

PROBLEM DIAGNOSIS & VISUALIZATION

Jiaqi Tan, Soila Kavulya, Xinghao Pan, Mike Kasick, Keith Bare,
Eugene Marinelli, Rajeev Gandhi

Priya Narasimhan

Carnegie Mellon University



My Background

- 15 years working in fault-tolerant middleware
- Developed transparent fault-tolerant middleware
 - Eternal, Immune
 - Standards: Fault-Tolerant CORBA, Portable Interceptors for CORBA
- Hard problems addressed
 - Strongly consistent replication & recovery
 - Application-transparent fault tolerance
 - Resolving conflicts between real-time & fault tolerance
 - Overcoming nondeterminism, unrealistic assumptions
 - Zero-downtime large-scale software upgrades

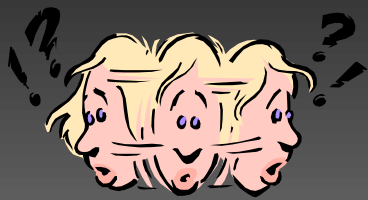
And Onto Automated Problem Diagnosis

- Diagnosing problems
 - Creates major headaches for administrators
 - Worsens as scale and system complexity grows
- Goal: automate it and get proactive
 - Failure detection and prediction
 - Problem determination (or “fingerprinting”)
- How: Instrumentation plus statistical analysis



Challenges in Problem Analysis

- Challenging in large-scale networked environment
 - Can have multiple failure manifestations with a single root cause
 - Can have multiple root causes for a single failure manifestation
 - Problems and/or their manifestations can “travel” among communicating components
 - A lot of information from multiple sources – what to use? what to discard?
- Automated fingerprinting
 - Automatically discover faulty node in a distributed system



Exploration

- Current explorations

- *Hadoop*

- Open-source implementation of Map/Reduce (Yahoo!), popular cloud-computing platform

- *PVFS*

- High-performance file system (Argonne National Labs)

- *Lustre*

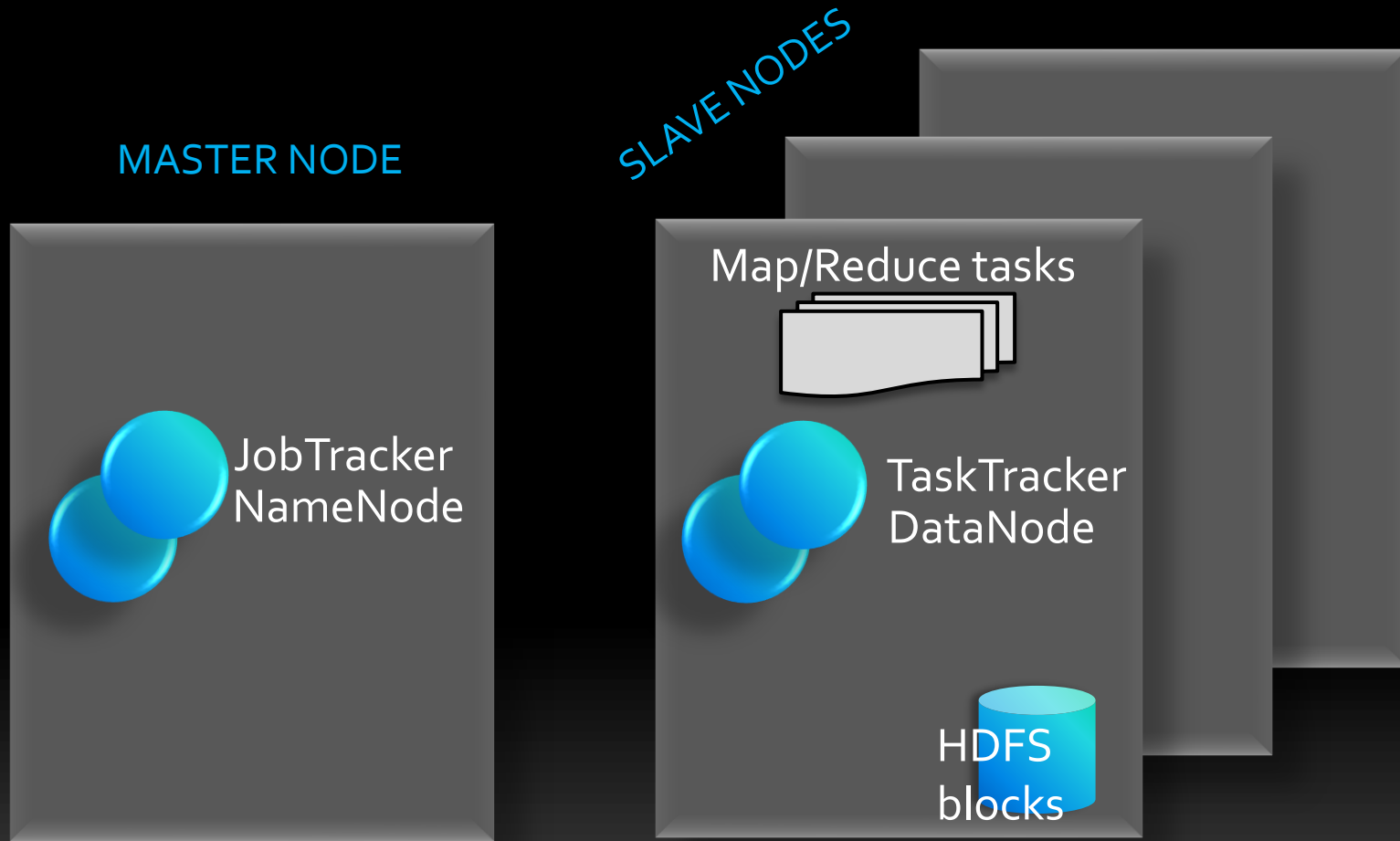
- High-performance file system (Sun Microsystems)

