



THE UNIVERSITY OF
CHICAGO



Scientific Workflow Systems for 21st Century, New Bottle or New Wine?

Ioan Raicu

Distributed Systems Laboratory
Computer Science Department
University of Chicago

Collaborators:

Yong Zhao, Microsoft Corporation

Ian Foster, University of Chicago and Argonne National Laboratory



Scientific Workflows (SWF 2008)
July 8th, 2008

Workflow Systems

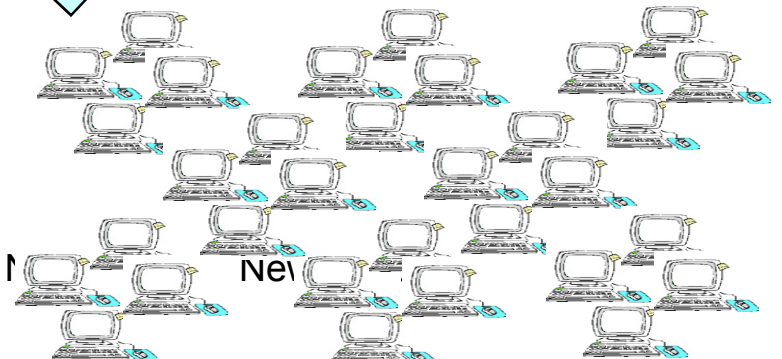
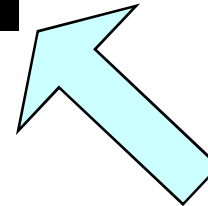
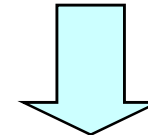
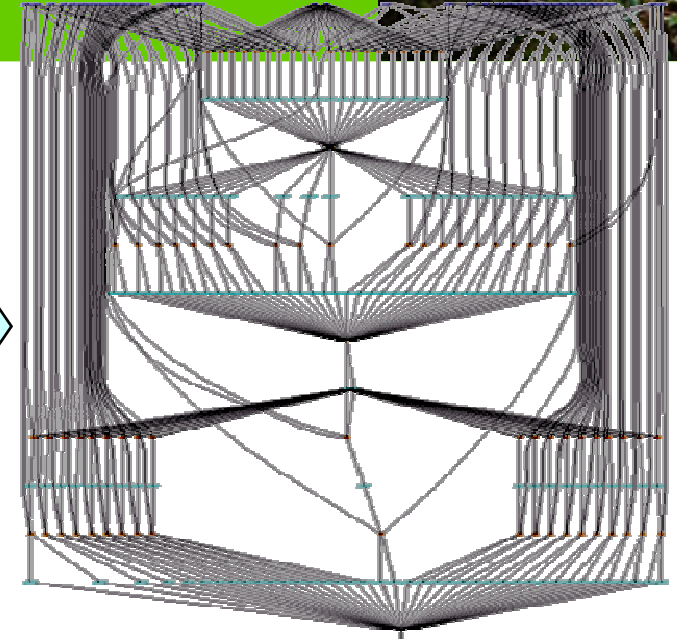
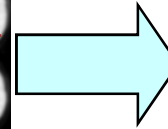
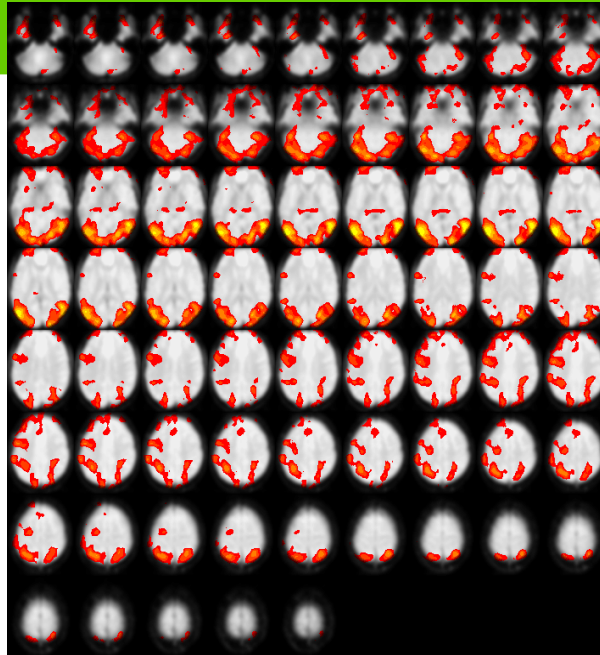
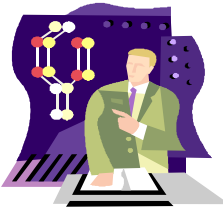


Describes computation components, ports and channels

Describes data and event flow

Coordinate the execution of the components

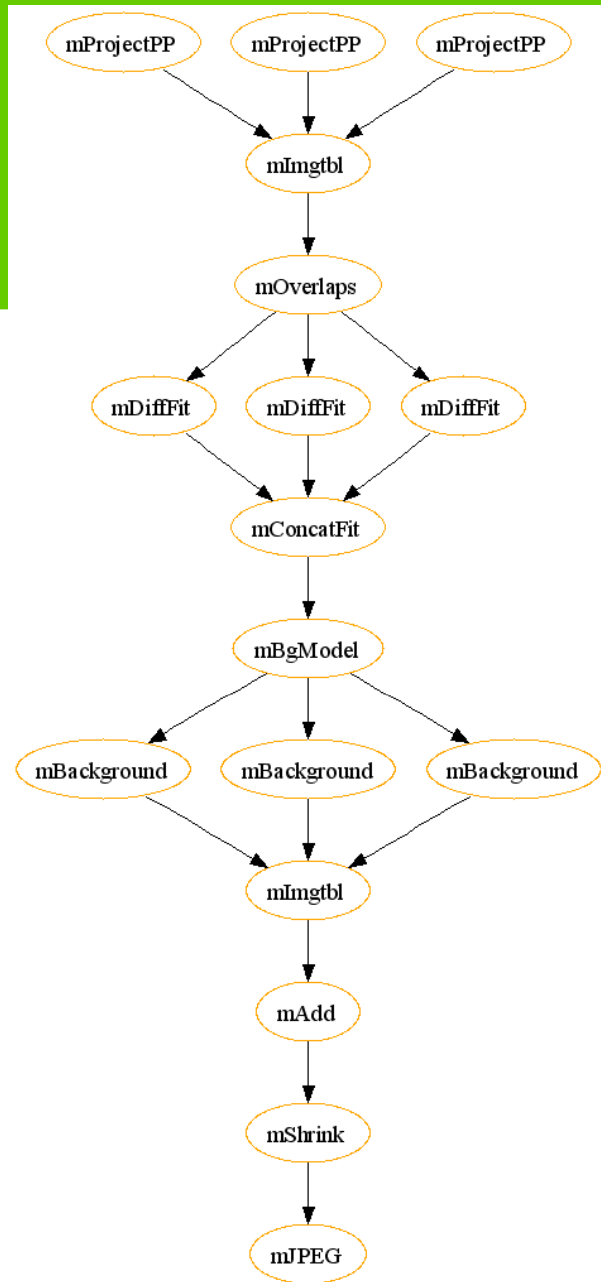
Functional MRI (fMRI)



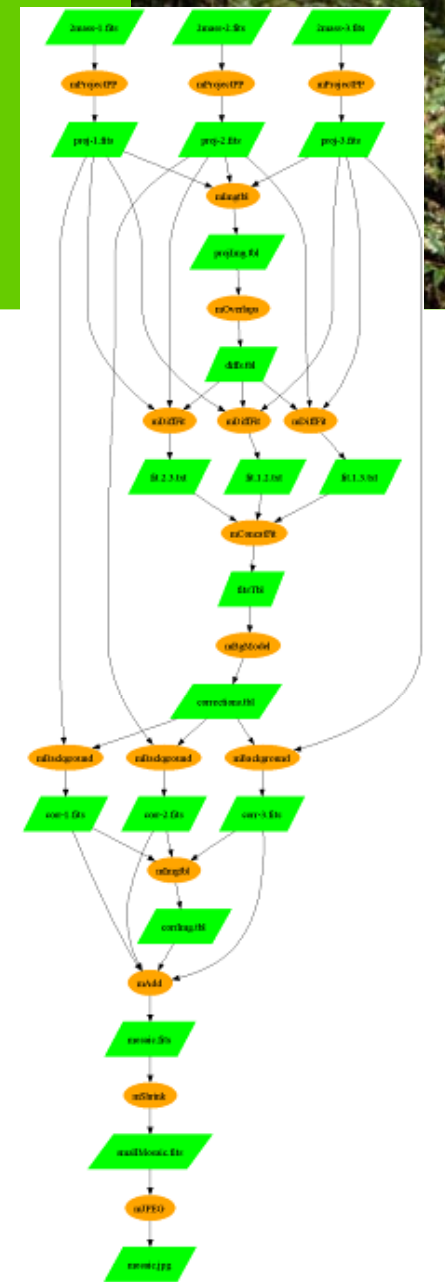
Century, 1

Net

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



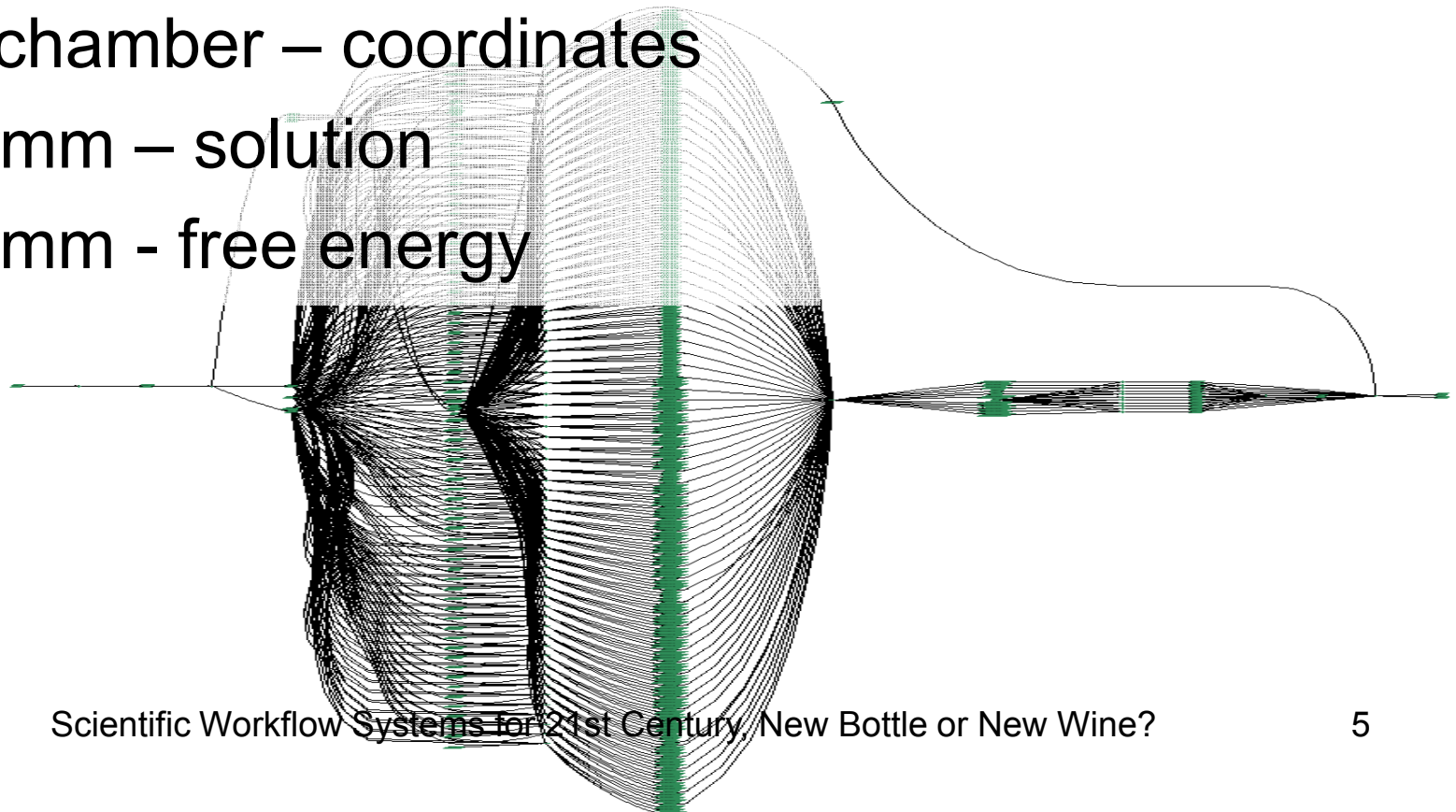
B. Berriman, J. Good (Caltech)
 J. Jacob, D. Katz (JPL)



