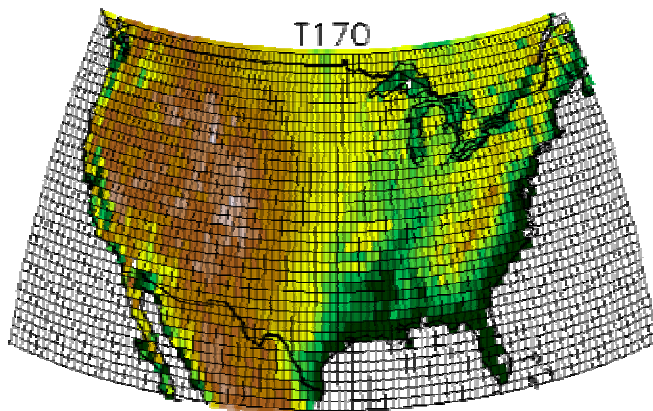
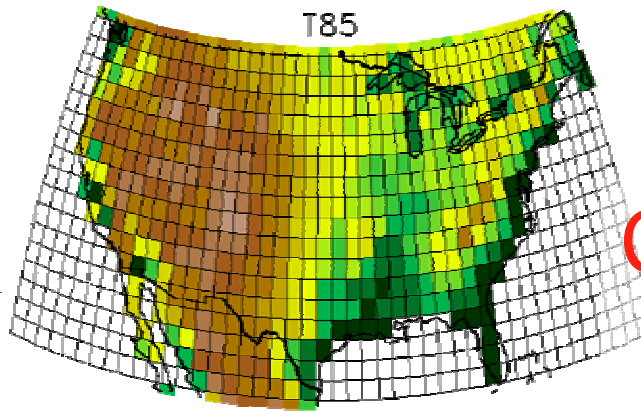
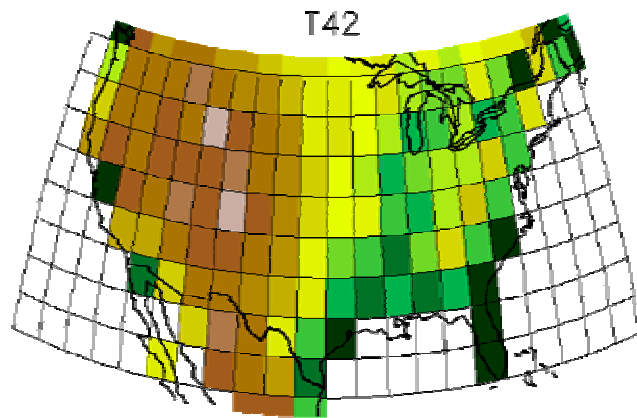
**GlobusWorld08**

May 14<sup>th</sup>, 2008



What will we do with 1+ Exaflops and 1M+ cores?

# 1) Tackle **Bigger and Bigger** Problems



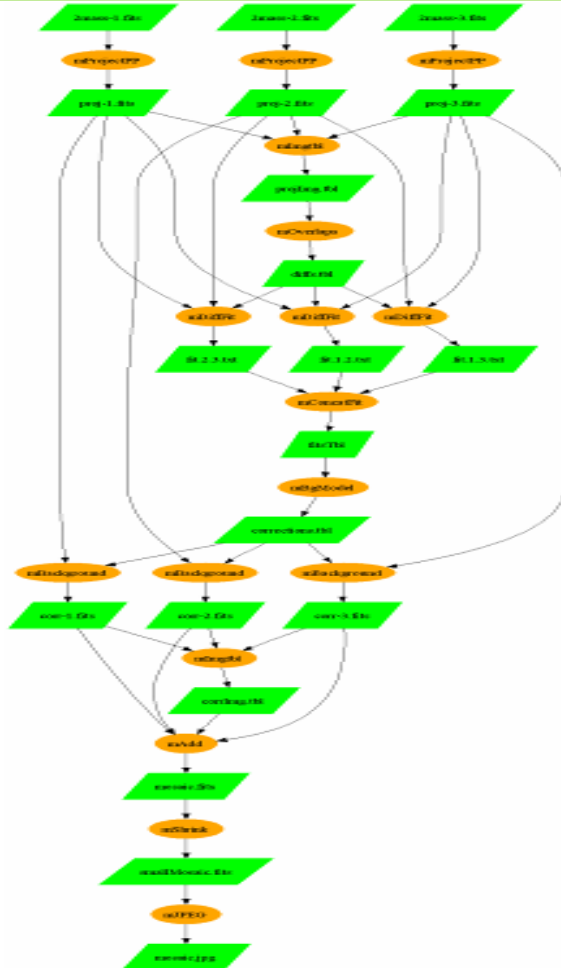
Computational  
Scientist  
as  
Hero



## 2) Tackle **Increasingly Complex Problems**



Computational  
Scientist  
as  
**Logistics  
Officer**

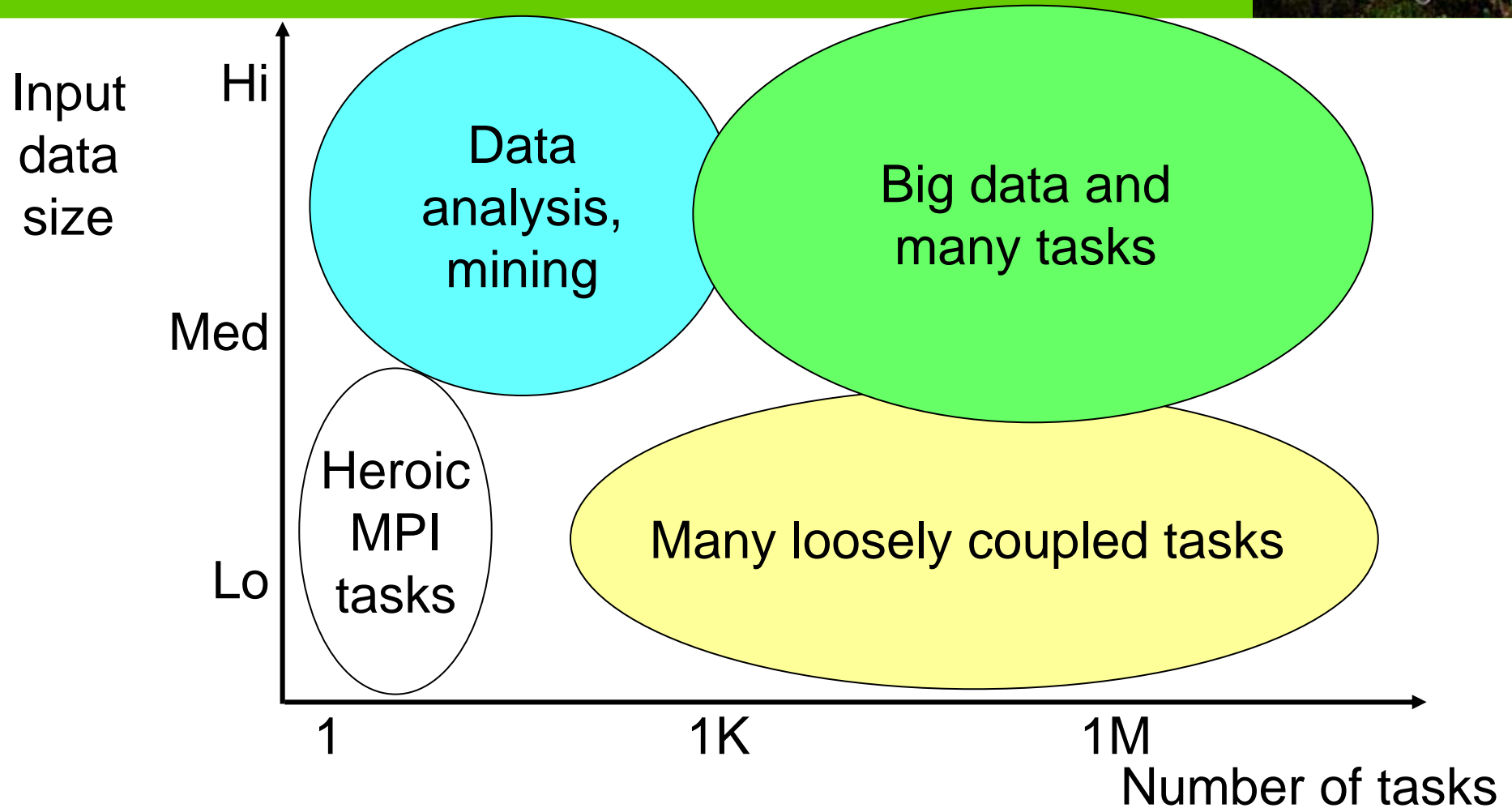


5/14/2008

Managing and Executing Loosely-Coupled Large-Scale Applications on Clusters, Grids, and Supercomputers