- Studied

- Various types of problems
 - Various kinds of instrumentation
 - Various kinds of data-analysis techniques



Hadoop 101



Why?

- Hadoop is fault-tolerant
 - Heartbeats: detect lost nodes
 - Speculative re-execution: recover work due to lost/laggard nodes
- Hadoop's fault-tolerance can mask performance problems
 - Nodes alive but slow
- Target failures for our diagnosis
 - Performance degradations (slow, hangs)

Goals, Non-Goals

- Diagnose faulty Master/Slave node to user/admin
 - Target production environment
 - Don't instrument Hadoop or applications additionally
 - Use Hadoop logs as-is (*white-box strategy*)
 - Use OS-level metrics (*black-box strategy*)
 - Work for various workloads and under workload changes
 - Support online and offline diagnosis
-
- Non-goals (for now)
 - Tracing problem to offending line of code

Target Hadoop Clusters

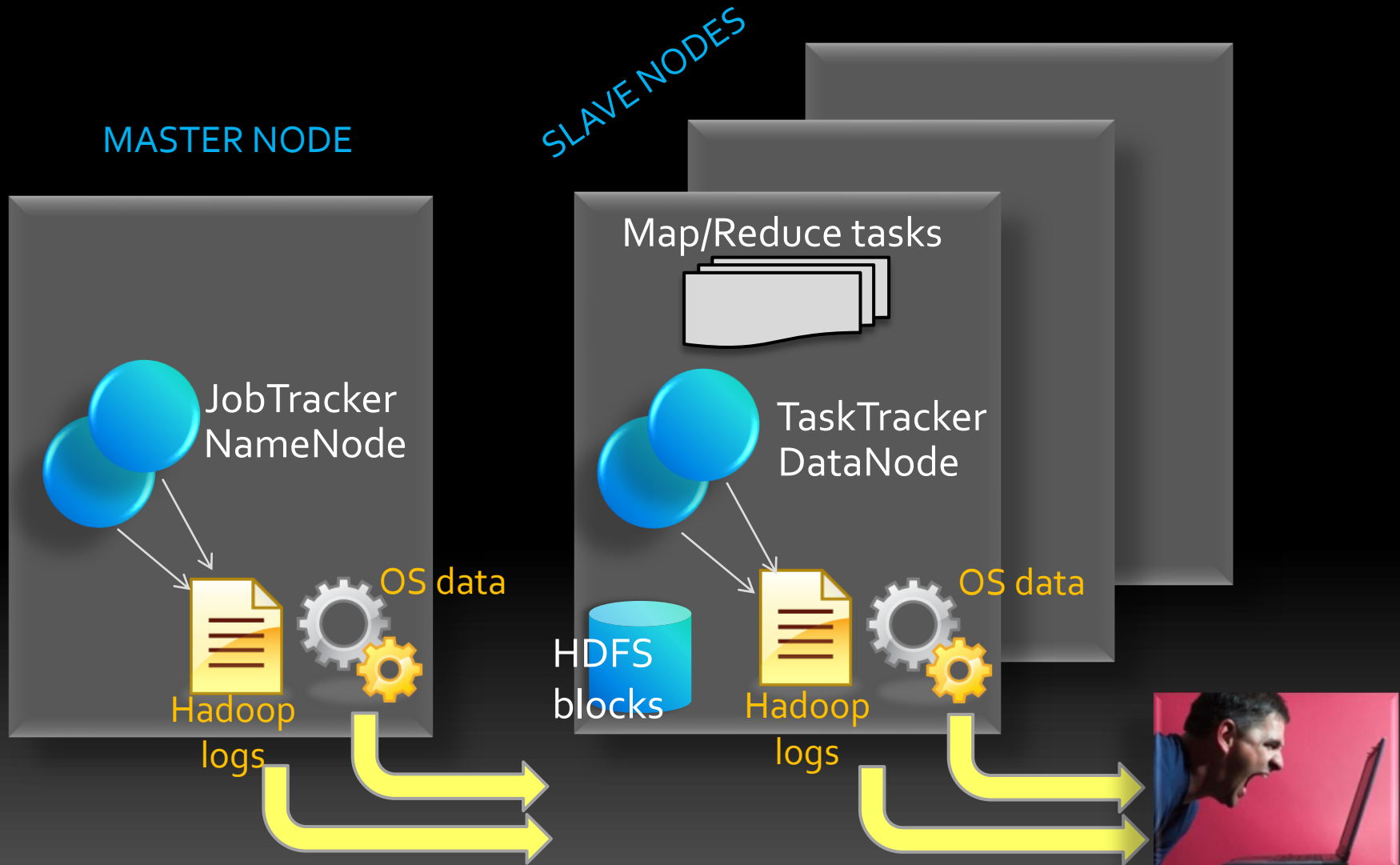
- 4000-processor Yahoo!'s M45 cluster
 - Production environment (managed by Yahoo!)
 - Offered to CMU as free cloud-computing resource
 - Diverse kinds of real workloads, problems in the wild
 - Massive machine-learning, language/machine-translation
 - Permission to harvest all logs and OS data each week
- 100-node Amazon's EC2 cluster
 - Production environment (managed by Amazon)
 - Commercial, pay-as-you-use cloud-computing resource
 - Workloads under our control, problems injected by us
 - gridmix, nutch, pig, sort, randwriter
 - Can harvest logs and OS data of only our workloads