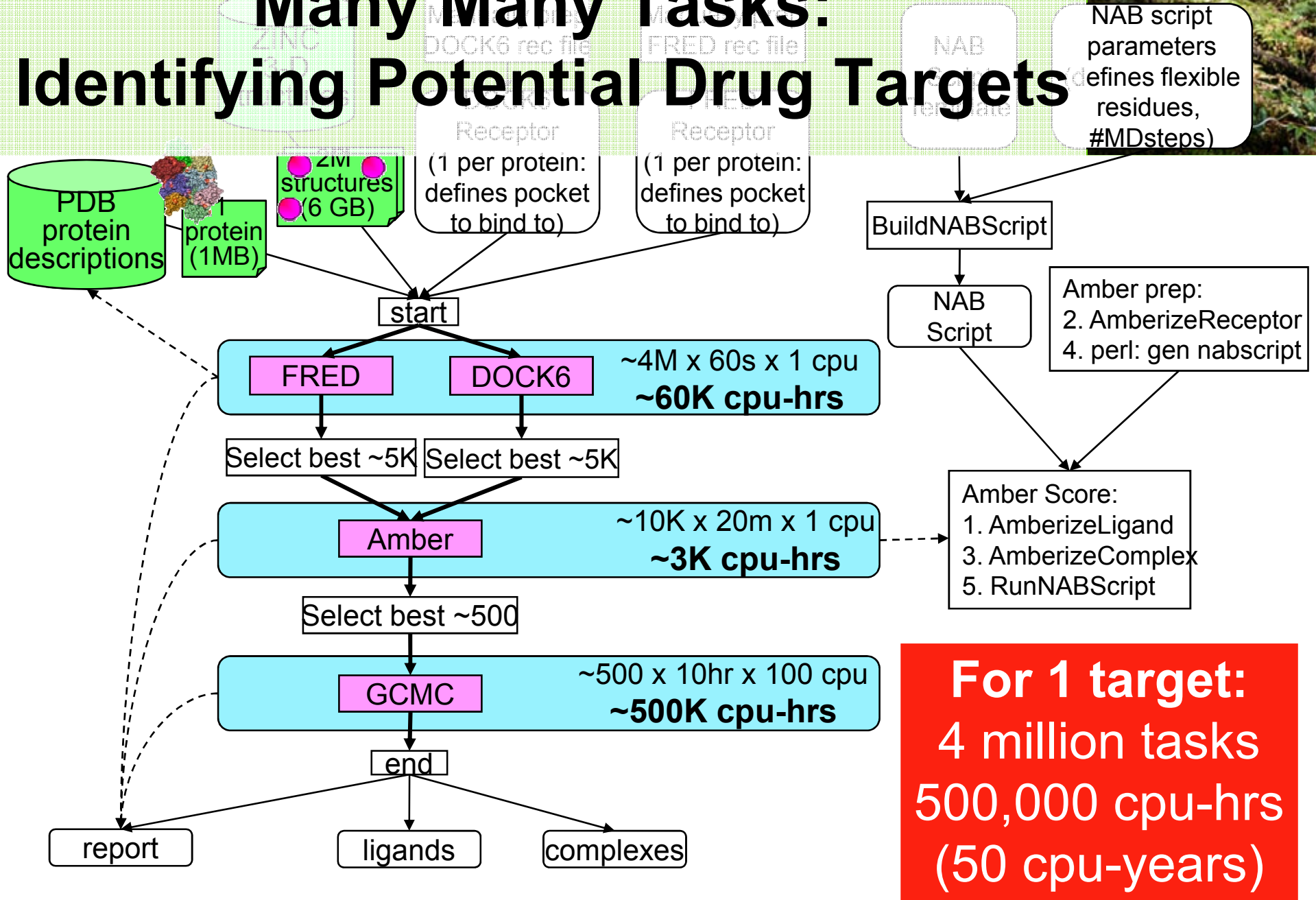
Molecular Dynamics



- Determination of free energies in aqueous solution
 - Antechamber – coordinates
 - Charmm – solution
 - Charmm - free energy



Many Many Tasks: Identifying Potential Drug Targets



Characterizing Scientific Workflows



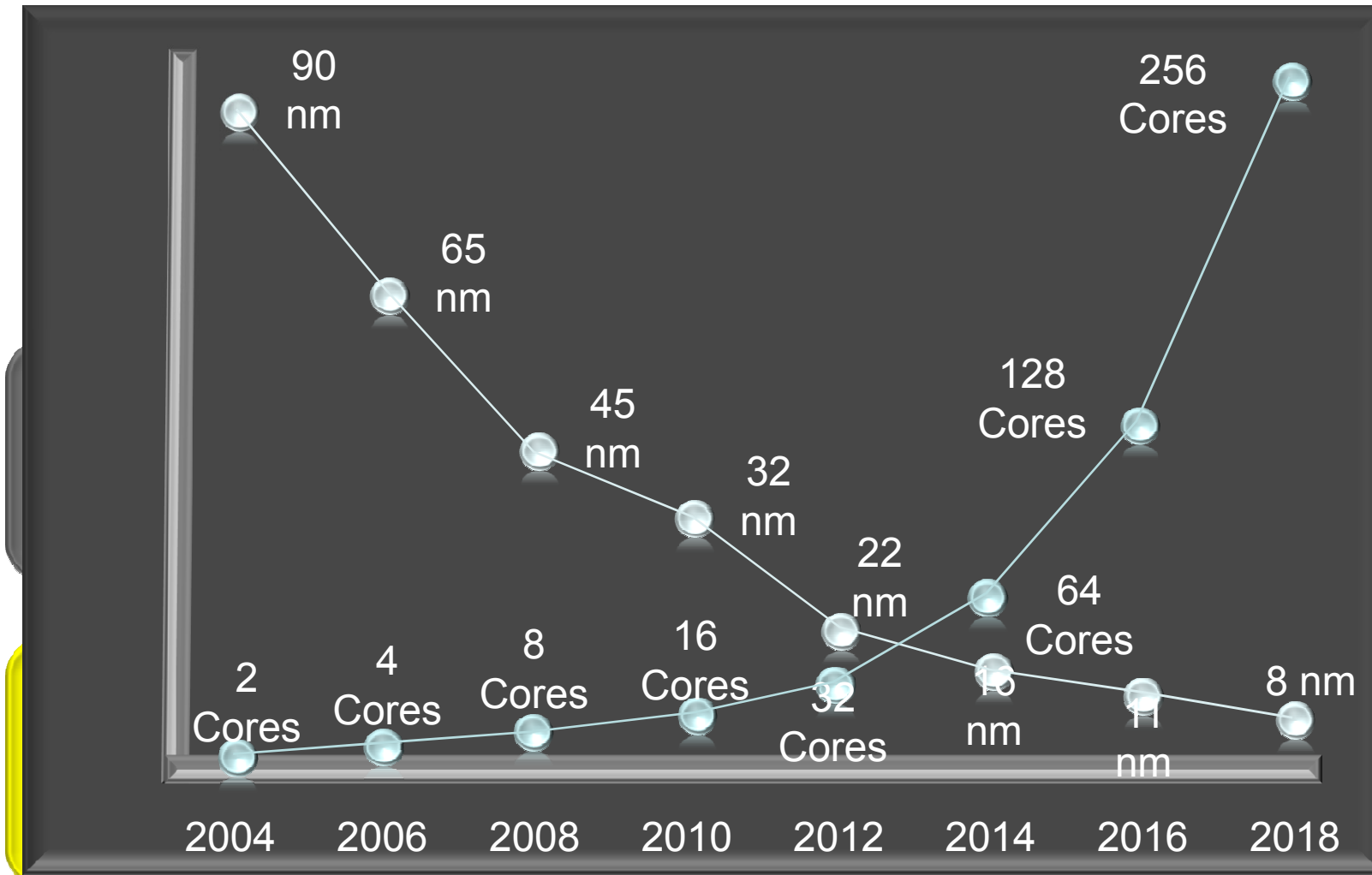
- Inherit *workflow system* definition +...
- Describe complex scientific procedures
- Automate data derivation processes
- High performance computing (HPC) to improve throughput and performance
- Provenance management and query

Characterizing Scientific Applications

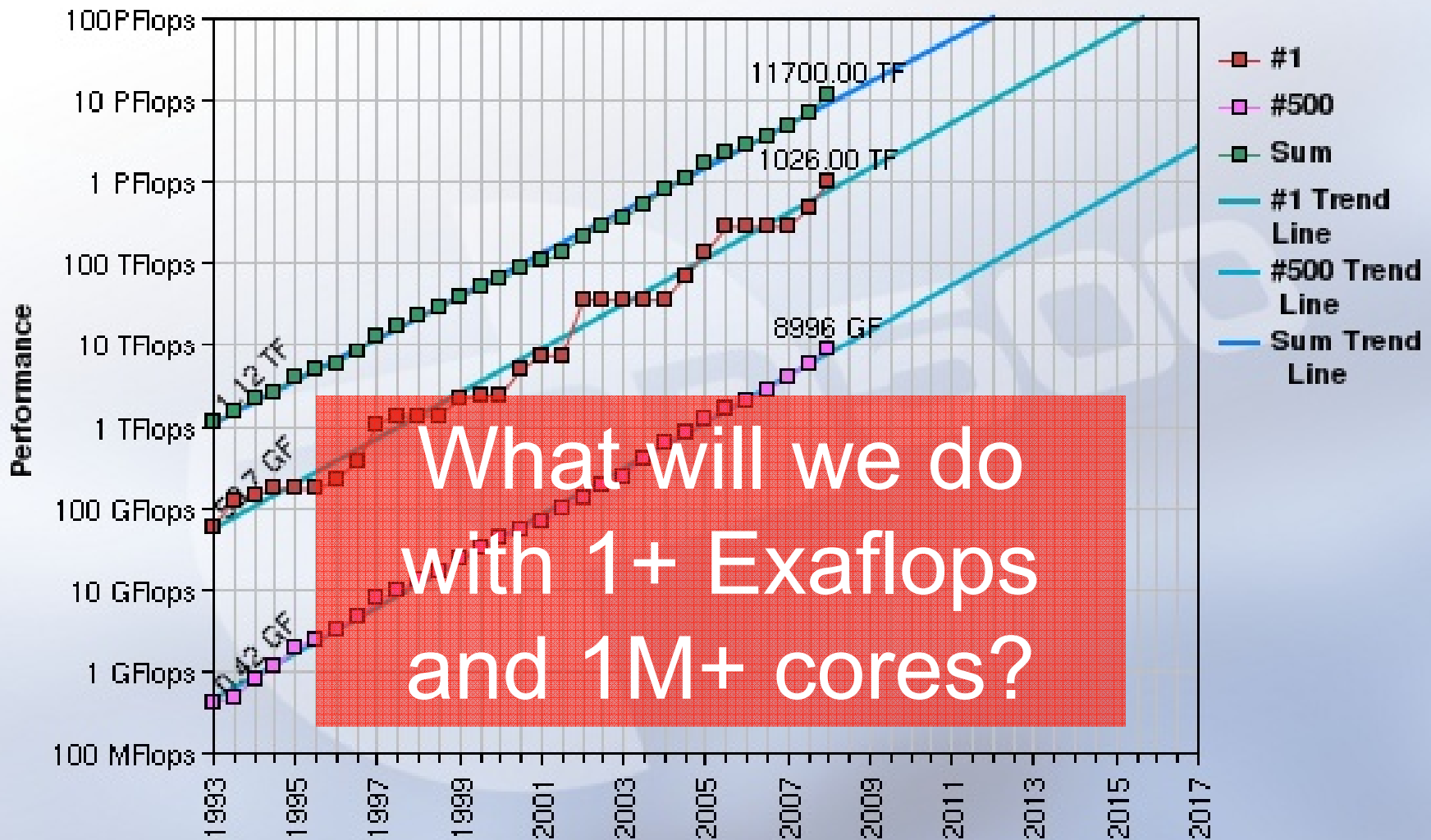


- Increasing in scale and complexity
 - Wide inherent parallelism
 - Multiple successive stages
 - Wide range of number of tasks
 - thousands to billions
- Potentially compute intensive
 - Widely varying task execution times
 - 100s of ms to 10s of hours
- Potentially data intensive
 - I/O operation rates and aggregate I/O rates
 - Meta-data creation and modification
 - Significant data re-use
 - Produced data is consumed by later stages