4

# Problem Types

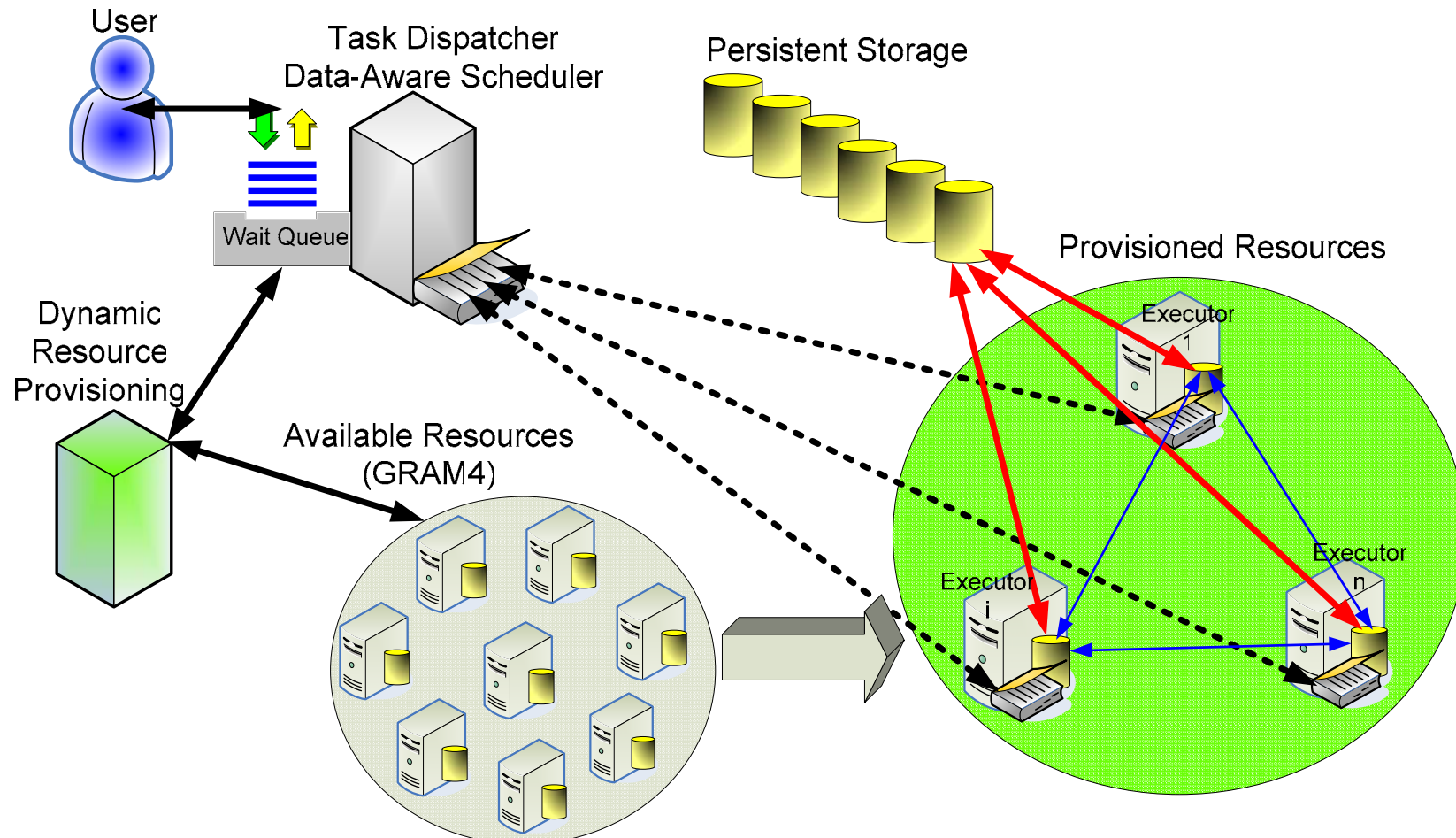


# Falkon: a Fast and Light-weight task executiON framework



- **Goal:** enable the *rapid and efficient* execution of many independent jobs on large compute clusters
- Combines three components:
  - a **streamlined task dispatcher** able to achieve order-of-magnitude higher task dispatch rates than conventional schedulers
  - **resource provisioning** through multi-level scheduling techniques
  - **data diffusion** and data-aware scheduling to leverage the co-located computational and storage resources

# Falkon Overview



5/14/2008

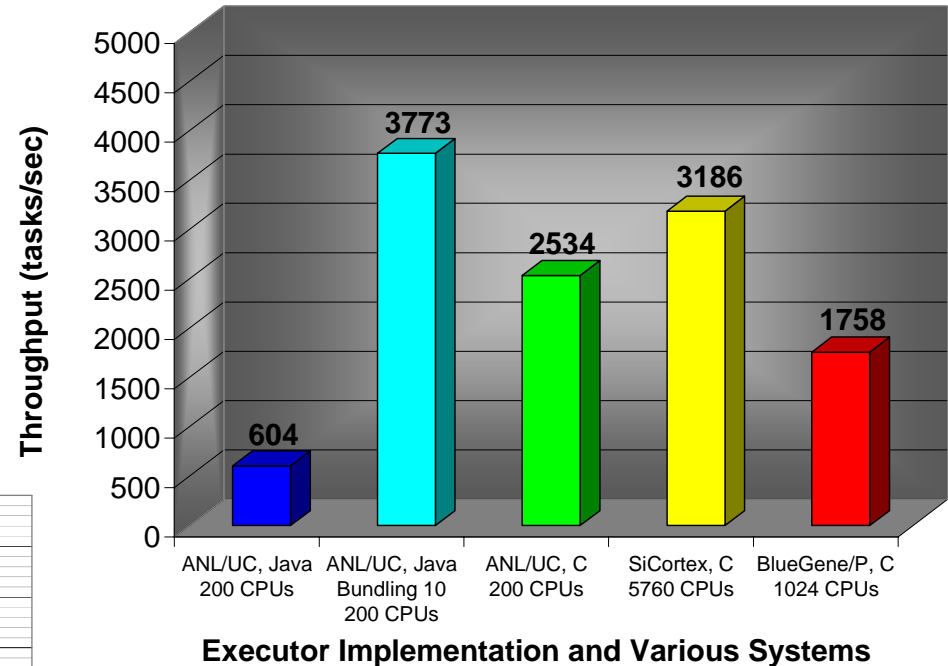
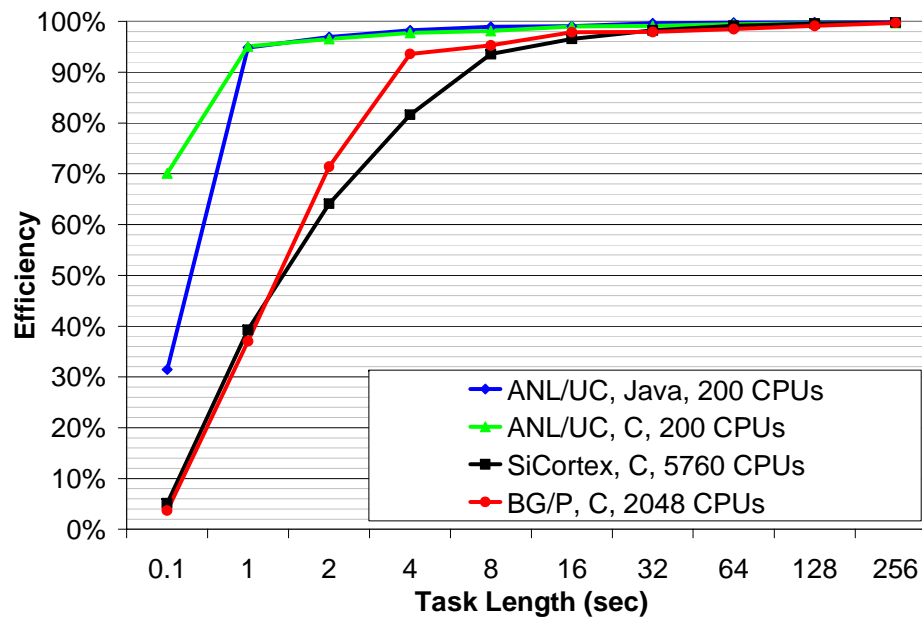
Managing and Executing Loosely-Coupled Large-Scale Applications on Clusters, Grids, and Supercomputers

7

# Dispatcher Throughput



- Fast:
  - Up to 3700 tasks/sec
- Scalable:
  - 54,000 processors
  - 1,500,000 tasks queued