Some Performance Problems Studied

	Fault	Description
Resource contention	CPU hog	External process uses 70% of CPU
	Packet-loss	5% or 50% of incoming packets dropped
	Disk hog	20GB file repeatedly written to
	Disk full	Disk full
Application bugs Source: Hadoop JIRA	HADOOP-1036	Maps hang due to unhandled exception
	HADOOP-1152	Reduces fail while copying map output
	HADOOP-2080	Reduces fail due to incorrect checksum
	HADOOP-2051	Jobs hang due to unhandled exception
	HADOOP-1255	Infinite loop at Nameode

Studied Hadoop Issue Tracker (JIRA) from Jan-Dec 2007

Hadoop: Instrumentation



How About Those Metrics?

- **White-box** metrics (from Hadoop logs)
 - Event-driven (based on Hadoop's activities)
 - *Durations*
 - Map-task durations, Reduce-task durations, ReduceCopy-durations, etc.
 - System-wide **dependencies** between tasks and data blocks
 - **Heartbeat** information: Heartbeat rates, Heartbeat-timestamp skew between the Master and Slave nodes
- **Black-box** metrics (from OS /proc)
 - 64 different time-driven metrics (sampled every second)
 - Memory used, context-switch rate, User-CPU usage, System-CPU usage, I/O wait time, run-queue size, number of bytes transmitted, number of bytes received, pages in, pages out, page faults

Log-Analysis Approach

- [S](#)ALSA: [A](#)nalyzing [L](#)ogs as [S](#)tate Machines [USENIX WASL 2008]
- Extract state-machine views of execution from Hadoop logs
 - Distributed control-flow view of logs
 - Distributed data-flow view of logs
- Diagnose failures based on statistics of these extracted views
 - Control-flow based diagnosis
 - Control-flow + data-flow based diagnosis
- Perform analysis incrementally so that we can support it online



Applying SALSA to Hadoop

Data-flow view:
transfer of data
to other nodes



Control-flow
view: state
orders, durations

Map

[t] Launch Map task
:
[t] Copy Map outputs
:
[t] Map task done

Map outputs to
Reduce tasks on
other nodes

[t] Launch Reduce task
:
[t] Reduce is idling, waiting for Map
outputs
:
[t] Repeat until all Map outputs copied
[t] Start Reduce Copy
(of completed Map output)
:
[t] Finish Reduce Copy
[t] Reduce Merge Copy