Many-Core Growth Rates

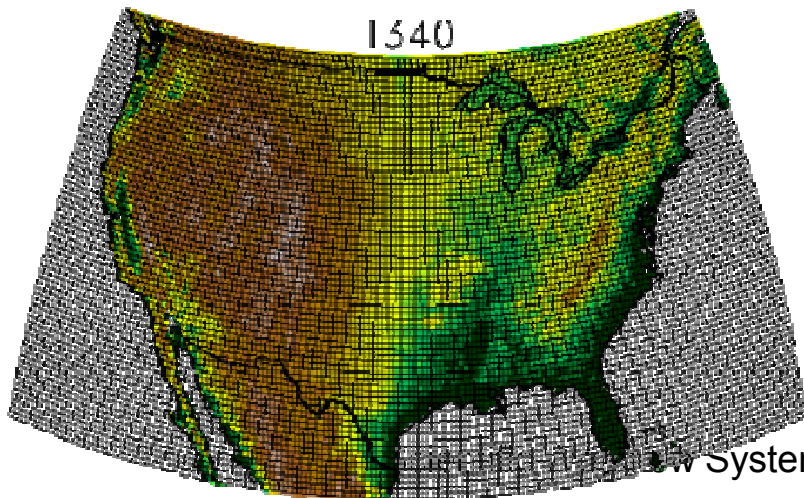
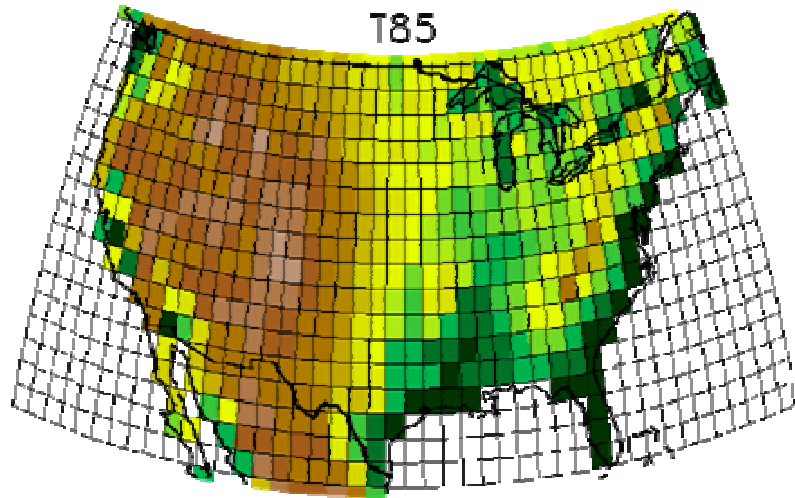


Projected Performance Development



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

1) Tackle **Bigger and Bigger** Problems

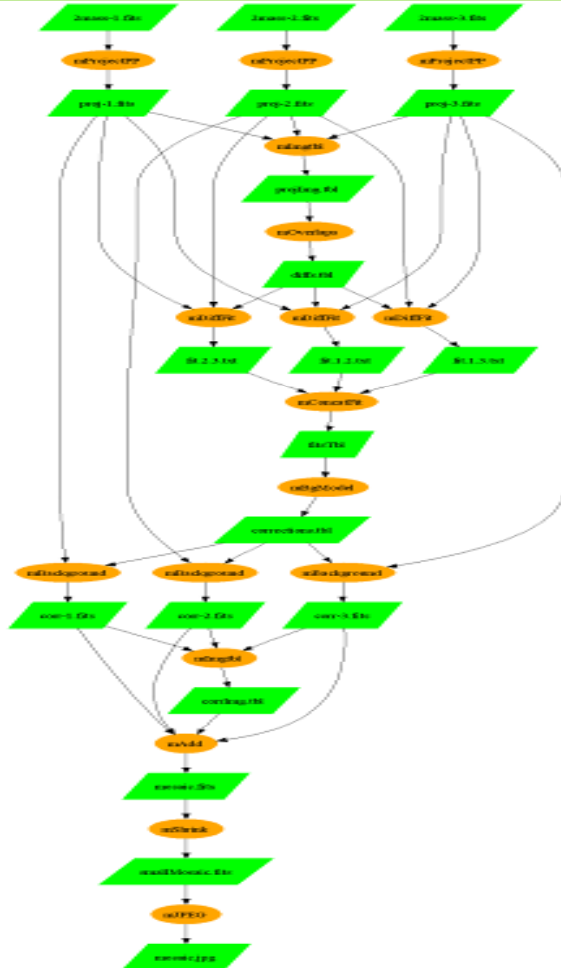


Computational
Scientist
as
Hero

2) Tackle **Increasingly Complex** Problems



Computational
Scientist
as
**Logistics
Officer**



“More Complex Problems”



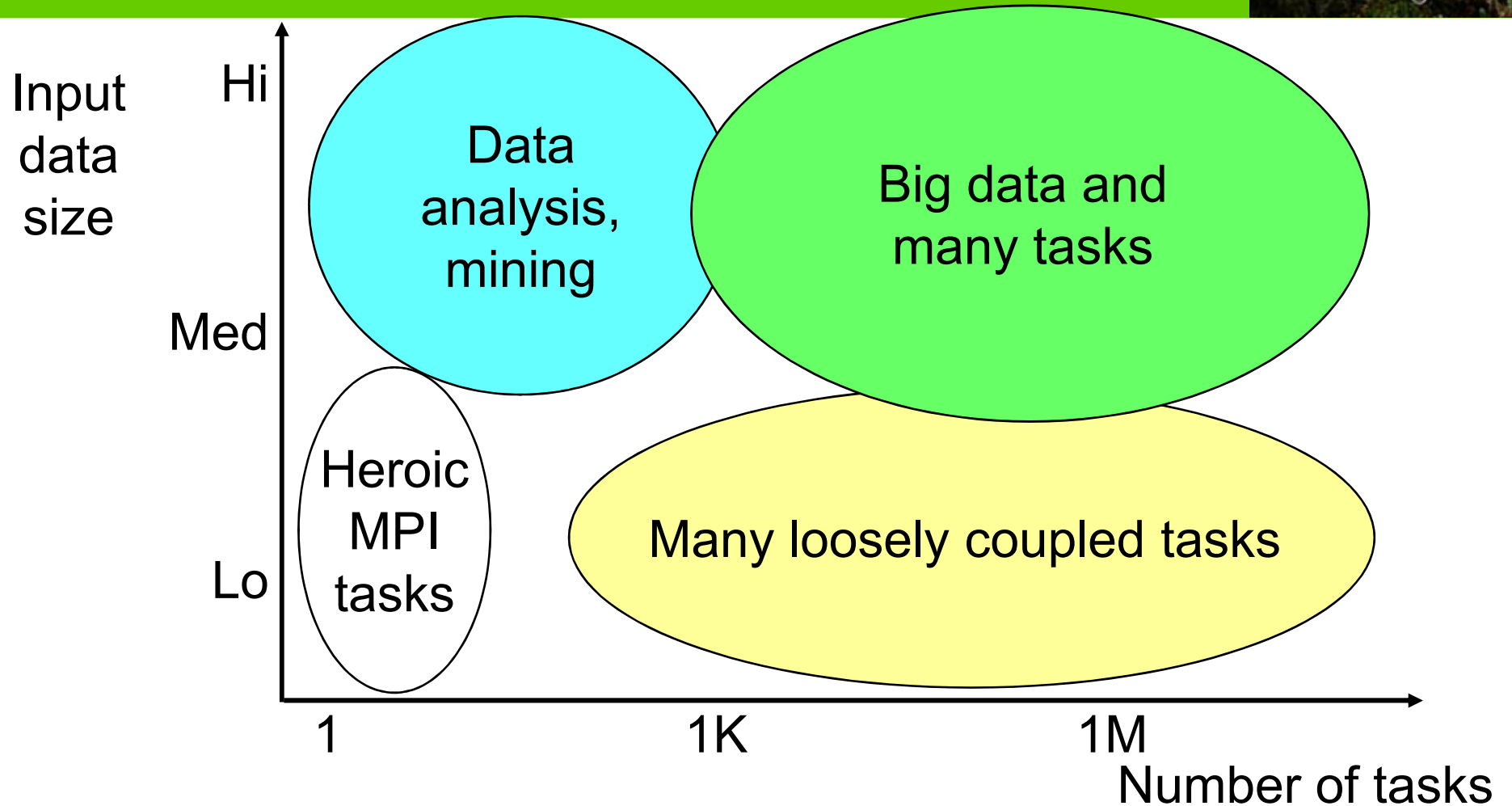
- Ensemble runs to quantify **climate model uncertainty**
- Identify **potential drug targets** by screening a database of ligand structures against target proteins
- Study **economic model sensitivity** to parameters
- Analyze **turbulence dataset** from many perspectives
- Perform **numerical optimization** to determine optimal resource assignment in energy problems
- Mine collection of data from **advanced light sources**
- Construct databases of computed properties of **chemical compounds**
- Analyze data from the **Large Hadron Collider**
- Analyze **log data** from 100,000-node parallel computations