- Efficient:
  - High efficiency with second long tasks on 1000s of processors

-Coupled Large-Scale Applications on and Supercomputers

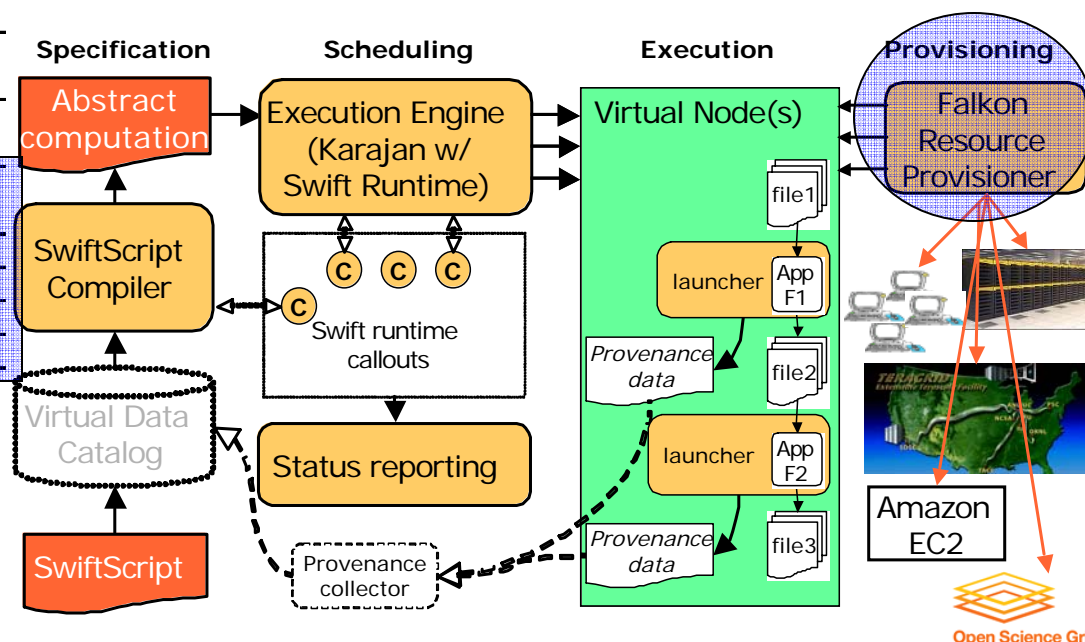


# Falkon Integration with Swift



Application	#Tasks/workflow	#Stages
ATLAS: High Energy Physics Event Simulation	500K	1
fMRI DBIC: AIRSN Image Processing	100s	12
FOAM: Ocean/Atmosphere Model	2000	3
GADU: Genomics	40K	4
HNL: fMRI Aphasia Study	500	4
NVO/NASA: Photorealistic Montage/Morphology	1000s	16
QuarkNet/I2U2: Physics Science Education	10s	3 ~ 6
RadCAD: Radiology Classifier Training	1000s	5
SIDGrid: EEG Wavelet Processing, Gaze Analysis	100s	20
SDSS: Coadd, Cluster Search	40K, 500K	2, 8
SDSS: Stacking, AstroPortal	10Ks ~ 100Ks	2 ~ 4
MolDyn: Molecular Dynamics	1Ks ~ 20Ks	8
MARS: Economic Modeling	1M~1B	1
DOCK: : Molecular Dynamics	1B	1

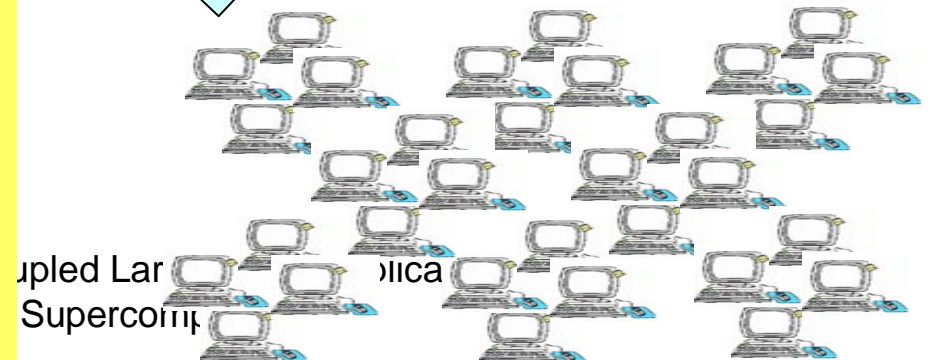
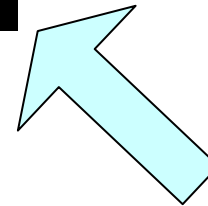
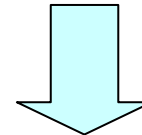
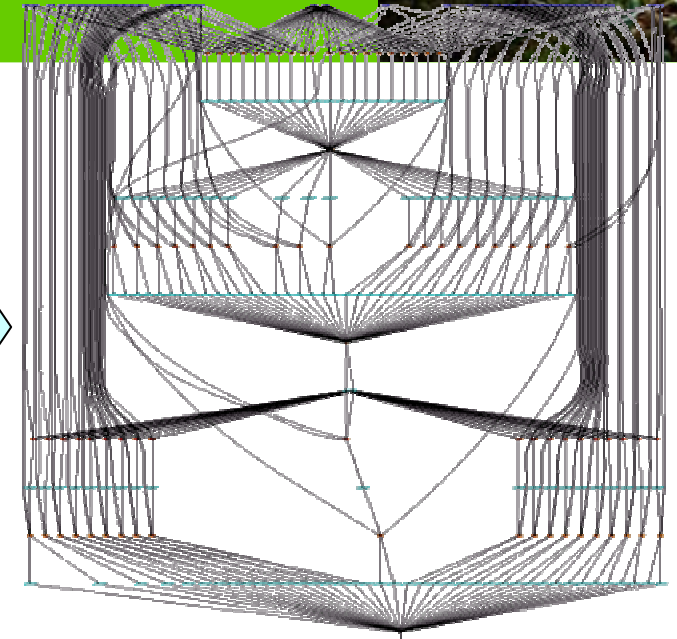
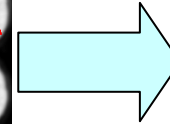
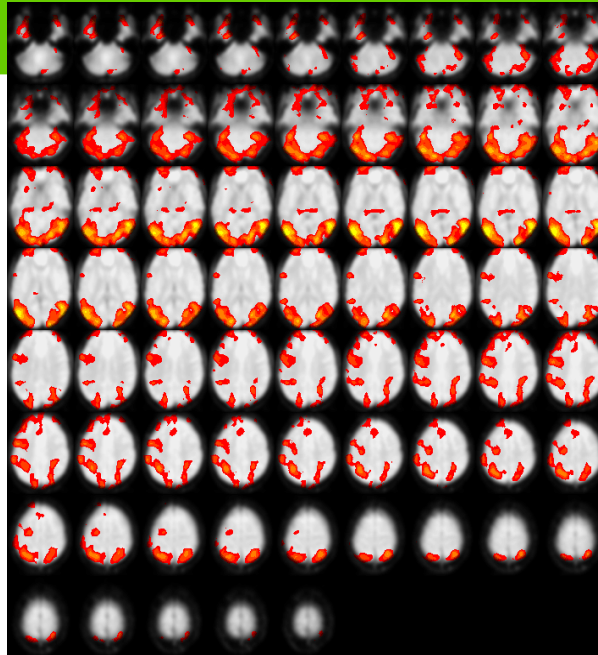
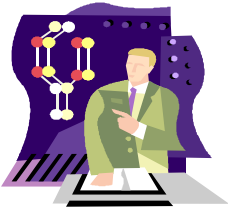
## Swift Architecture



5/14/2008

Managing and Executing Loosely-Coupled Large-Scale Applications on Clusters, Grids, and Supercomputers

# Functional MRI (fMRI)



upled Lar  
Supercomp

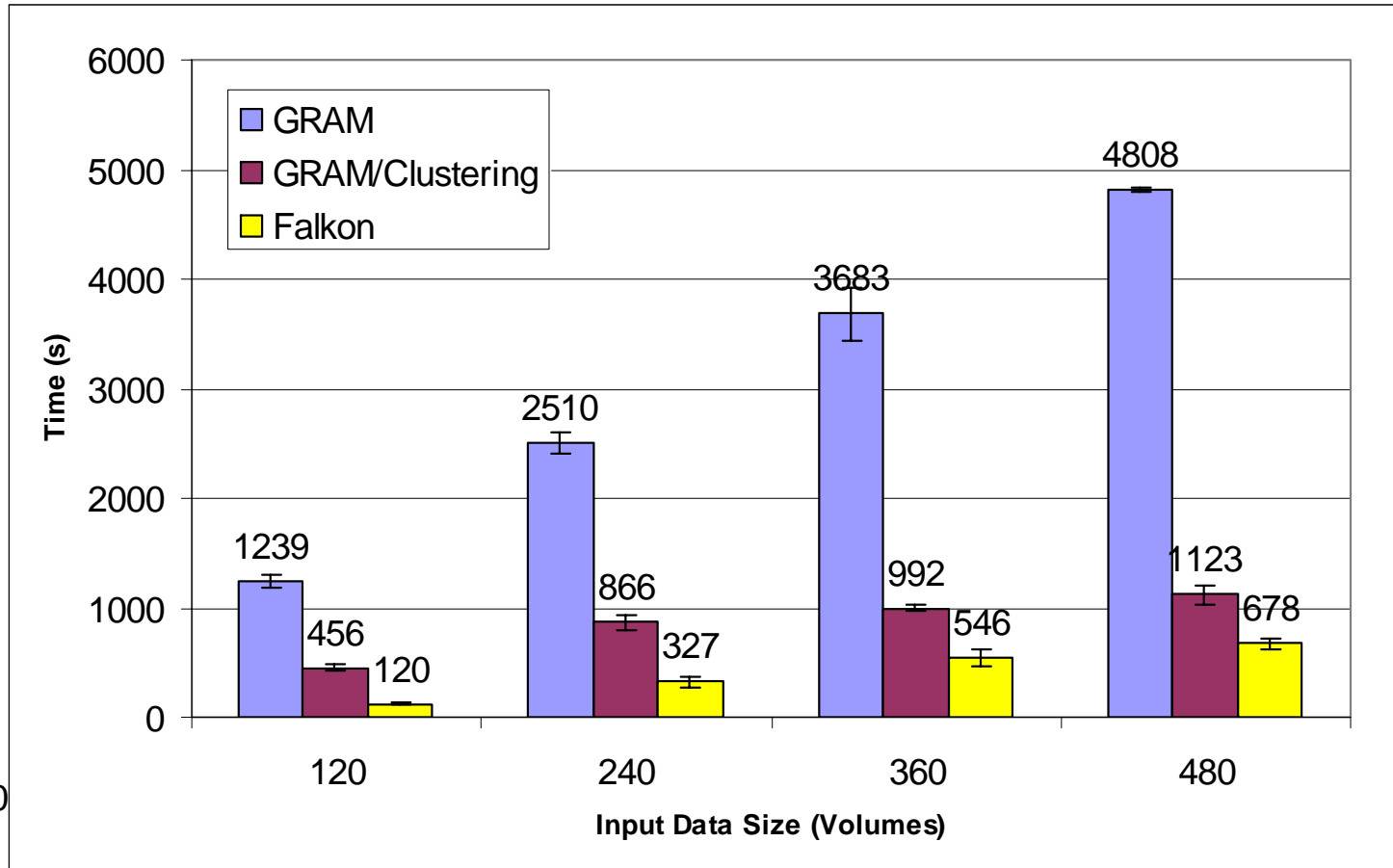
ica

- Wide range of analyses
  - Testing, interactive analysis, production runs
  - Data mining
  - Parameter studies