Incoming Map outputs
for this Reduce task

Reduce
Idle

Reduce
Copy

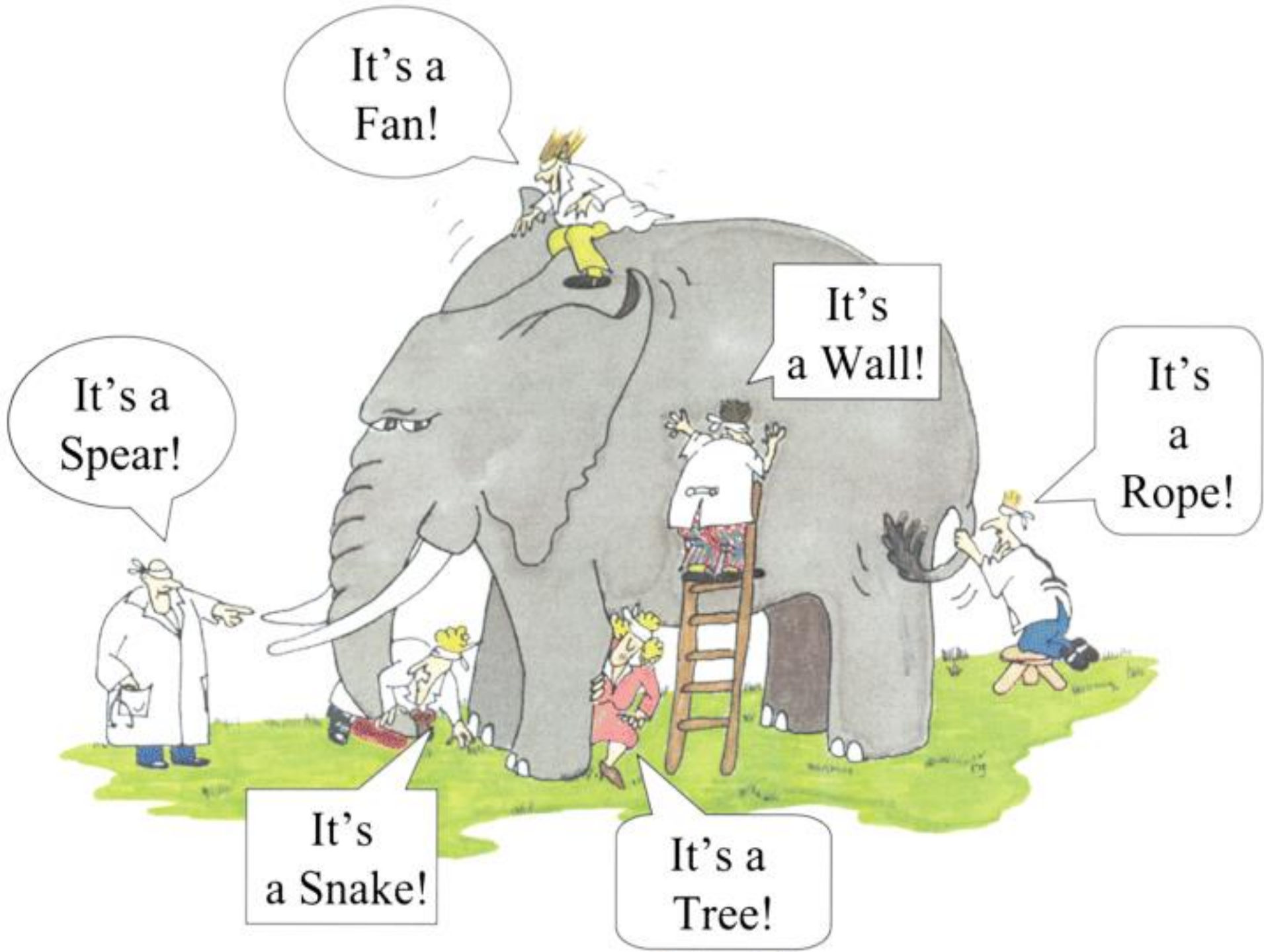
Reduce
Merge
Copy

Distributed Control+Data Flow

- Distributed control-flow
 - Causal flow of task execution across cluster nodes, i.e., Reduces waiting on Maps via Shuffles
- Distributed data-flow
 - Data paths of Map outputs shuffled to Reduces
 - HDFS data blocks read into and written out of jobs
- **Job-centric data-flows:** Fused Control+Data Flows
 - Correlate paths of data and execution
 - Create conjoined causal paths from data source before, to data destination after, processing
 - Helps to trace correlated performance problems

What Else Do We Do?

- Analyze black-box data with similar intuition
 - Derive PDFs and use a clustering approach
 - Distinct behavior profiles of metric correlations
 - Compare them across nodes
 - Technique called Ganesha [*HotMetrics 2009*]
- Analyze heartbeat traffic
 - Compare heartbeat durations across nodes
 - Compare heartbeat-timestamp skews across nodes
- Different metrics, different viewpoints, different algorithms