Programming Model Issues

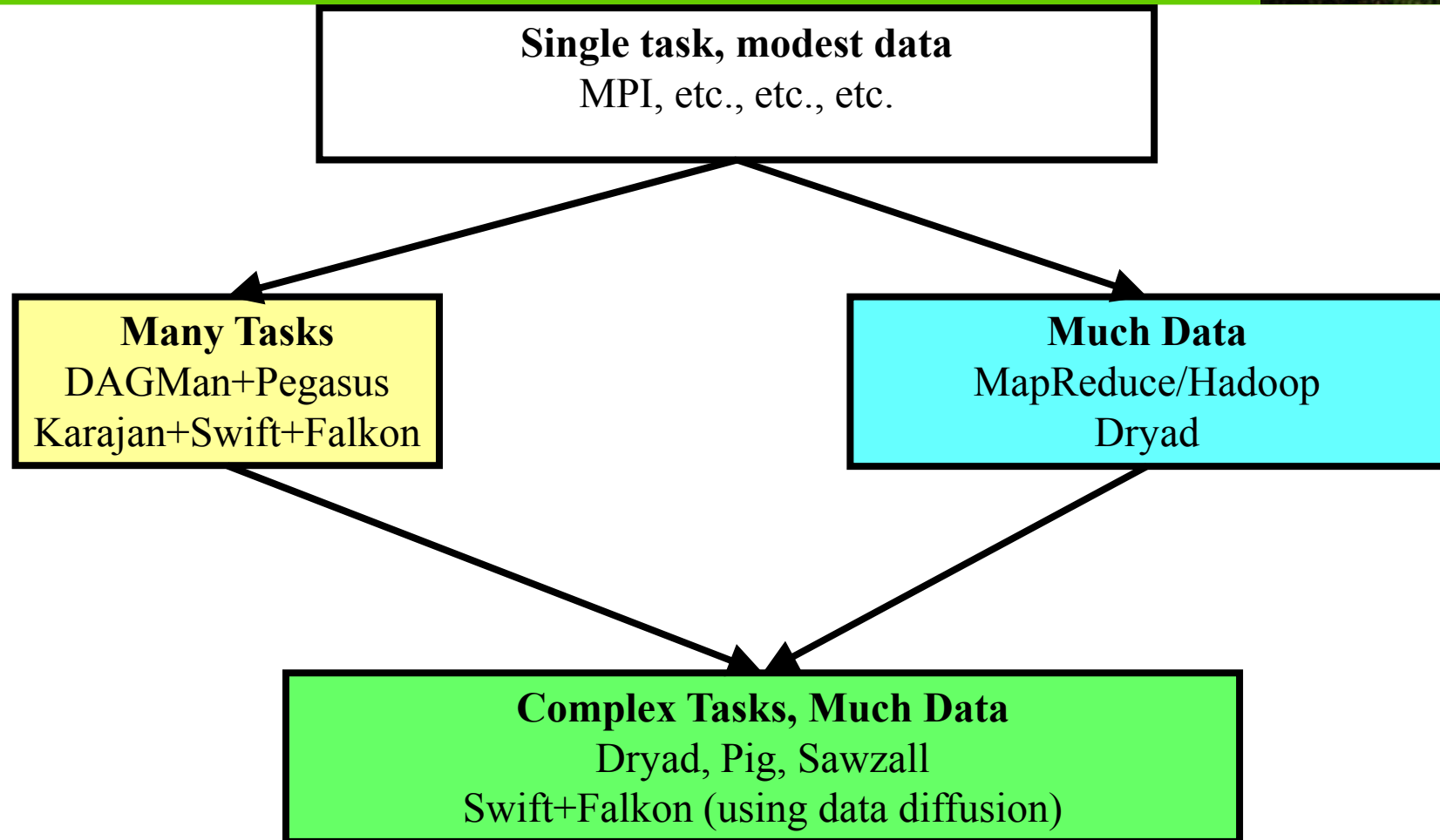


- Massive **task parallelism**
- Massive **data parallelism**
- Integrating **black box applications**
- Complex **task dependencies** (task graphs)
- **Failure**, and other execution management issues
- **Data management**: input, intermediate, output
- **Dynamic computations** (task graphs)
- **Dynamic data access** to large, diverse datasets
- **Long-running** computations
- Documenting **provenance** of data products

Problem Types



An Incomplete and Simplistic View of Programming Models and Tools



Scientific Workflow Systems for 21st Century, New Bottle or New Wine?

Major Challenges to Large Scale Scientific Computation



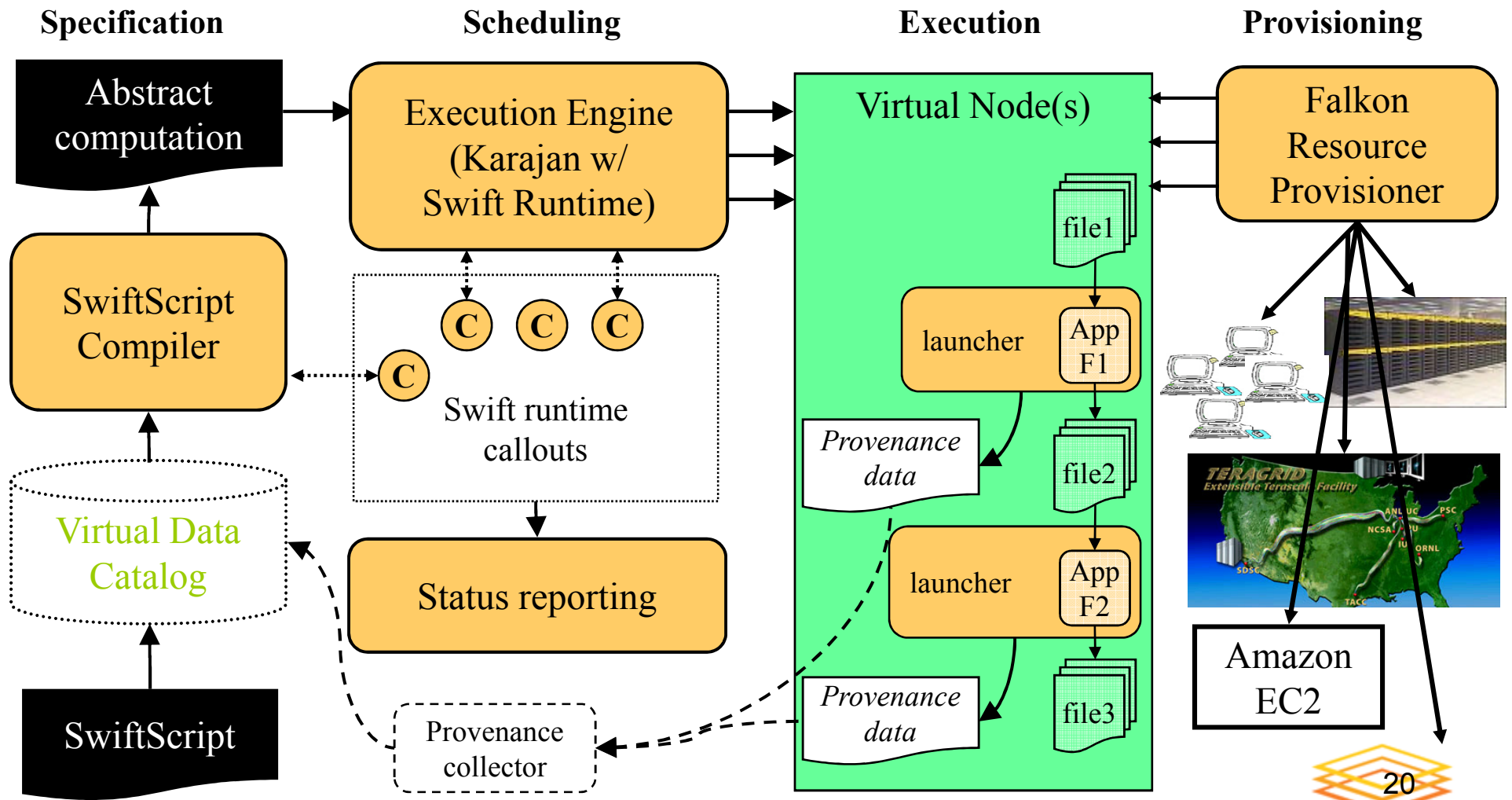
- Managing heterogeneous scientific data
 - Idiosyncratic layouts & formats (file sys, db, spreadsheet, XML, etc...)
- Describing complex science problems
- Coordinating distributed diverse computation procedures
 - Executables, scripts, Web services
- Scheduling & executing numerous tasks reliably and efficiently
 - Large quantity of data (Petabytes/year)
 - Large number of parallel/dependent tasks ($10^3 \sim 10^6$ tasks)
- Organizing, archiving and tracking
 - Datasets, procedures, workflows, provenance

Existing and emerging workflow technologies



	Star-P	Mapreduce	Kepler	Triana	Tavana	DAGMan	MS Workflow Foundation	XPDL	BPEL	Swift
Scales to Grids	+	+	-	-	-	++	-	-	-	++
Typing	+	-	+	-	-	-	++	++	++	++
Iteration	+	+	+	-	-	-	+	-	-/+	++
Scripting	++	-	-	-	+	+	+	-	-	++
Dataset Mapping	-	-	-	-	-	-	-	-	-	+
Service Interop	-	-	+	-	-	-	-	+	-	+
Subflow/comp.	+	-	+	+	-	-	+	+	-	+
Provenance	-	-	+	-	+	-	+	-	-	+
Open source	-	-	+	+	+	+	-	+	+	+

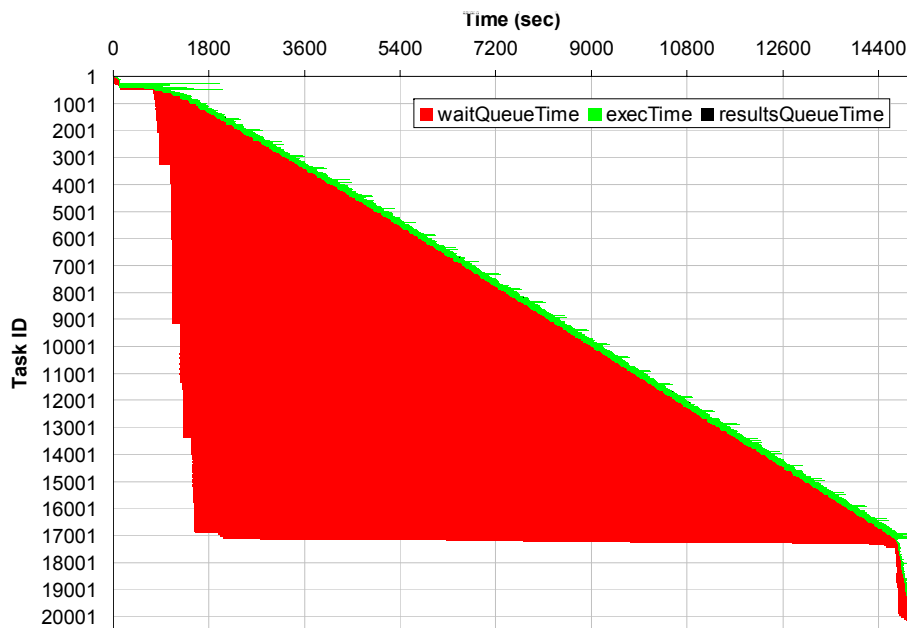
Swift Architecture



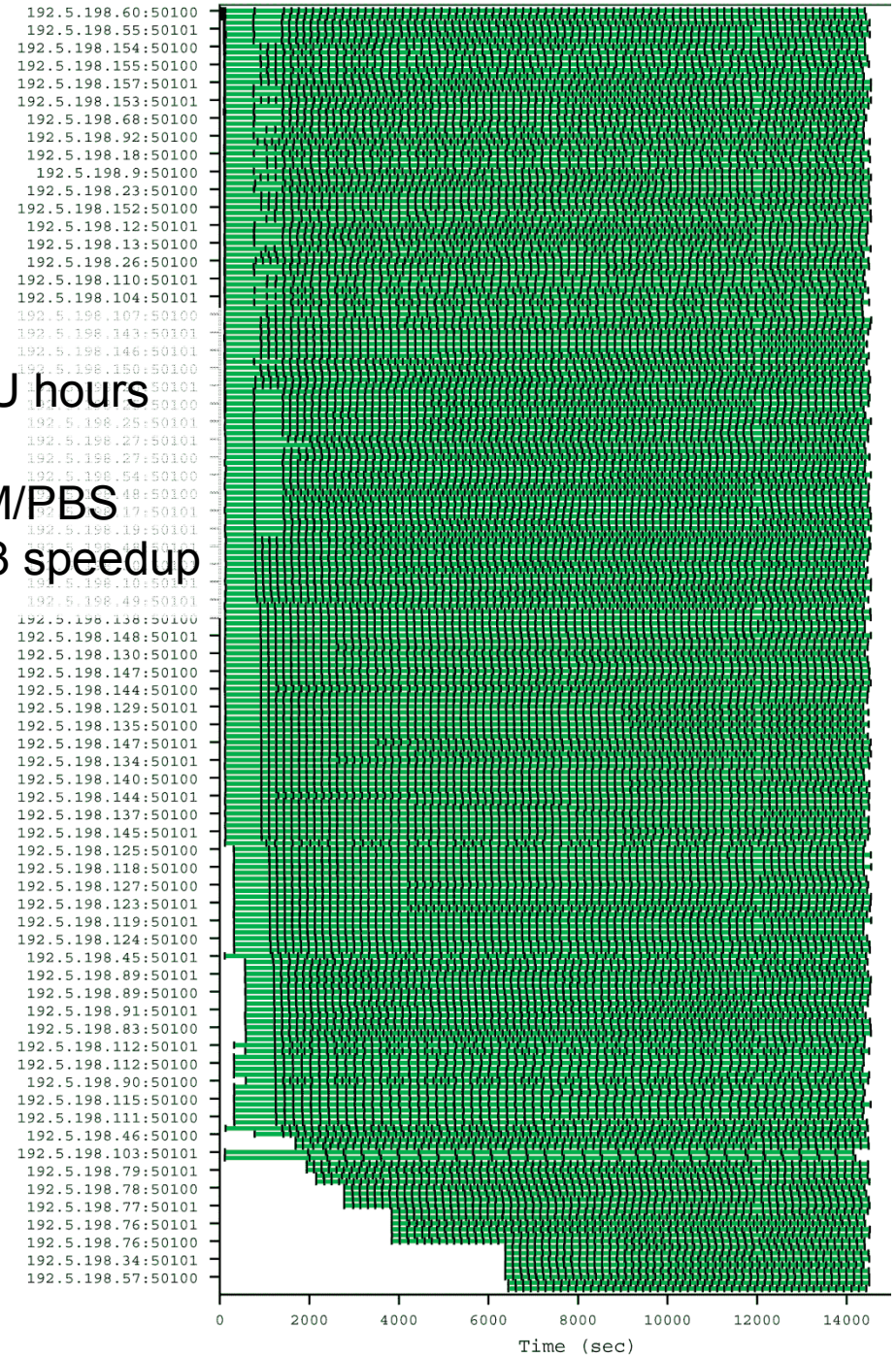
Scientific Workflow Systems for 21st Century, New Bottle or New Wine?