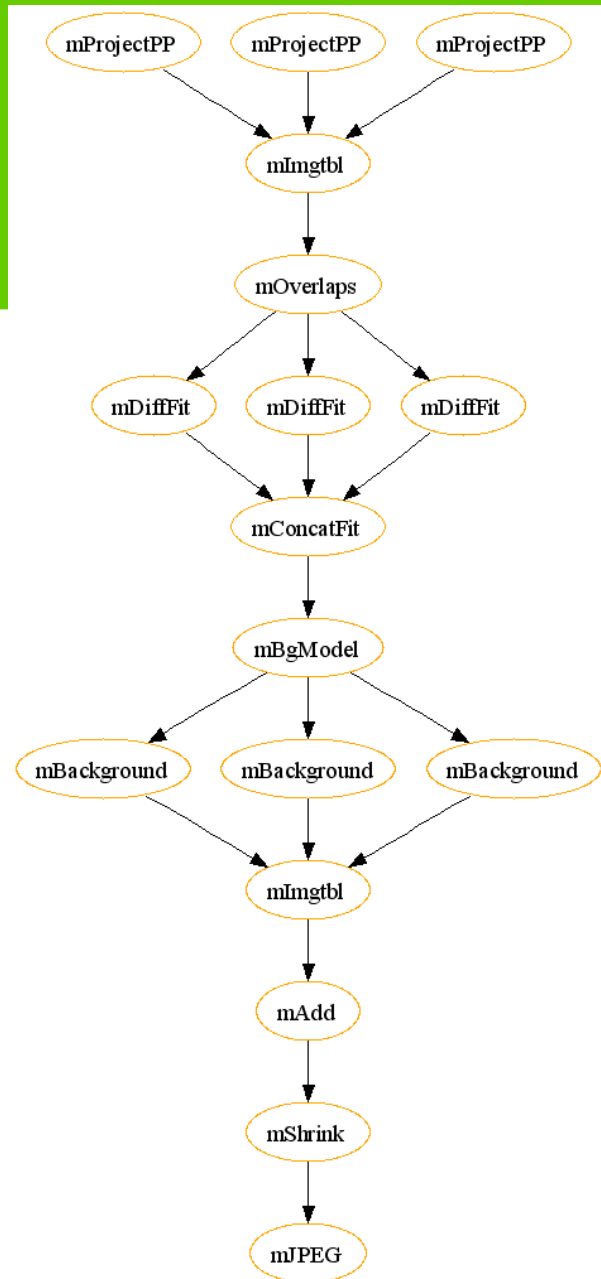
# fMRI Application



- GRAM vs. Falcon: 85%~90% lower run time
- GRAM/Clustering vs. Falcon: 40%~74% lower run time



5/14/20



B. Berriman, J. Good (Caltech)  
 J. Jacob, D. Katz (JPL)



5/14/2008

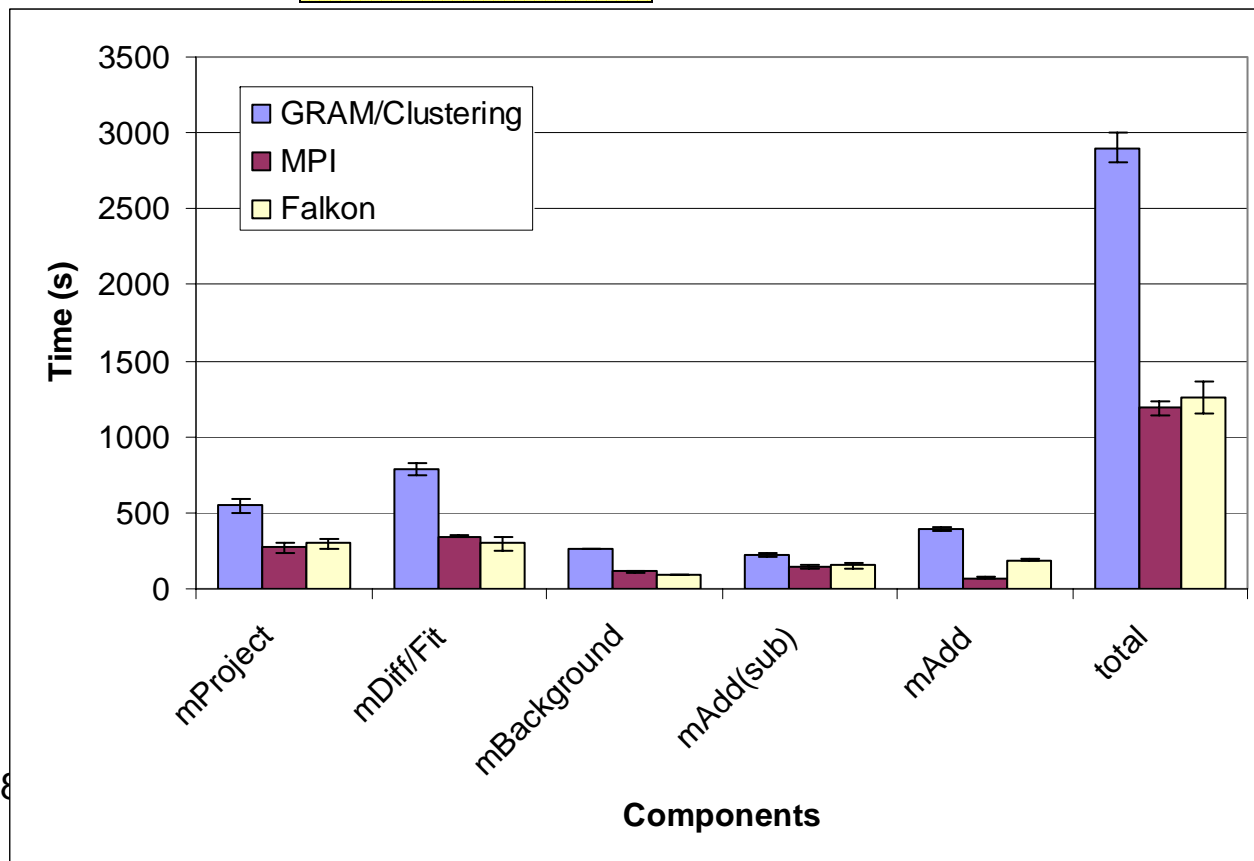
Managing and Executing Loosely-Coupled Large-Scale Applications on  
 Clusters, Grids, and Supercomputers

12

# Montage Application

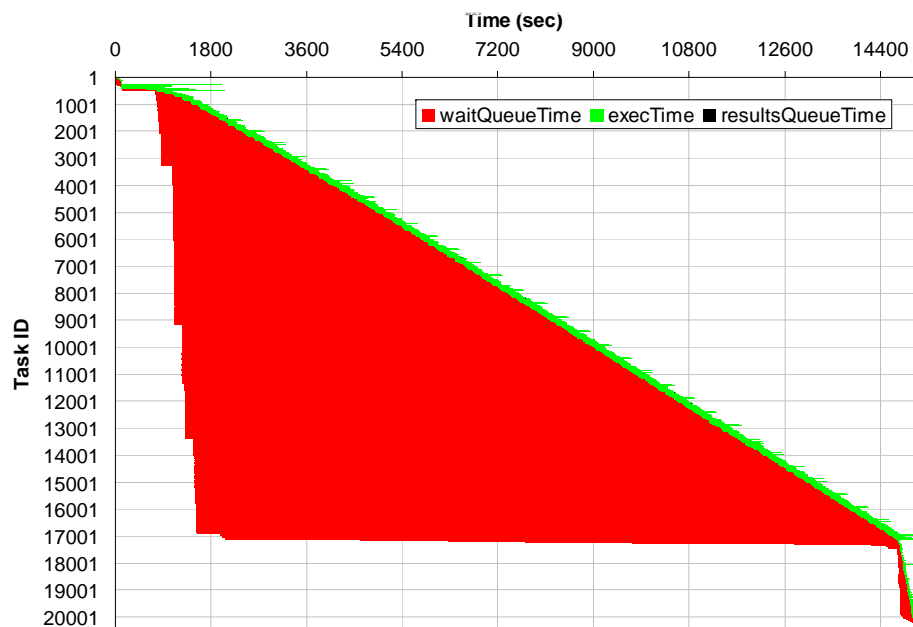


- GRAM/Clustering vs. Falcon: **57%** lower application run time
- MPI\* vs. Falcon: **4%** higher application run time
- \* MPI should be **lower bound**