Putting the Elephant Together

JobTracker
Durations
views

TaskTracker
heartbeat
timestamps

TaskTracker
Durations
views

JobTracker
heartbeat
timestamps

Job-centric
data flows

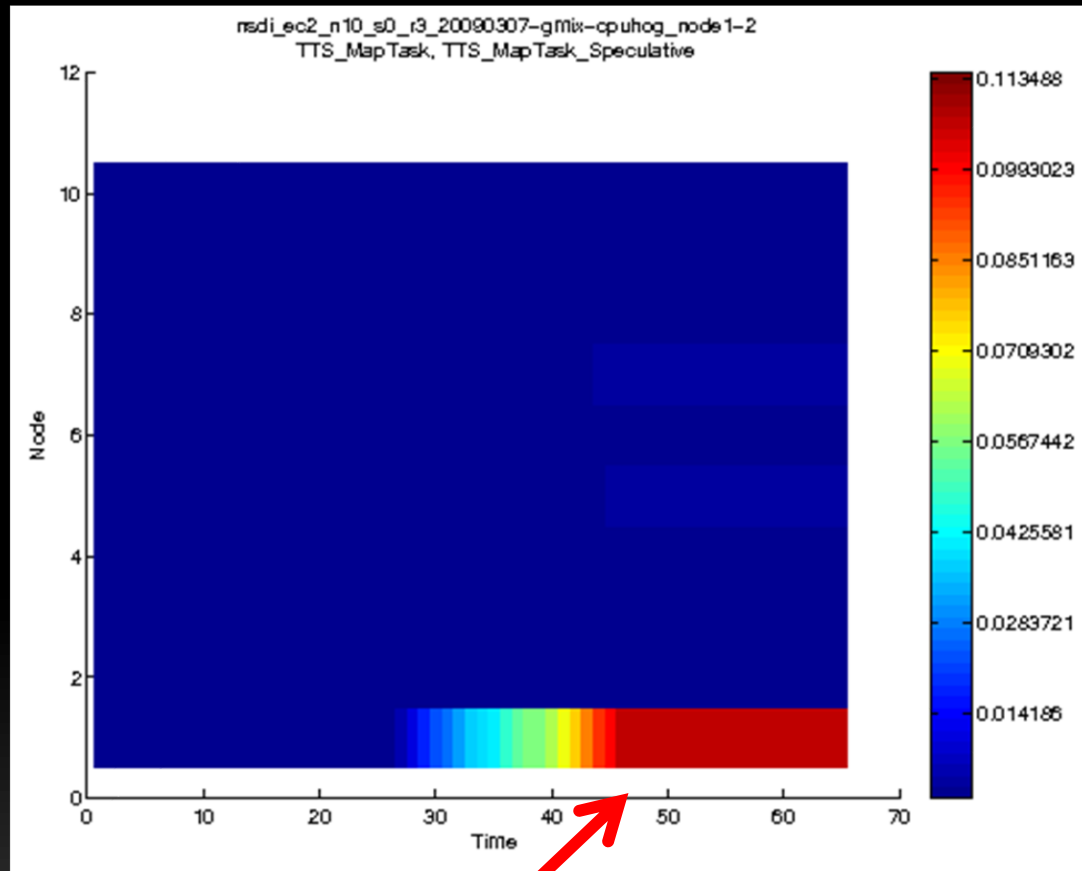
Black-box
resource
usage

BliMEy: Blind Men and the Elephant Framework
[CMU-CS-09-135]

Visualization

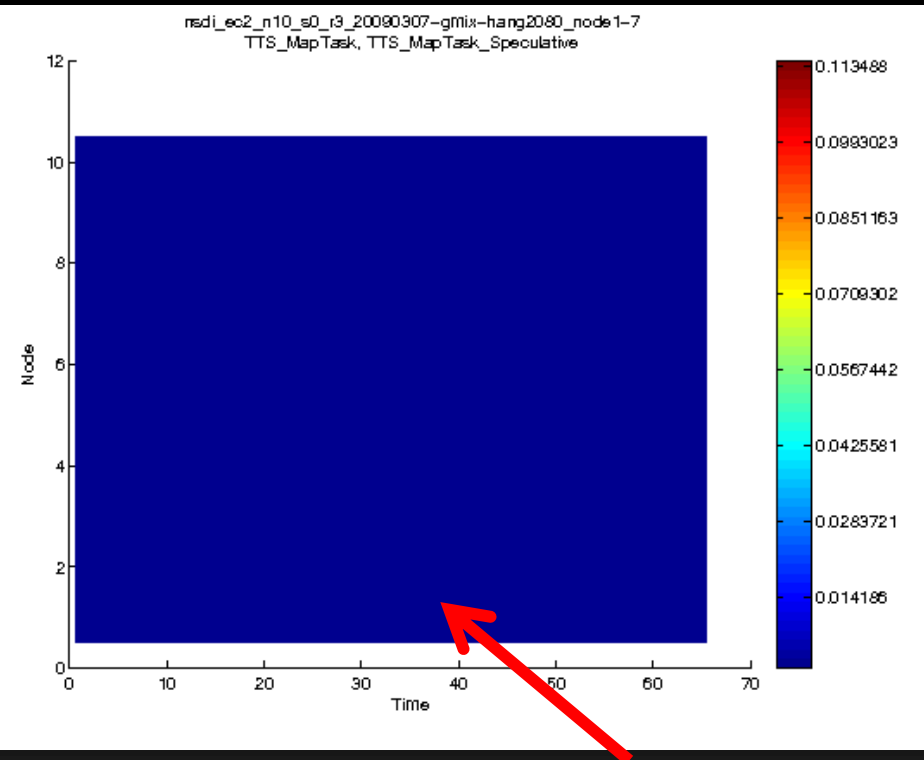
- To uncover Hadoop's execution in an insightful way
- To reveal outcome of diagnosis on sight
- To allow developers/admins to get a handle as the system scales
- Value to programmers [*HotCloud 2009*]
 - Allows them to spot issues that might assist them in restructuring their code
 - Allows them to spot faulty nodes

Visualization(*heatmaps*)