MolDyn Application

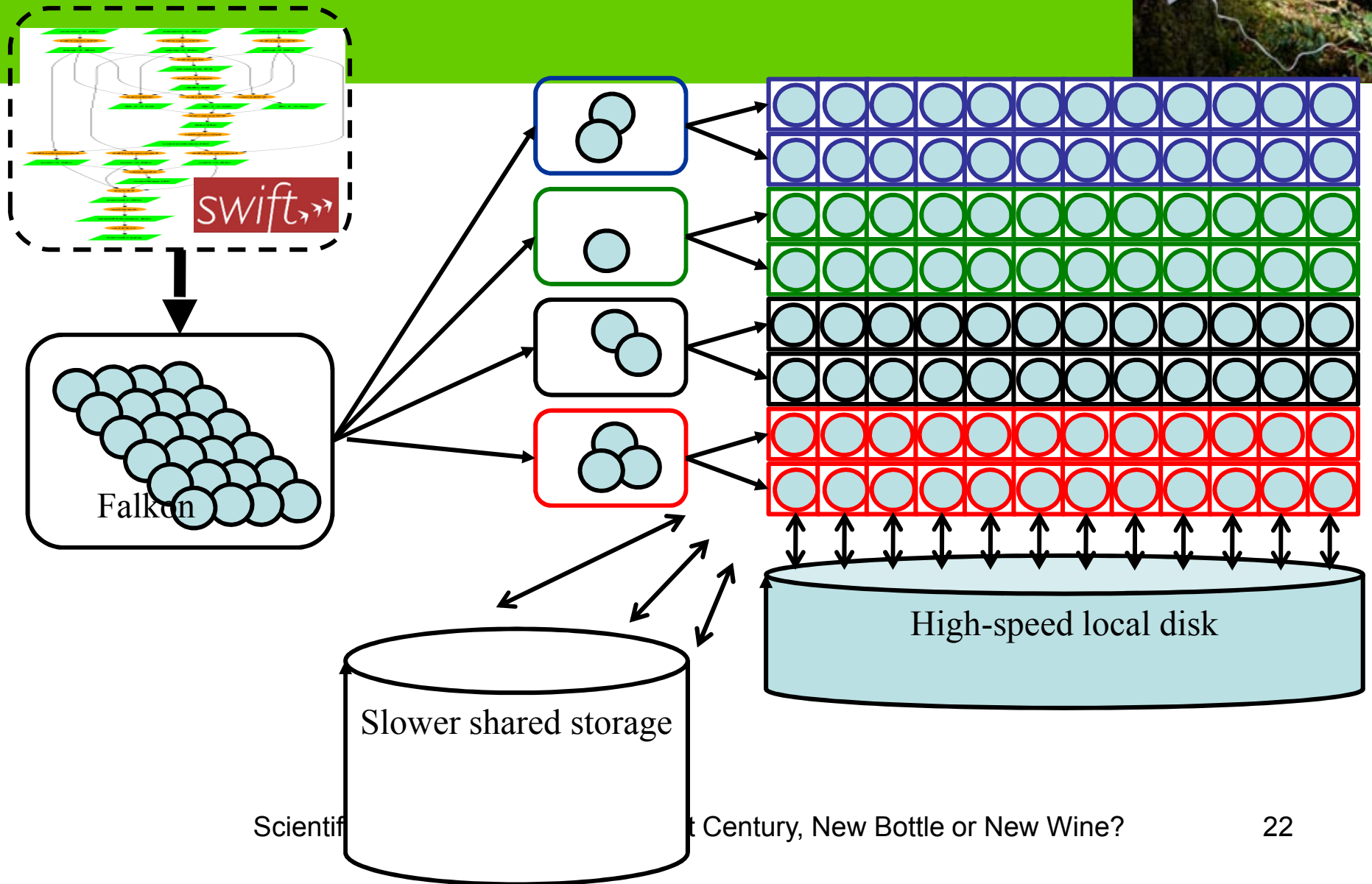
- 244 molecules → 20497 jobs
- 15091 seconds on 216 CPUs → 867.1 CPU hours
- Efficiency: **99.8%**
- Speedup: 206.9x → 8.2x faster than GRAM/PBS
- 50 molecules w/ GRAM (4201 jobs) → 25.3 speedup



Scientific Workflow Systems for 21st



Managing 120K CPUs

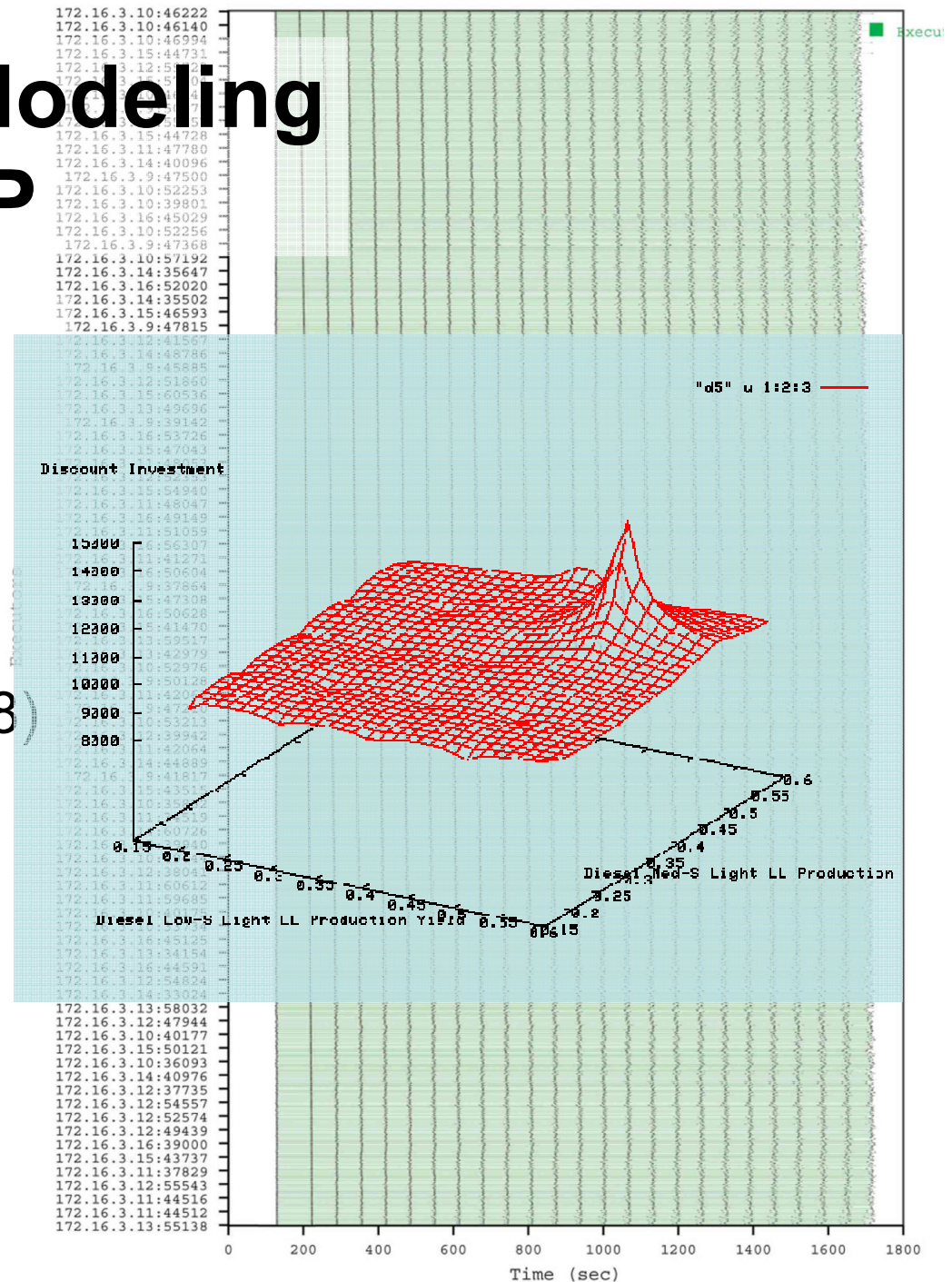


MARS Economic Modeling on IBM BG/P

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



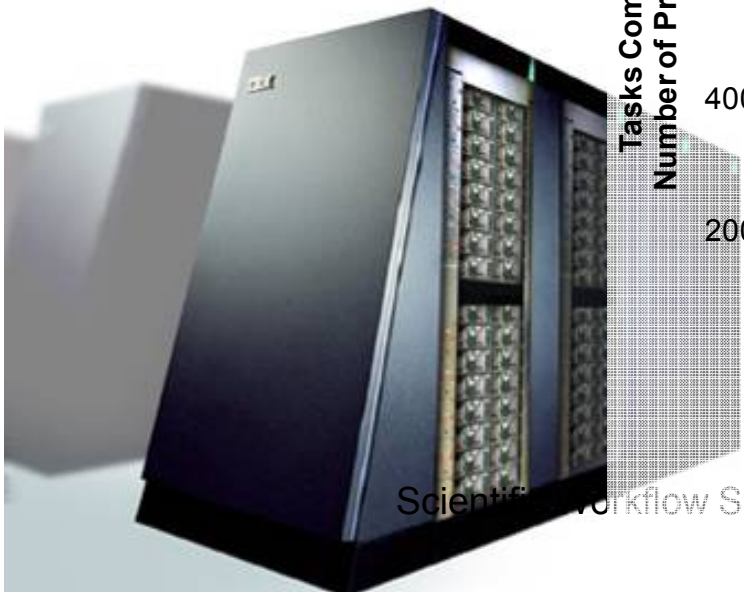
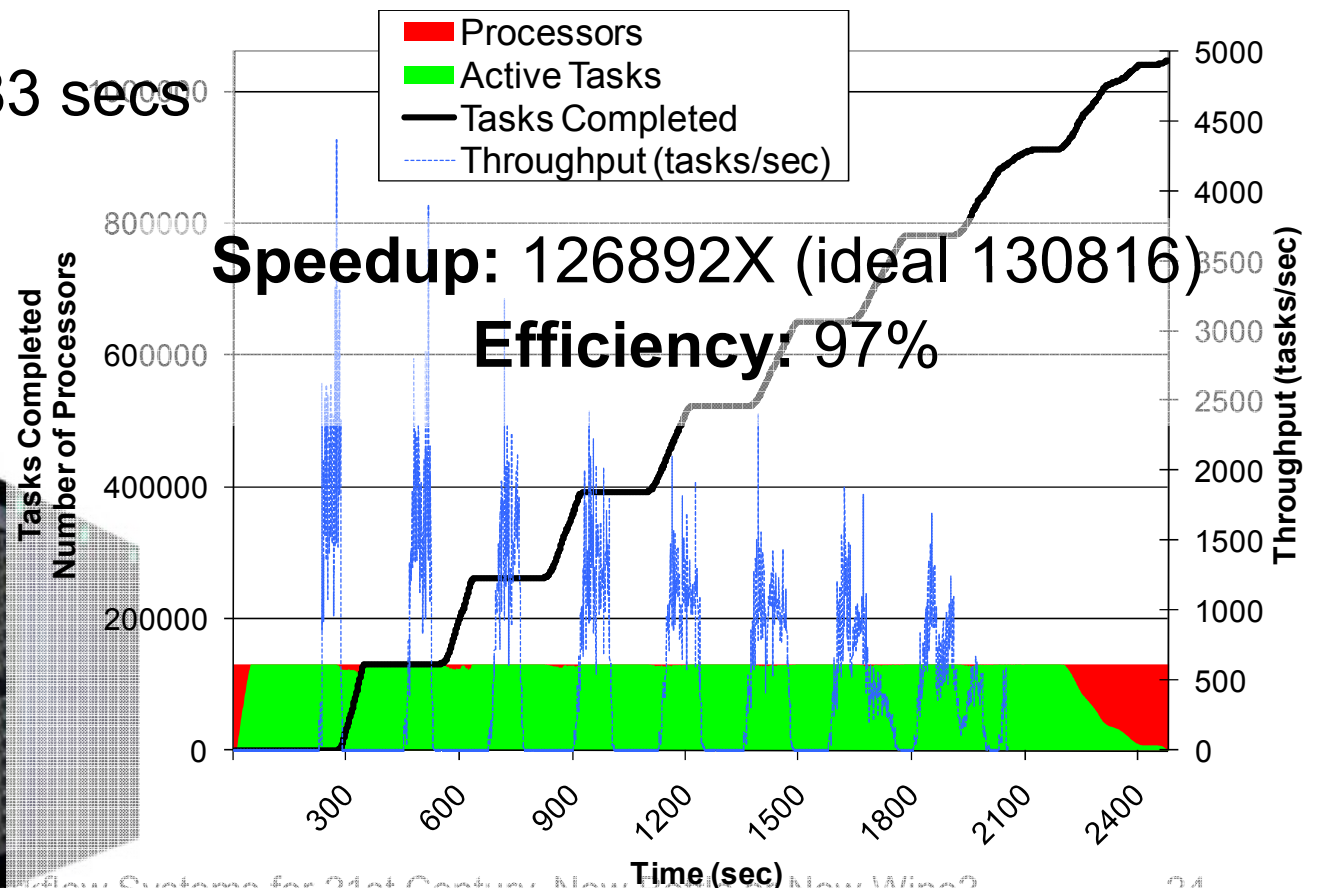
systems 1