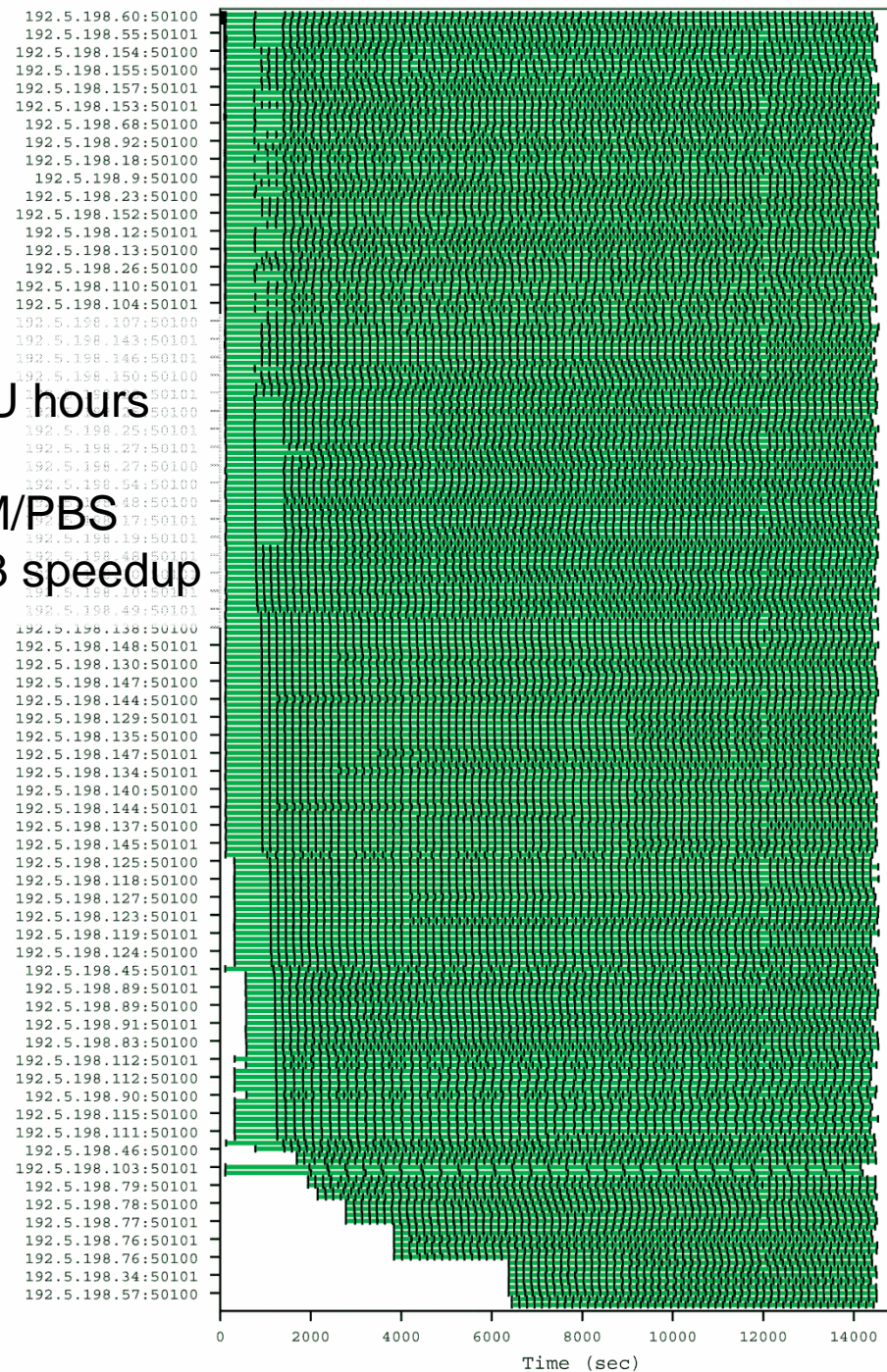
# MolDyn Application

- 244 molecules → 20497 jobs
- 15091 seconds on 216 CPUs → 867.1 CPU hours
- Efficiency: 99.8%
- Speedup: **206.9x → 8.2x** faster than GRAM/PBS
- 50 molecules w/ GRAM (4201 jobs) → 25.3 speedup



5/14/2008

Managing and Executing Loosely-Cc  
Clusters, Grids, and

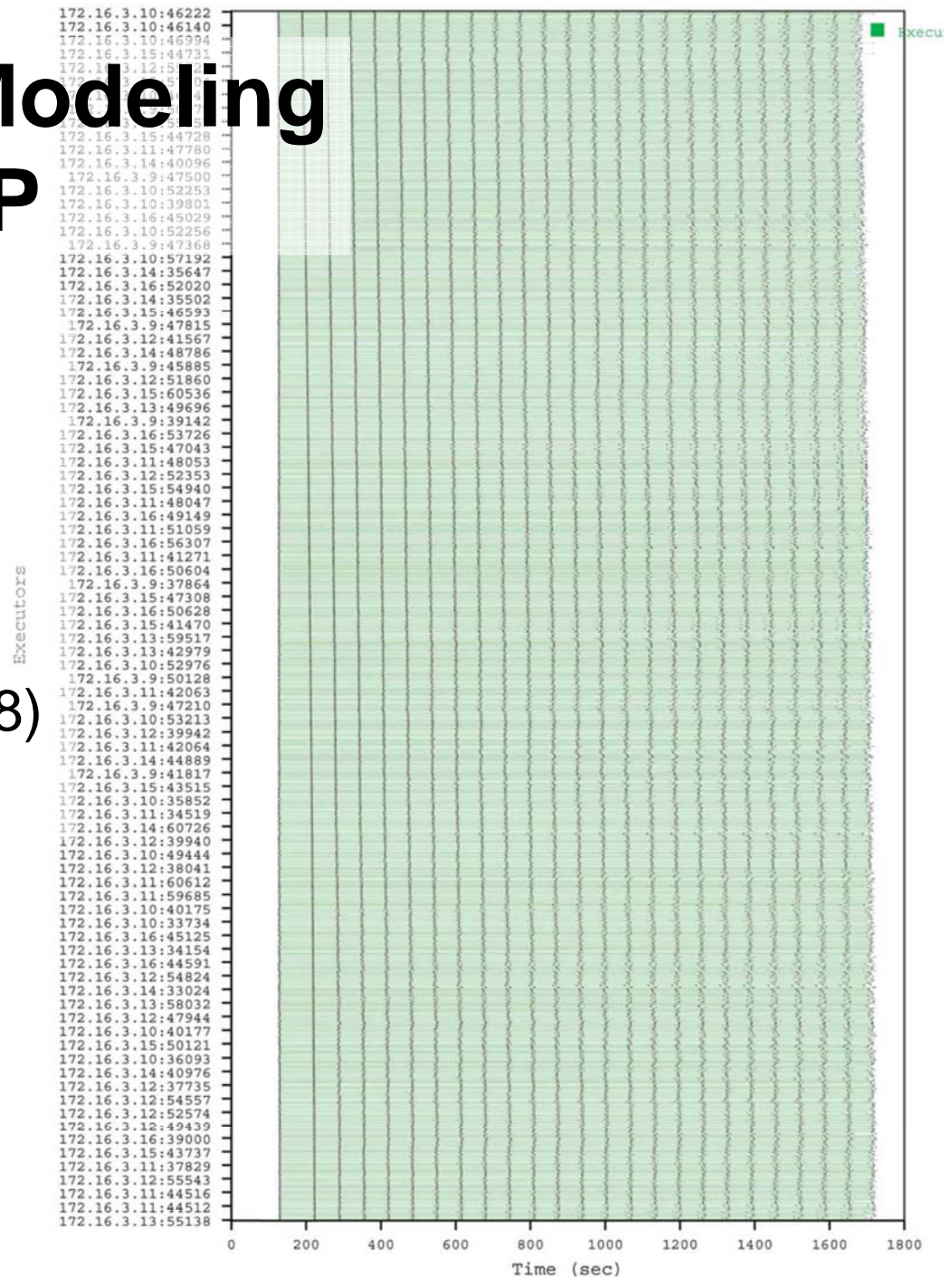


# MARS Economic Modeling on IBM BG/P

- CPU Cores: 2048
- Tasks: 49152
- Micro-tasks: 7077888
- Elapsed time: 1601 secs
- CPU Hours: 894
- Speedup: 1993X (ideal 2048)
- Efficiency: 97.3%



Operating Loc  
sters, Gr



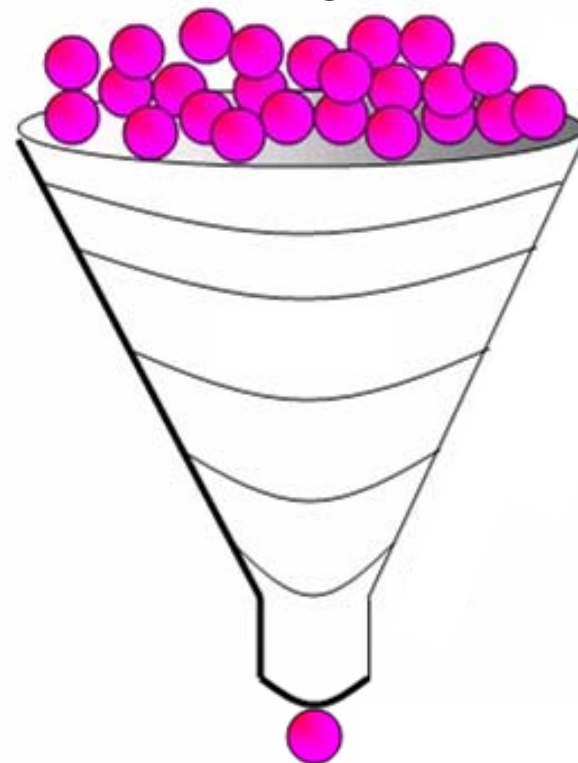
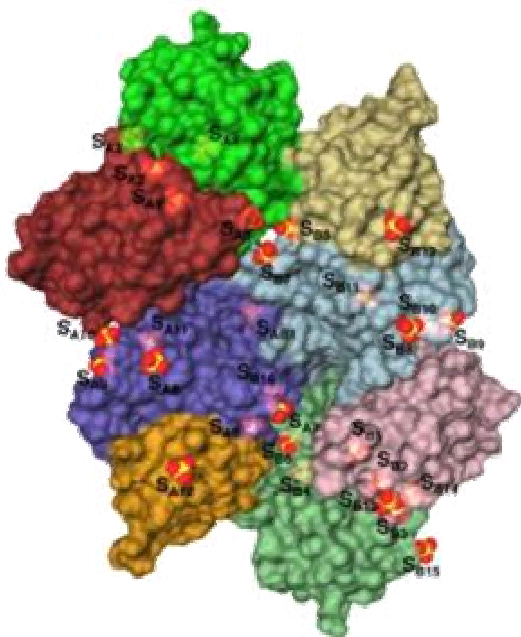
# Many Many Tasks: Identifying Potential Drug Targets