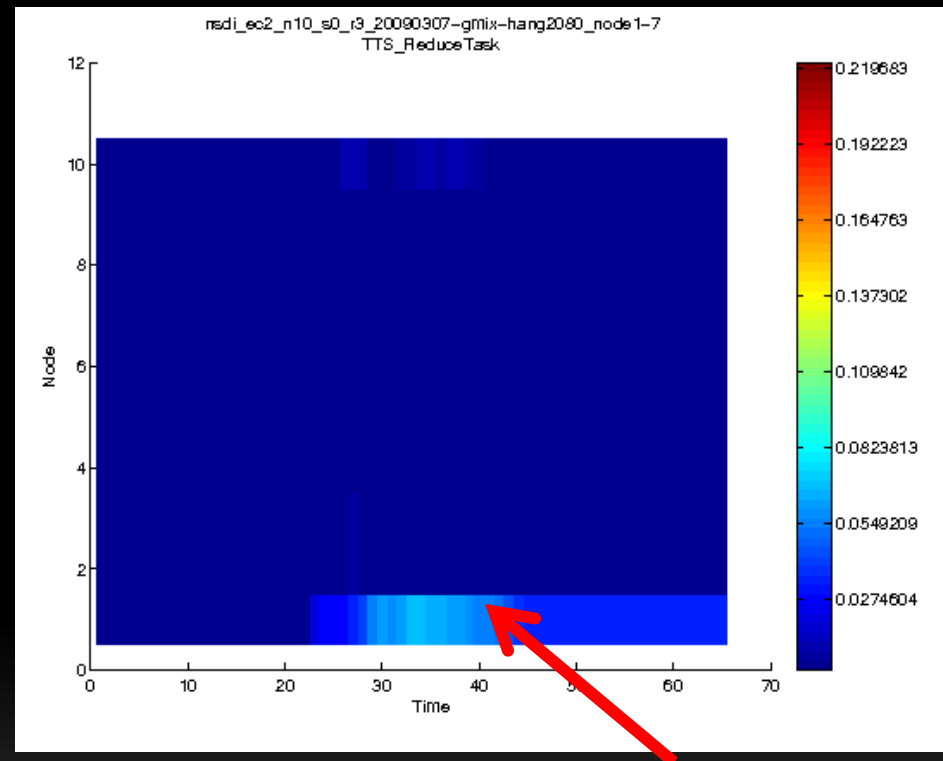


CPU Hog on node 1
visible on Map-task durations

Visualization (*heatmaps*)