MARS Economic Modeling on IBM BG/P (128K CPUs)



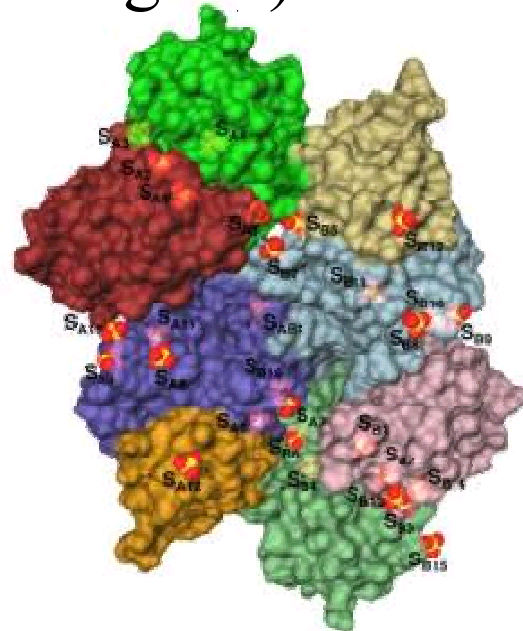
- CPU Cores: 130816
- Tasks: 1048576
- Elapsed time: 2483 secs
- CPU Years: 9.3



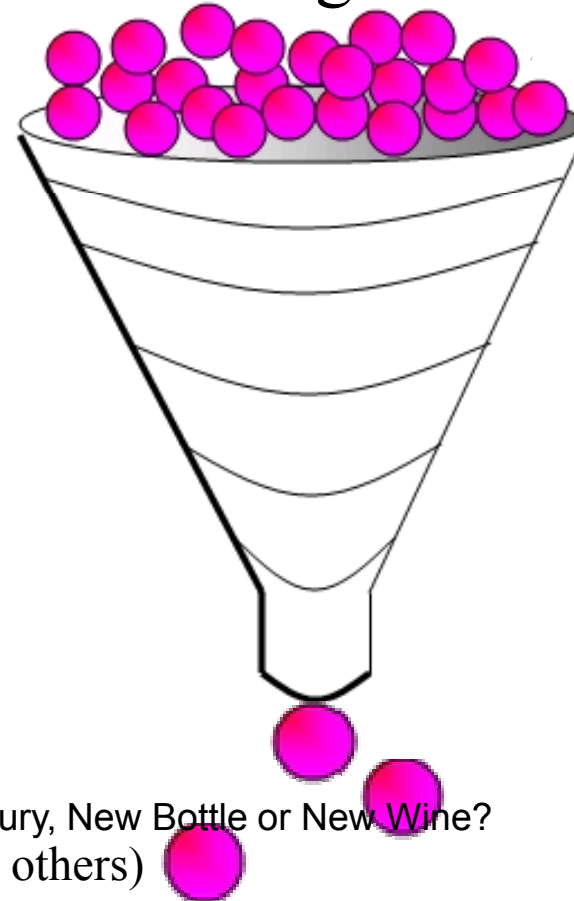
Many Many Tasks: Identifying Potential Drug Targets



Protein
target(s) x



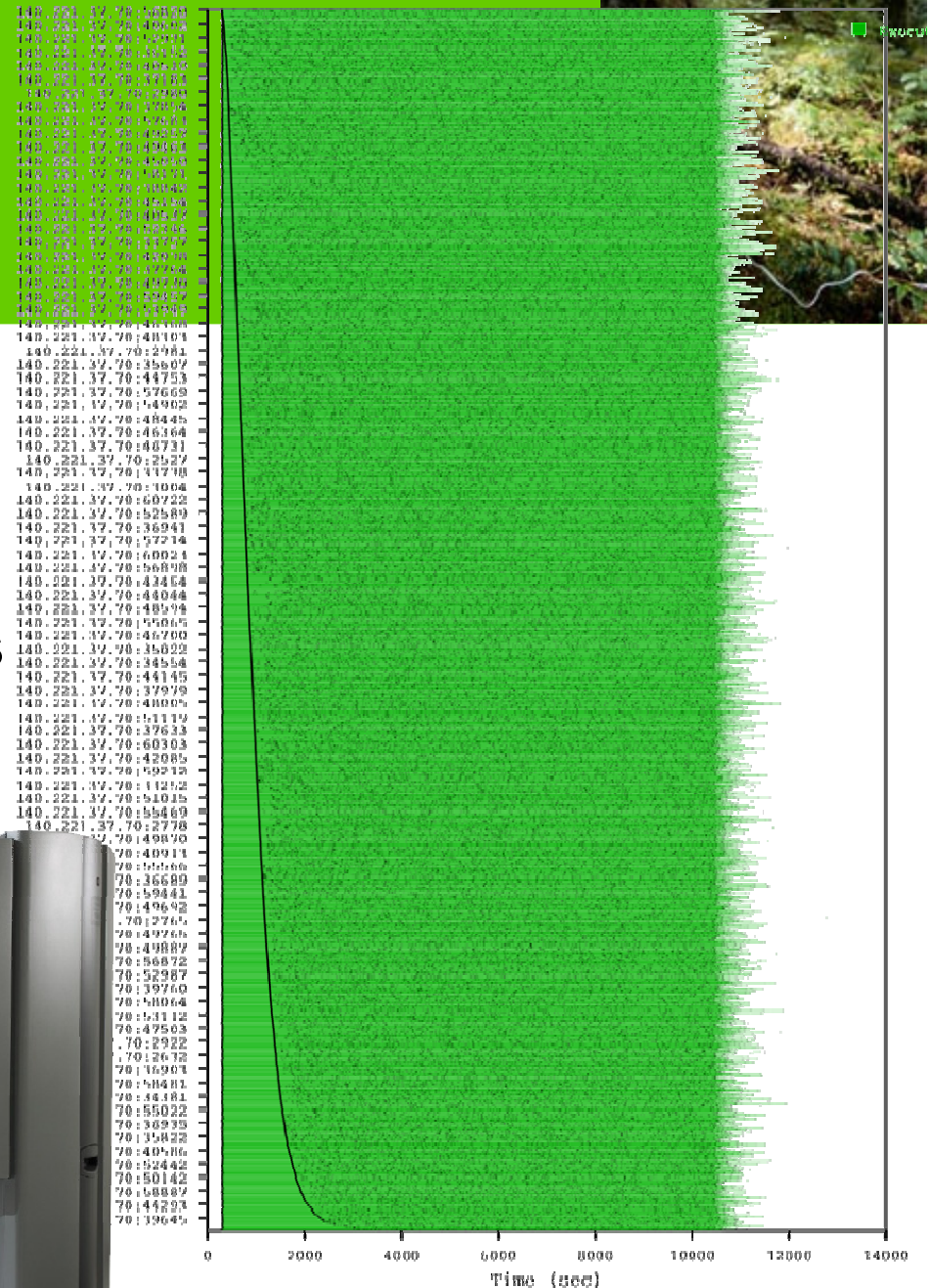
2M+ ligands



Scientific Workflow Systems for 21st Century, New Bottle or New Wine?
(Mike Kubal, Benoit Roux, and others)

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%



DOCK on the BG/P



CPU cores: 118784

Tasks: 934803

Elapsed time: 2.01 hours

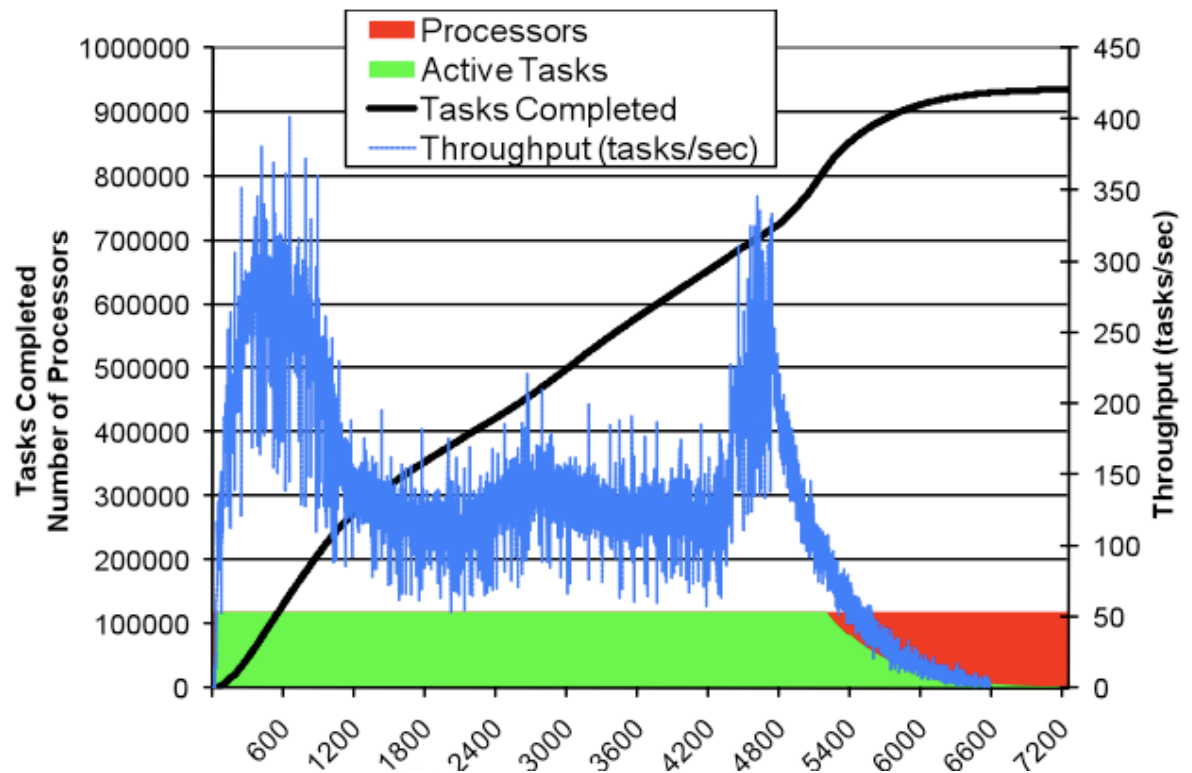
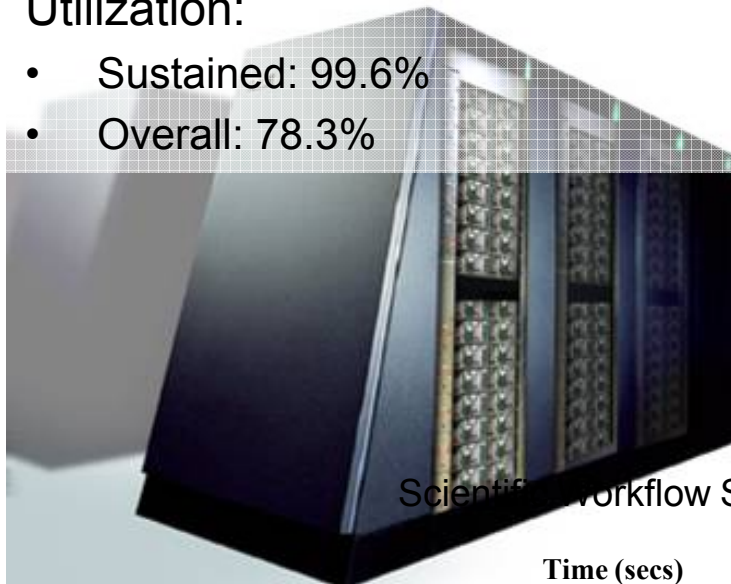
Compute time: 21.43 CPU years