200+ Protein  
target(s)

x

5M+ ligands



(Mike Kubal, Benoit Roux, and others)

5/14/2008

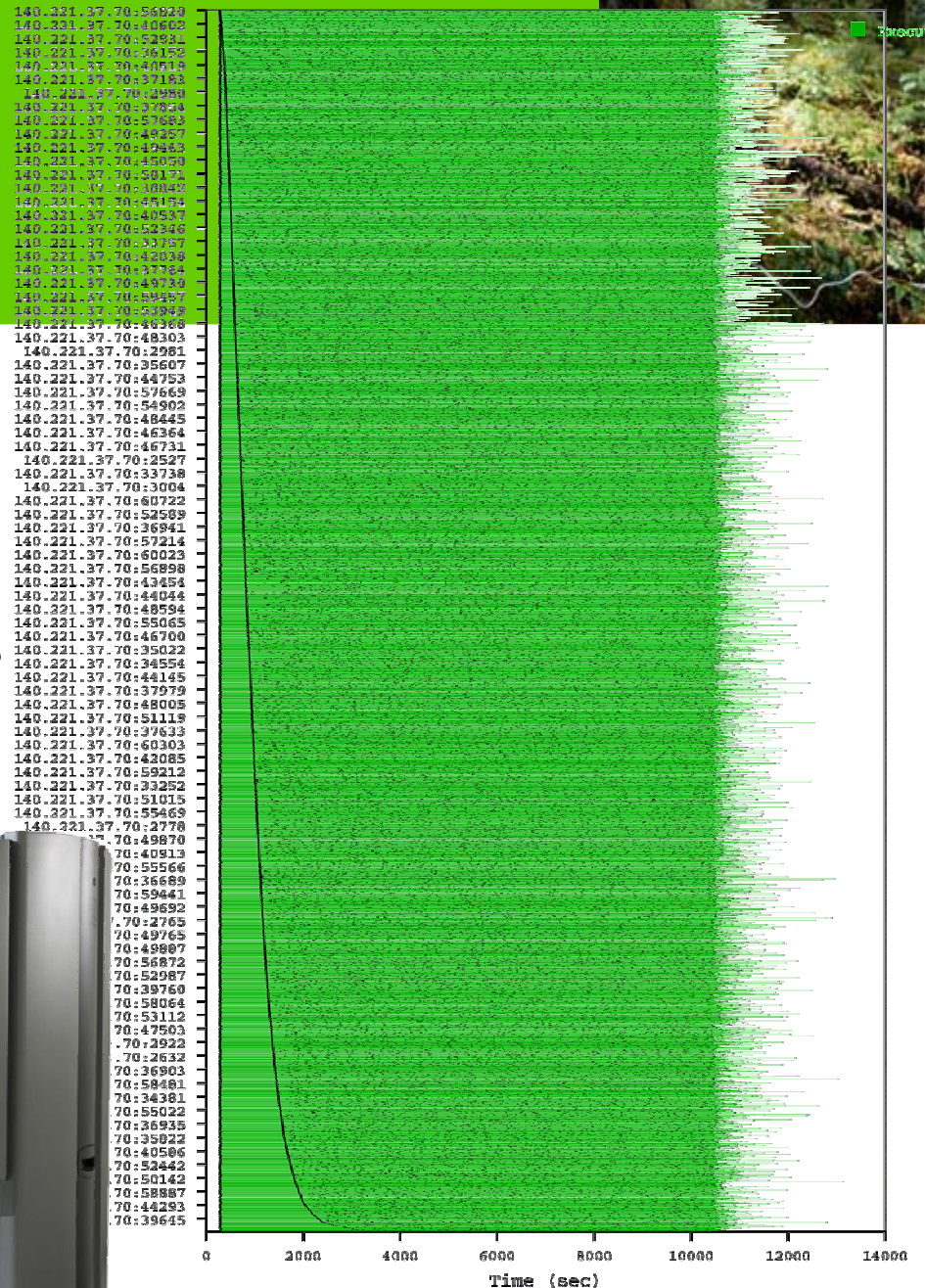
Managing and Executing Loosely-Coupled Large-Scale Applications on  
Clusters, Grids, and Supercomputers

16

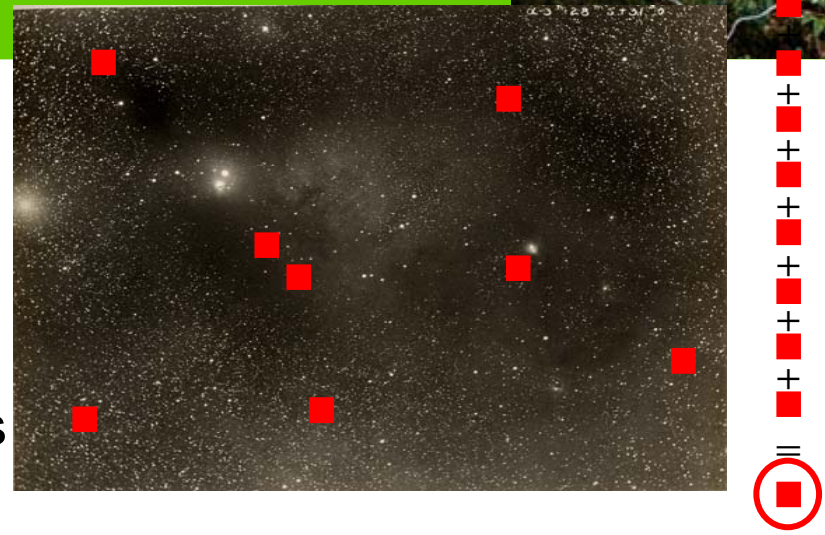


# DOCK on SiCortex

- CPU cores: 5760
- Tasks: 92160
- Elapsed time: 12821 sec
- Compute time: 1.94 CPU years
- Average task time: 660.3 sec
- Speedup: 5650X (ideal 5760)
- Efficiency: 98.2%

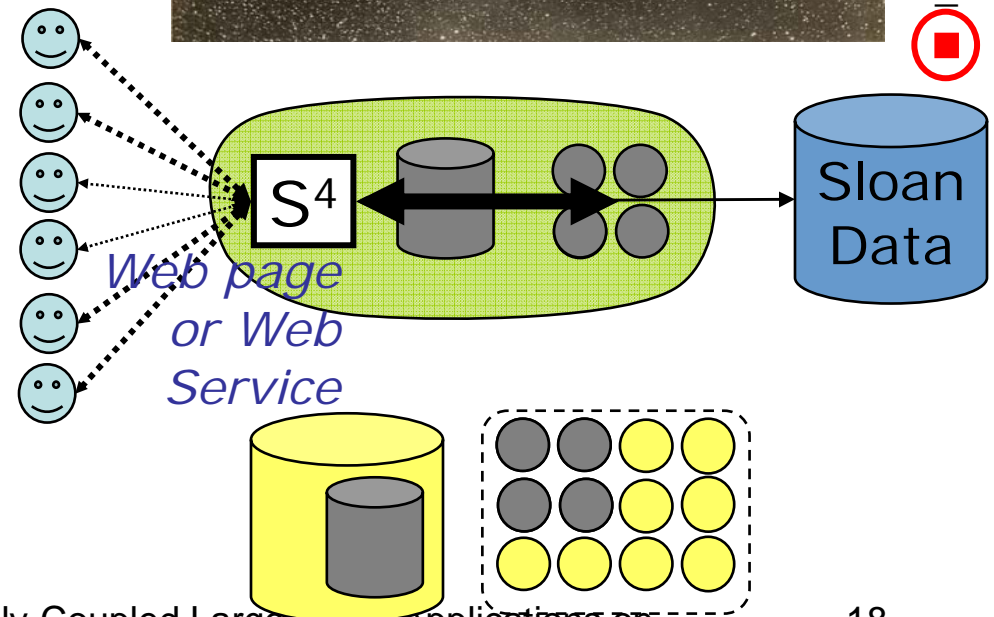


# AstroPortal Stacking Service



- Purpose
  - On-demand “stacks” of random locations within ~10TB dataset
- Challenge
  - Rapid access to 10-10K “random” files
  - Time-varying load
- Sample Workloads