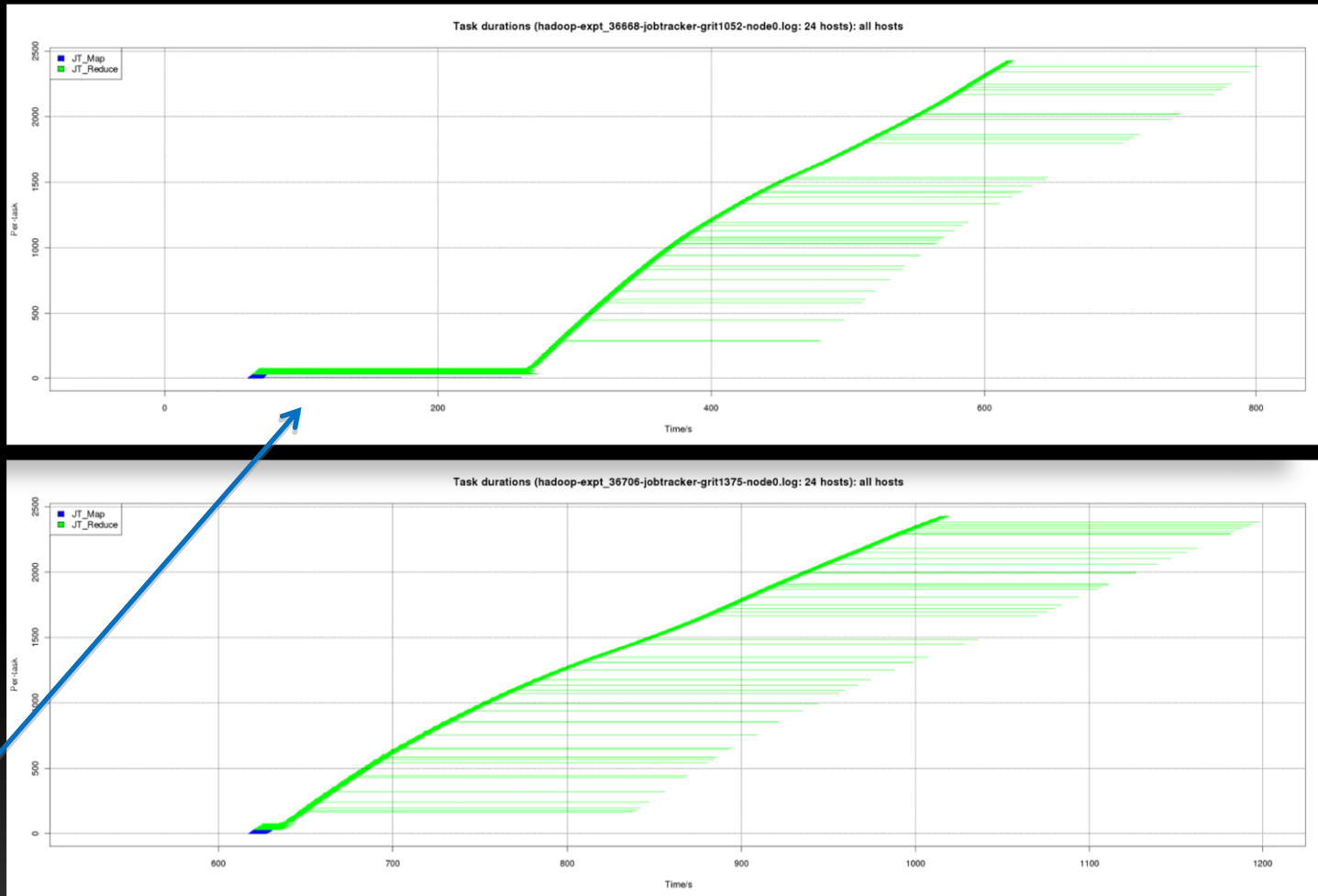


TaskTracker hang on node 1
not visible on Map-task durations



TaskTracker hang on node 1
visible on Reduce-task
durations

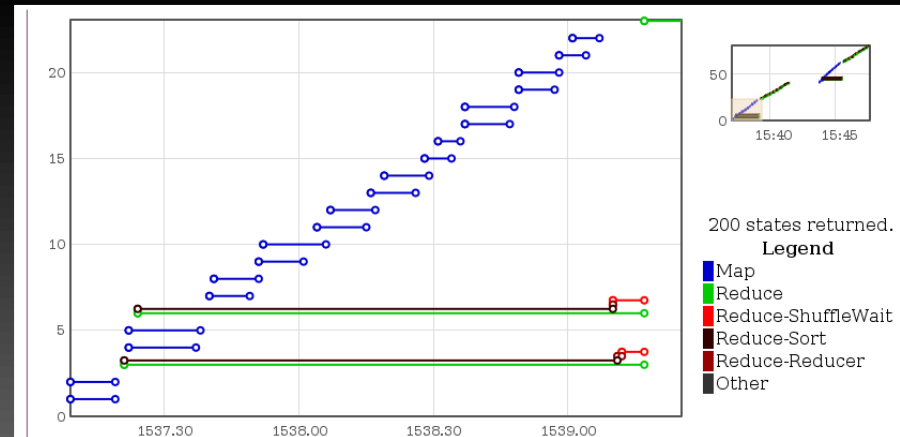
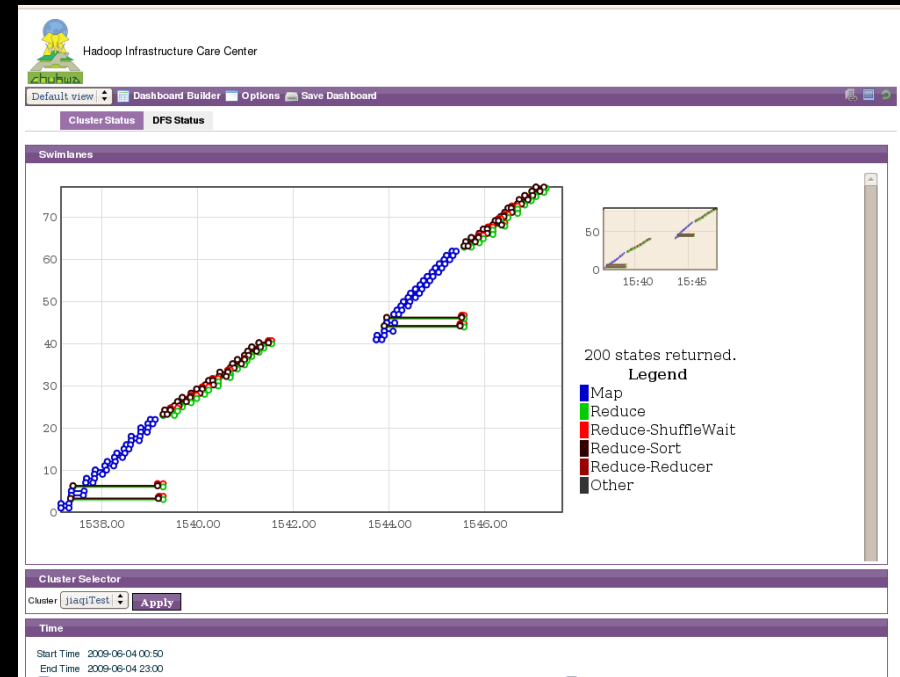
Visualizations (*swimLanes*)



Long-tailed map
Delaying overall
job completion time

Current Developments

- State-machine extraction + visualization being implemented for the Hadoop Chukwa project
 - Collaboration with Yahoo!
- Web-based visualization widgets for HICC (Hadoop Infrastructure Care Center)
- “Swimlanes” currently available in Chukwa trunk (CHUKWA-279)