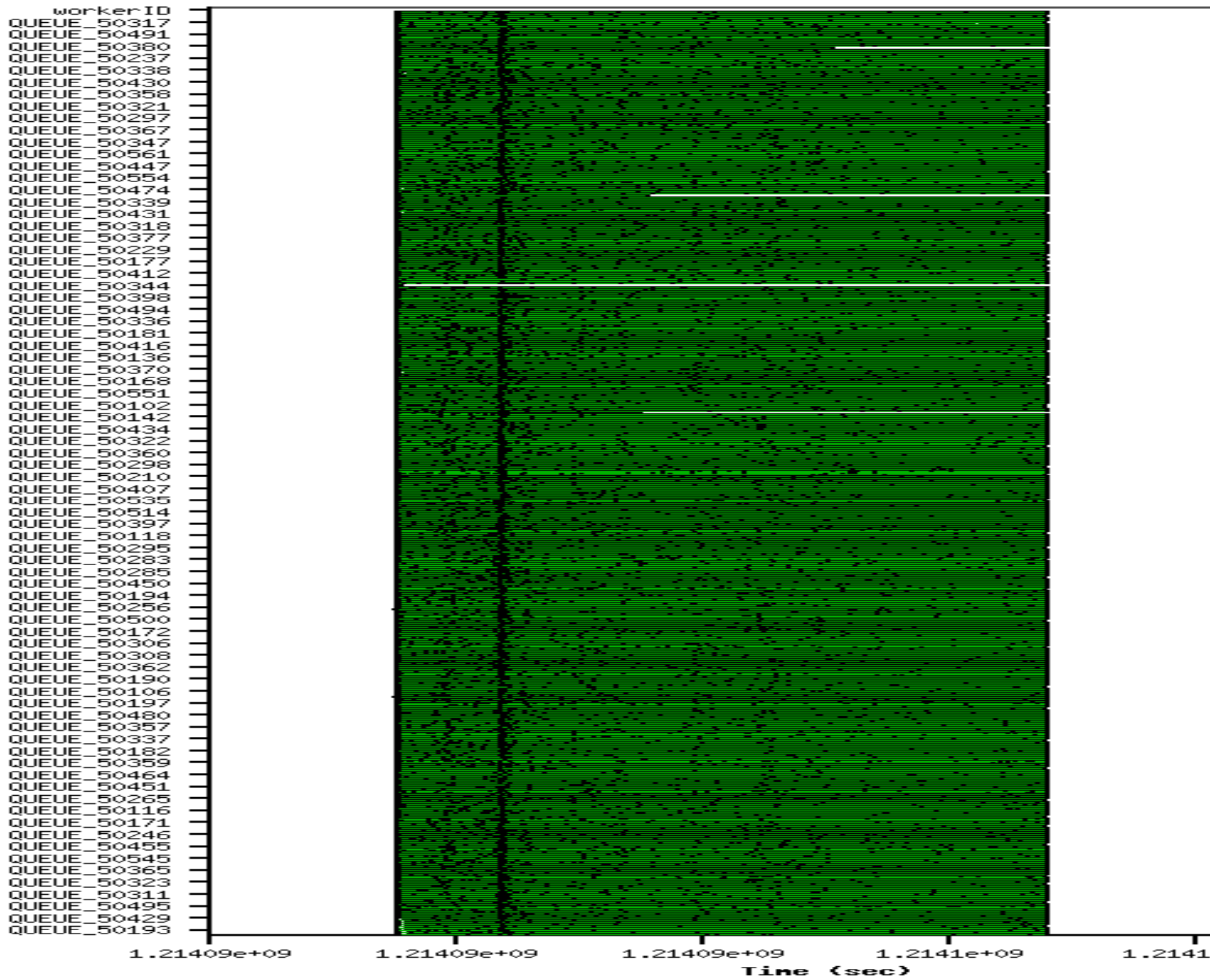
Average task time: 667 sec

Relative Efficiency: 99.7%
(from 16 to 32 racks)

Utilization:

- Sustained: 99.6%
- Overall: 78.3%

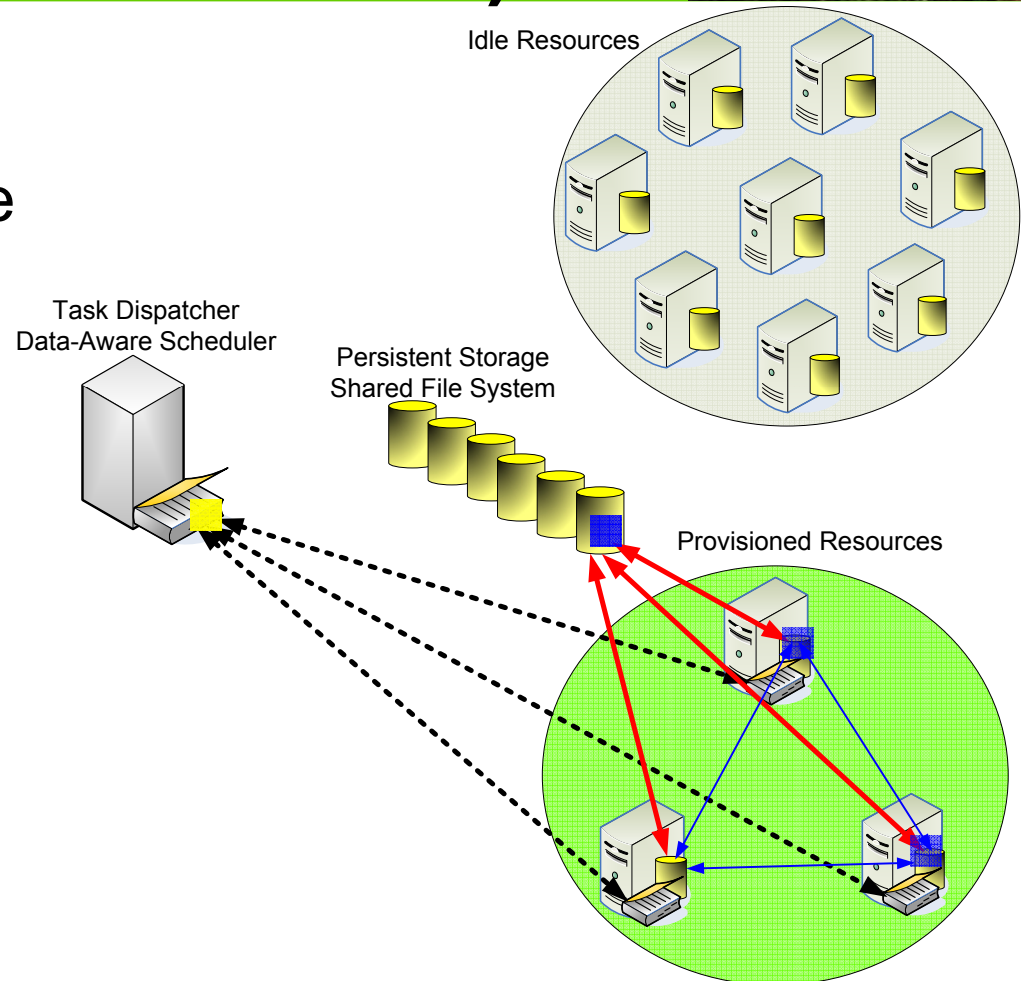




Support for Data Intensive Applications (Falkon and Data Diffusion)



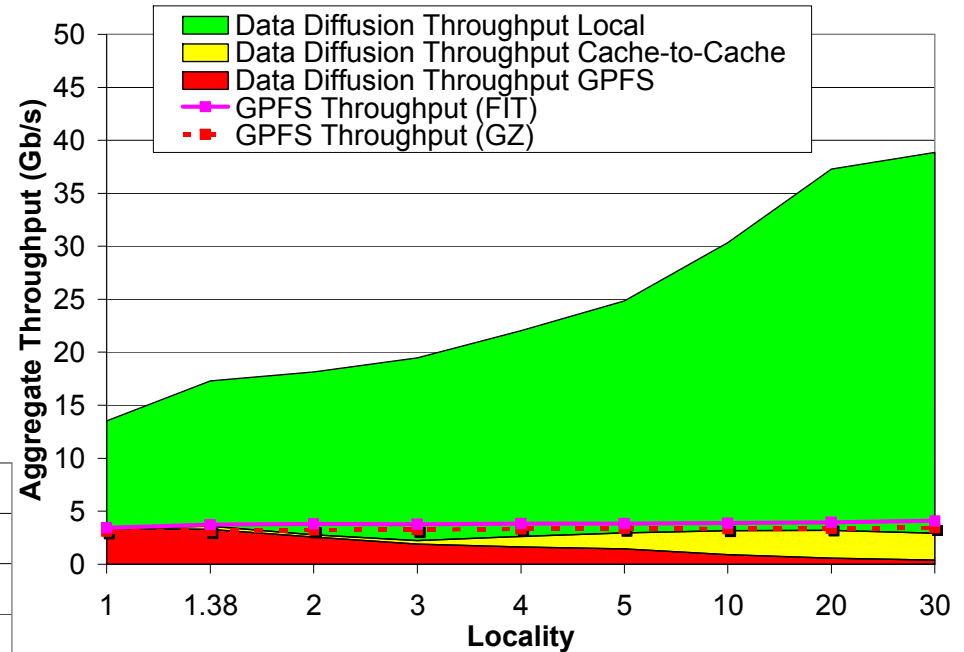
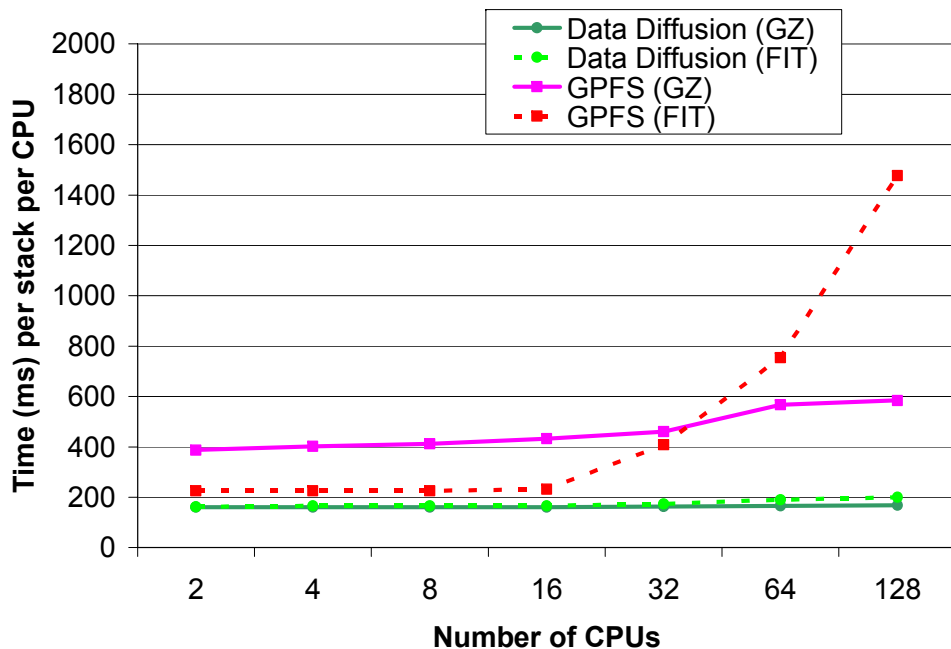
- Resource acquired in response to demand
- Data and applications diffuse from archival storage to newly acquired resources
- Resource “caching” allows faster responses to subsequent requests
 - Cache Eviction Strategies: RANDOM, FIFO, LRU, LFU
- Resources are released when demand drops