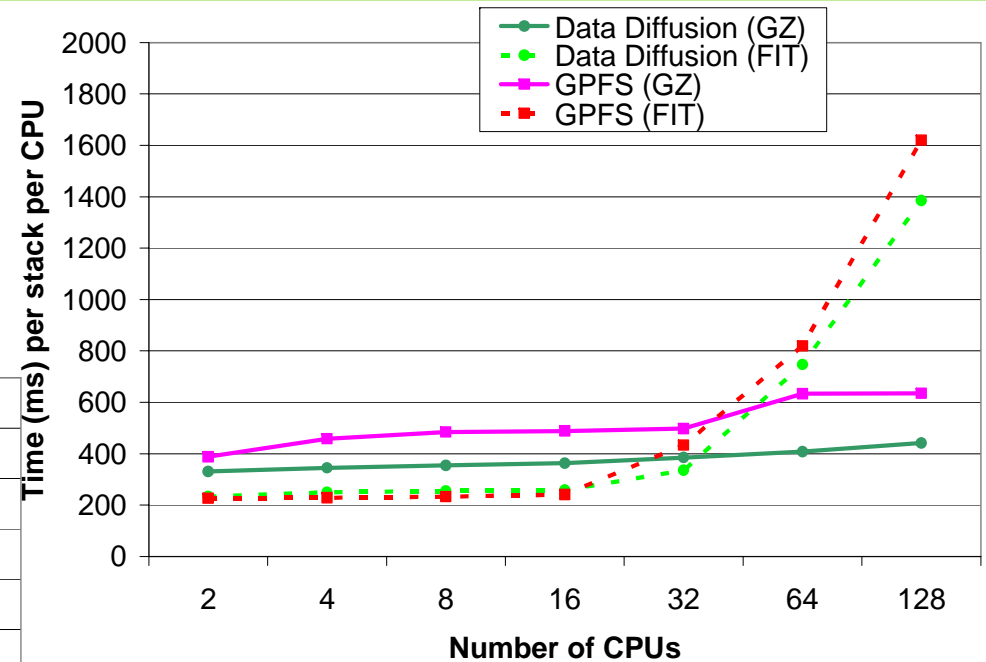
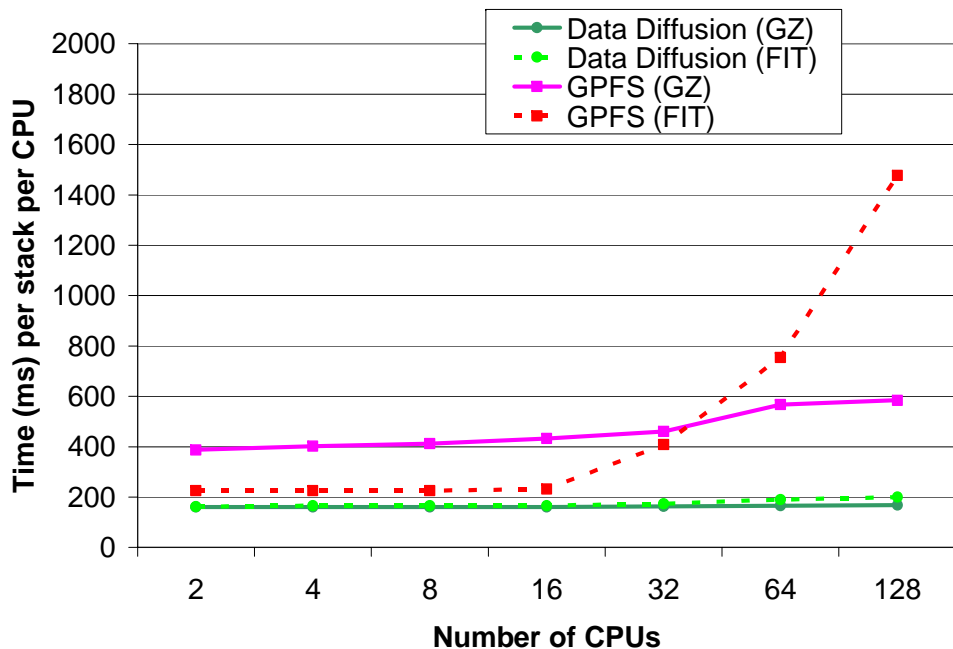
Locality	Number of Objects	Number of Files
1	111700	111700
1.38	154345	111699
2	97999	49000
3	88857	29620
4	76575	19145
5	60590	12120
10	46480	4650
20	40460	2025
30	23695	790



# AstroPortal Stacking Service with Data Diffusion



Low data locality →  
– Similar (but better)  
performance to GPFS

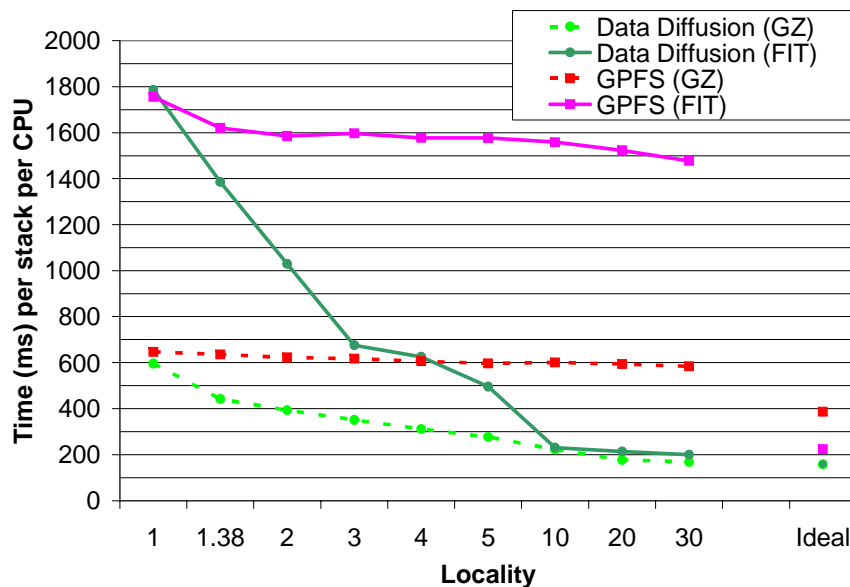
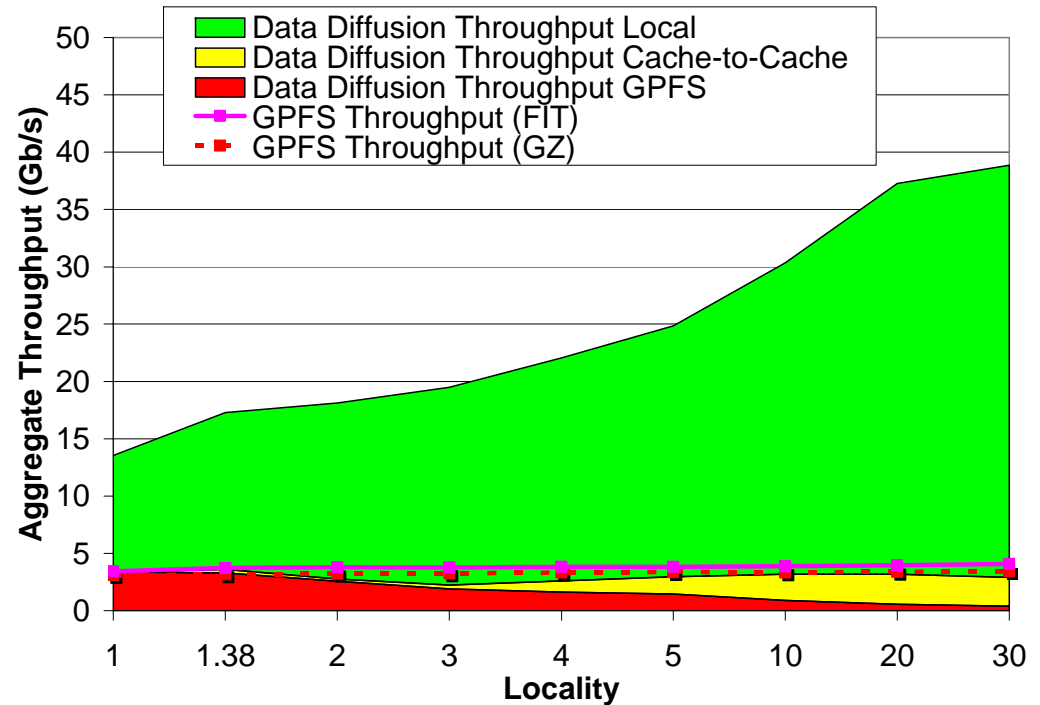


← High data locality  
– Near perfect scalability

# AstroPortal Stacking Service with Data Diffusion



- Aggregate throughput:
  - 39Gb/s
  - 10X higher than GPFS
- Reduced load on GPFS
  - 0.49Gb/s
  - 1/10 of the original load