Briefly: Diagnosis for PVFS

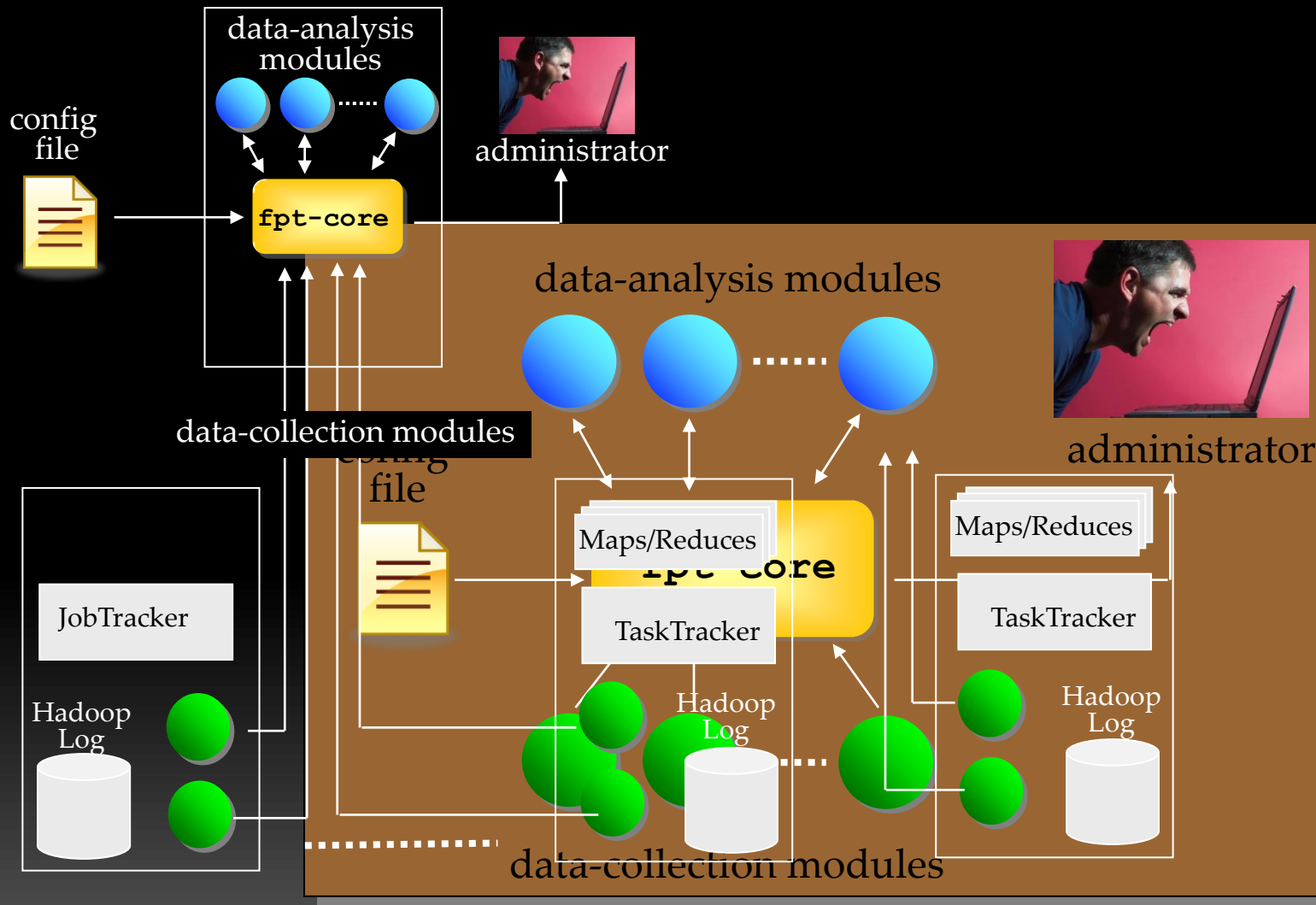
- PVFS: High-performance file system
- Focus on different kinds of problems
 - Disk-related and network-related problems
 - Motivated by anecdotal experiences of Argonne National Labs researchers
- Exploit the peer-similarity principle across I/O servers and metadata servers
- Diagnosis based on different instrumentation
 - Black-box OS performance metrics, system-call traces

Briefly: Online Fingerprinting

- **ASDF**: Automated System for Diagnosing Failures
 - Can incorporate any number of different data sources
 - Can use any number of analysis techniques to process this data
- Can support online or offline analyses for Hadoop
- Currently plugging in our white-box & black-box algorithms



Example: ASDF for Hadoop



Briefly: Diagnosis-Driven Recovery

- The Problem: Naïve recovery might not mitigate the problem
 - Problem might escalate
 - Problem might still remain and show up later
- What if recovery were informed by diagnosis, in real-time?
 - Preemptively migrate tasks to non-faulty servers, based on monitored data

Hard Problems

- Understanding the limits of black-box fingerprinting
 - What failures are outside the reach of a black-box approach?
 - What are the limits of “peer” comparison?
 - What other kinds of black-box instrumentation exist?
- Scalability
 - Scaling to run across large systems and understanding “growing pains”
- Visualization
 - Helping system administrators visualize problem diagnosis
- Trade-offs
 - More instrumentation and more frequent data can improve accuracy of diagnosis, but at what performance cost?
- Virtualized environments
 - Do these environments help/hurt problem diagnosis?

Summary

- Automated problem diagnosis
- Current targets: Hadoop, PVFS, Lustre
- Initial set of failures
 - Real-world bug databases, problems in the wild
- Short-term: Transition techniques into Hadoop code-base working with Yahoo!
- Long-term
 - Scalability, scalability, scalability,
 - Expand fault study
 - Improve visualization, working with users
- Additional details
 - *USENIX WASL 2008* (white-box log analysis)
 - *USENIX HotCloud 2009* (visualization)
 - *USENIX HotMetrics 2009* (black-box metric analysis)
 - *HotDep 2009* (black-box analysis for PVFS)

FOR MORE INFORMATION:

[HTTP://WWW.ECE.CMU.EDU/~FINGERPOINTING](http://www.ece.cmu.edu/~fingerprinting)

priya@cs.cmu.edu