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



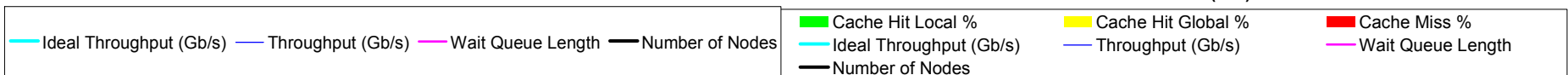
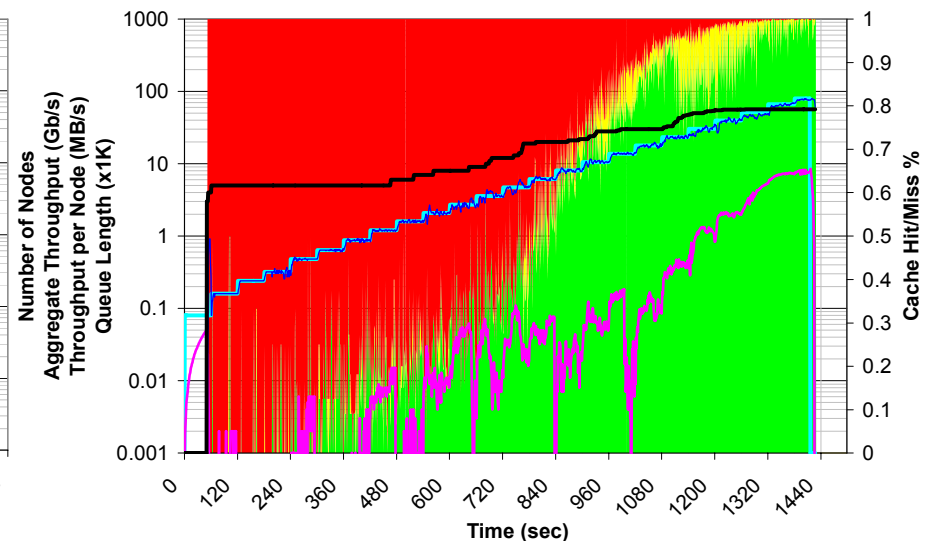
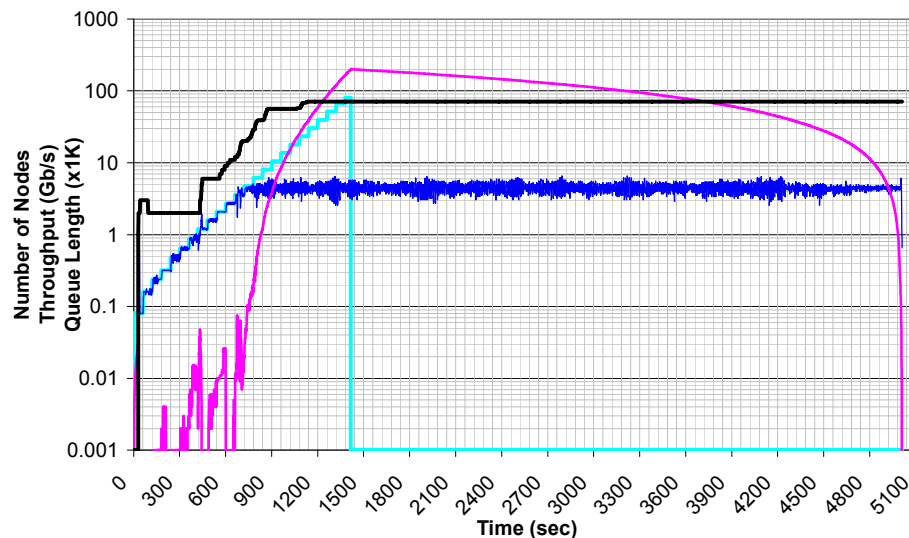
← High data locality
– Near perfect scalability

st Century, New Bottle or New Wine?

Data Diffusion: Data-Intensive Workload



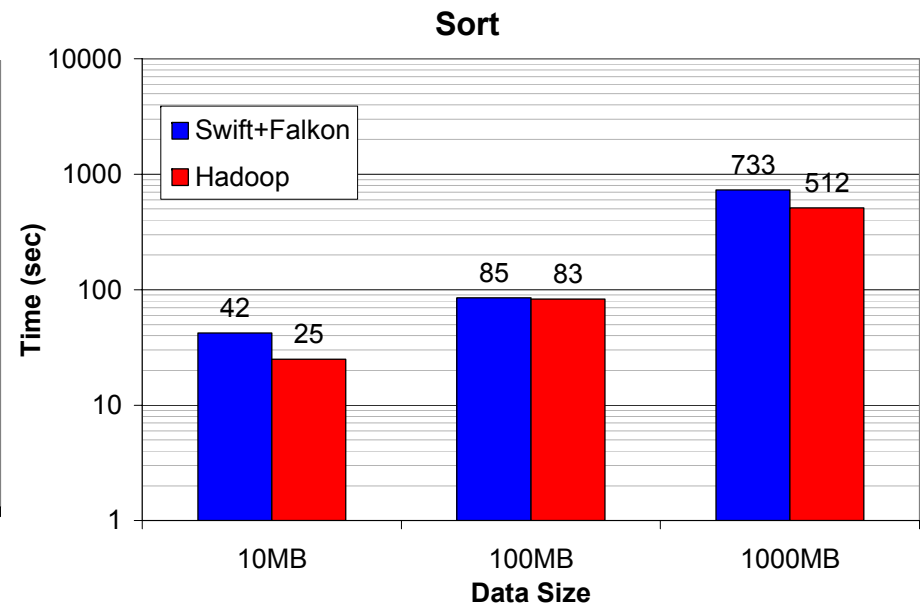
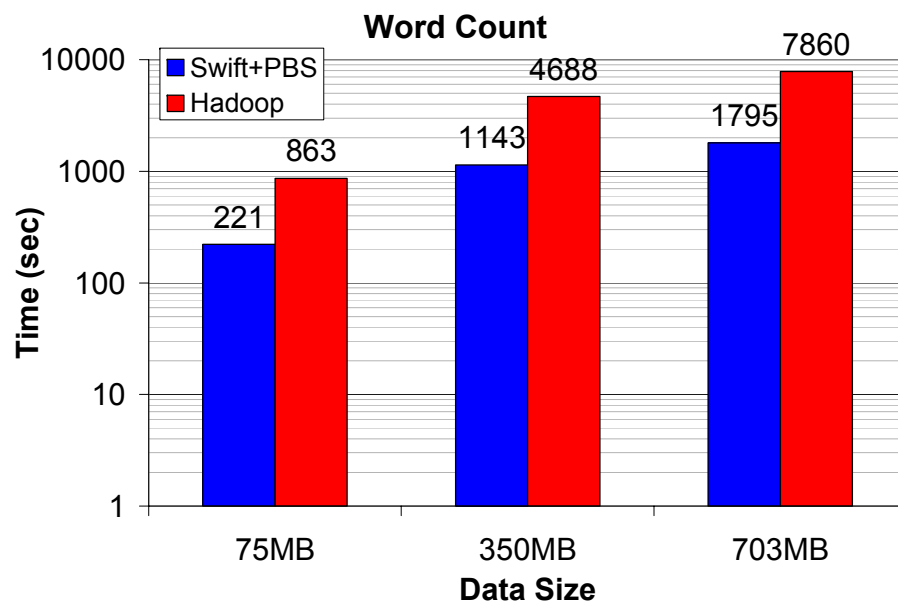
- 250K tasks on 128 processors
 - 10MB read, 10ms compute
- Comparing GPFS with data diffusion
 - 5011 sec vs. 1427 sec (ideal is 1415 sec)



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
 - Make use of idle resources from “supercomputers” via backfilling
 - Costs to run “supercomputers” per FLOP is among the best
 - BG/P: 0.35 gigaflops/watt (**higher is better**)
 - SiCortex: 0.32 gigaflops/watt
 - BG/L: 0.23 gigaflops/watt
 - x86-based HPC systems: an order of magnitude lower
- Loosely coupled apps do not require specialized system software
- Shared file systems are good for all applications
 - They don’t scale proportionally with the compute resources
 - Data intensive applications don’t perform and scale well

Features Scientific Workflow Systems should Have!



- **Parallelism**
 - Support for both explicit and implicit parallelism
- **Performance and Scalability**
 - Million to billions of tasks
 - Handle 100s~1000s of tasks/sec
- **Data management**
 - Reduce reliance on shared file systems
 - Scale with processing power
 - Data-aware scheduling
- **Reliability**
 - Self healing
 - Efficient and scalable monitoring
- **Provenance**

Solutions

(we have experience with)



- **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

MTAGS

Workshop on Many-Task Computing on Grids and Supercomputers

co-located with ACM/IEEE SC08 (International Conference for High Performance, Networking, Storage and Analysis)
Austin, Texas -- November 17th, 2008

[Home](#)

[Call for Papers](#)

[Program Committee](#)

[Important Dates](#)

[Paper Submission](#)

[Venue](#)

[Registration](#)

[Workshop Program](#)

Important Dates

Papers Due:	August 15th, 2008
Notification of Acceptance:	October 1st, 2008
Camera Ready Papers Due:	October 15th, 2008
Workshop Date:	November 17th, 2008

Committee Members

Workshop Chairs

Yong Zhao, Microsoft
Ian Foster, University of Chicago & Argonne National Laboratory
Ioan Raicu, University of Chicago

Technical Committee

Ian Foster, University of Chicago & Argonne National Laboratory
Dan Ardelean, Google
Bob Grossman, University of Illinois at Chicago
Indranil Gupta, University of Illinois at Urbana Champaign
Tevfik Kosar, Louisiana State University

Handling Megajobs BOF at SC08



- More and more people need to run thousands to millions of closely related jobs that are associated with individual projects. Scientists seek convenient means to specify and manage many jobs, arranging inputs, aggregating outputs, identifying successful and failed jobs and repairing failures. System administrators seek methods to process extraordinary numbers of jobs for multiple users without overwhelming queuing systems or disrupting fair-share usage policies. And, grid developers are producing a new generation of queuing and scheduling systems as well as auxiliary systems for use with existing queuing and scheduling systems. This Birds-of-feather session provides a venue for the exchange of information about processing large numbers of jobs. Short presentations of an invited sample of projects will be followed by discussion.
- For more information, contact:
 - Marlon Pierce: mpierce@cs.indiana.edu
 - Dick Repasky: rrepasky@indiana.edu
 - Ioan Raicu: iraicu@cs.uchicago.edu