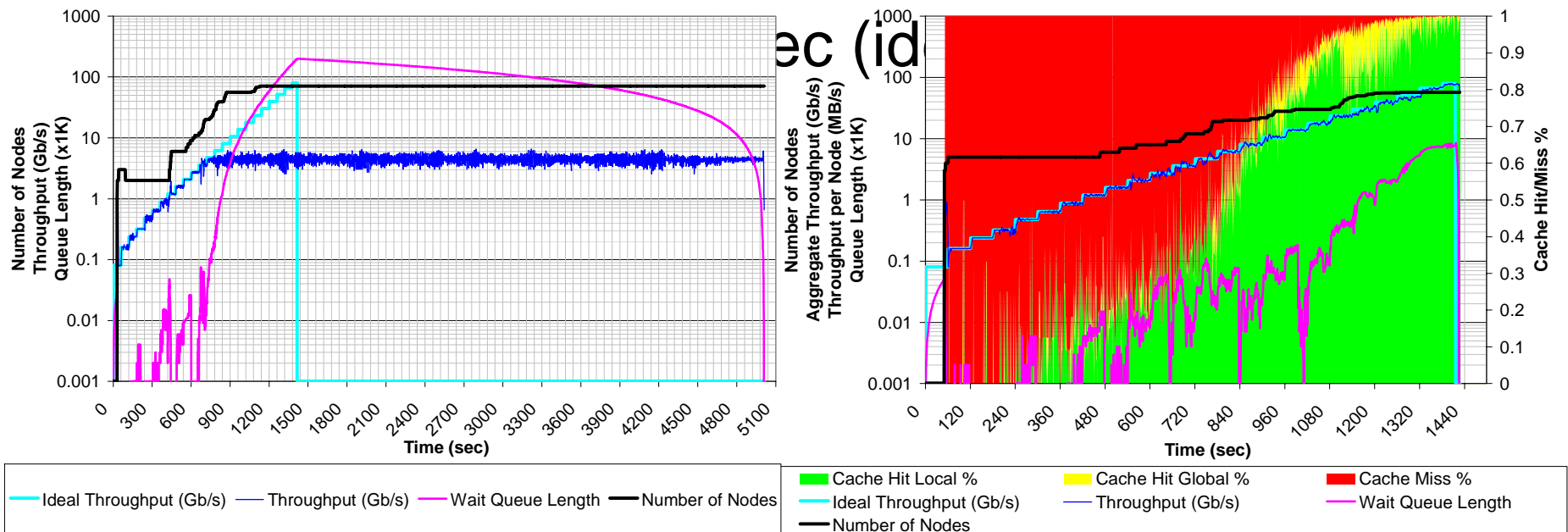


- Big performance gains as locality increases

# Data Diffusion: Data-Intensive Workload



- 250K tasks on 128 processors
  - 10MB read, 10ms compute
- Comparing GPFS with data diffusion



5/14/2008

# Data Diffusion: Data-Intensive Workload

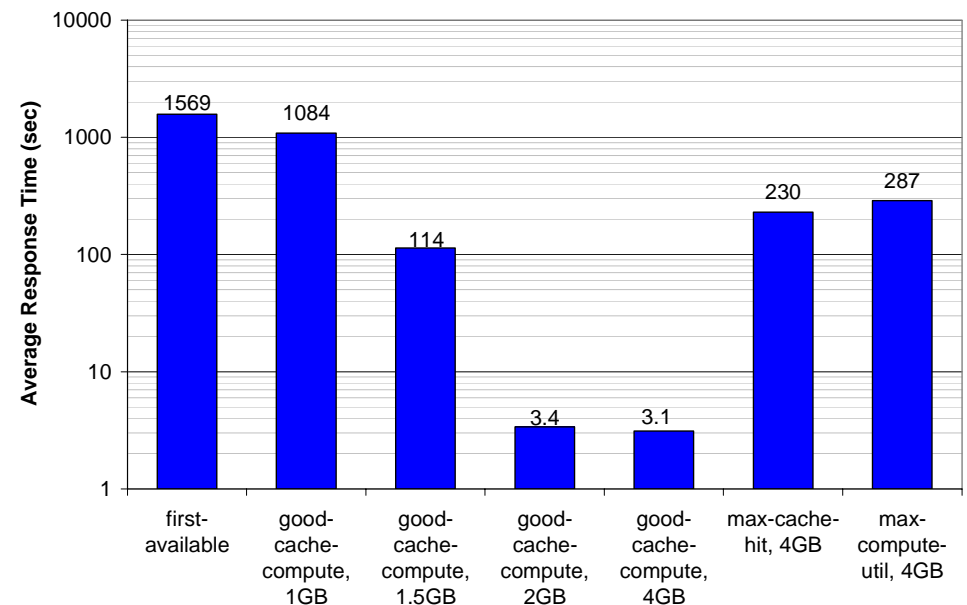


## ← Throughput:

- Average: 14Gb/s vs 4Gb/s
- Peak: 100Gb/s vs. 6Gb/s

## Response Time →

– 3 sec vs 1569 sec → 506X



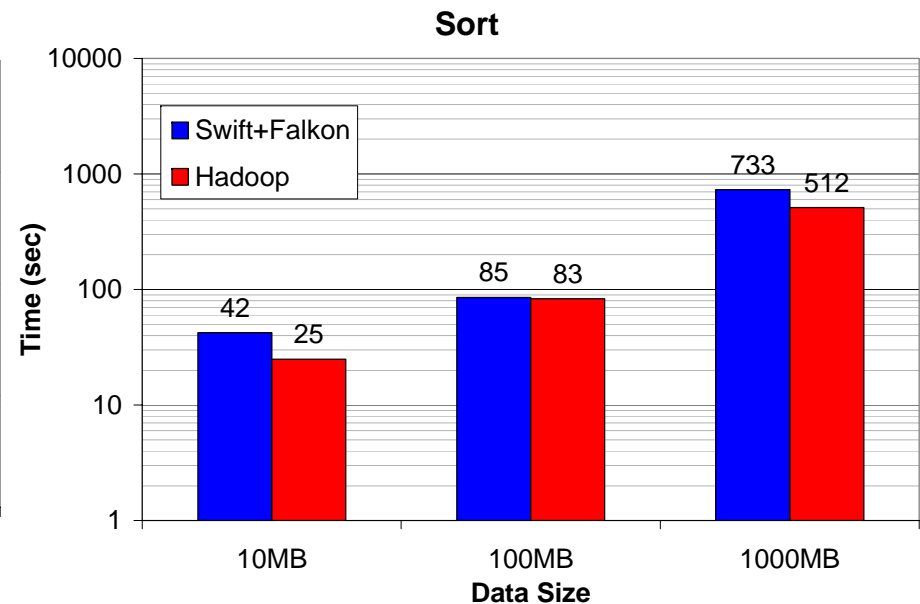
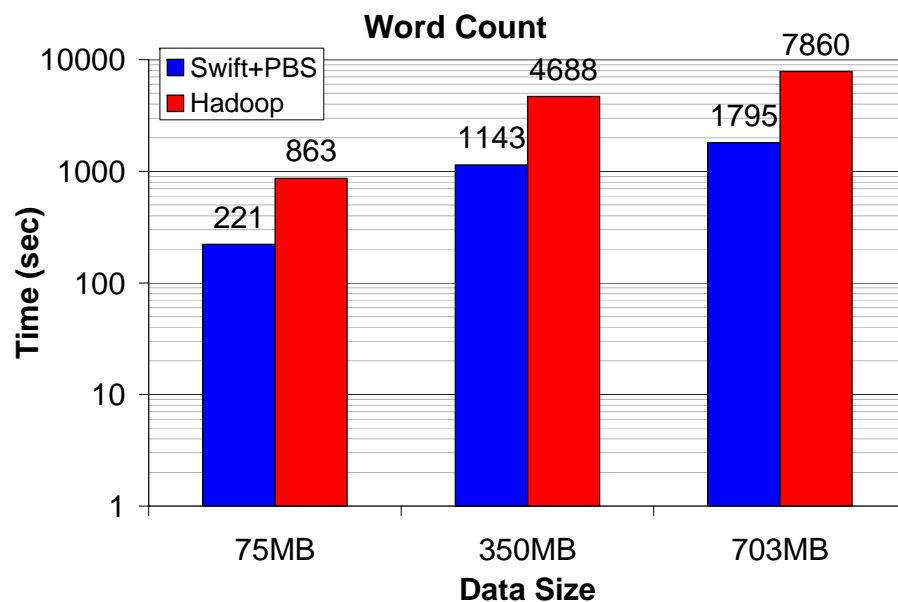
5/14/2008

Managing and Executing Loosely  
Clusters, Grids,

# Hadoop vs. Swift



- Classic benchmarks for MapReduce
  - Word Count
  - Sort
- Swift performs similar or better than Hadoop (on 32 processors)



# Mythbusting



- ~~Embarrassingly~~ Happily parallel apps are trivial to run
  - Logistical problems can be tremendous
- Loosely coupled apps do not require “supercomputers”
  - Total computational requirements can be enormous
  - Individual tasks may be tightly coupled
  - Workloads frequently involve large amounts of I/O
- Loosely coupled apps do not require specialized system software
- Shared file systems are good all around solutions
  - They don’t scale proportionally with the compute resources



# Solutions



- **Falkon**
  - A Fast and Light-weight task executiON framework
  - Globus Incubator Project
  - <http://dev.globus.org/wiki/Incubator/Falkon>
- **Swift**
  - Parallel programming tool for rapid and reliable specification, execution, and management of large-scale science workflows
  - <http://www.ci.uchicago.edu/swift/index.php>
- **Environments:**
  - *Clusters*: TeraPort (TP)
  - *Grids*: Open Science Grid (OSG), TeraGrid (TG)
  - *Specialized large machines*: SiCortex 5732
  - *Supercomputers*: IBM BlueGene/P (BG/P)

# More Information



- More information:
  - Personal research page: <http://people.cs.uchicago.edu/~iraicu/>
  - Falkon: <http://dev.globus.org/wiki/Incubator/Falkon>
  - Swift: <http://www.ci.uchicago.edu/swift/index.php>
- Collaborators (relevant to this proposal):
  - Ian Foster, The University of Chicago & Argonne National Laboratory
  - Alex Szalay, The Johns Hopkins University
  - Yong Zhao, Microsoft
  - Mike Wilde, Computation Institute, University of Chicago & Argonne National Laboratory
  - Catalin Dumitrescu, Fermi National Laboratory
  - Zhao Zhang, The University of Chicago
  - Jerry C. Yan, NASA, Ames Research Center
- Funding:
  - NASA: Ames Research Center, Graduate Student Research Program (GSRP)
  - DOE: Mathematical, Information, and Computational Sciences Division subprogram of the Office of Advanced Scientific Computing Research, Office of Science, U.S. Dept. of Energy
  - NSF: TeraGrid