



---

# **A Preliminary Report on Red Storm RAS Performance**

**(RAS = Reliability, Availability, and Serviceability)**

**Jon Stearley, Robert Ballance**

**{jrstear,raballa}@sandia.gov**

**Sandia National Laboratories**

**May 8, 2006**

**Cray Users Group (CUG-06)**



# Are these numbers Comparable?

---

Systems	CPUs	Reliability & Availability
ASCI Q	8,192	<b>MTBI: 6.5 hrs.</b> 114 unplanned outages/month. ◆ HW outage sources: storage, CPU, memory.
ASCI White	8,192	<b>MTBF: 5 hrs. (2001) and 40 hrs. (2003).</b> ◆ HW outage sources: storage, CPU, 3 <sup>rd</sup> -party HW.
NERSC Seaborg	6,656	<b>MTBI: 14 days. MTTR: 3.3 hrs.</b> ◆ SW is the main outage source. <b>Availability: 98.74%.</b>
PSC Lemieux	3,016	<b>MTBI: 9.7 hrs.</b> <b>Availability: 98.33%.</b>
Google	~15,000	<b>20 reboots/day; 2-3% machines replaced/year.</b> ◆ HW outage sources: storage, memory. <b>Availability: ~100%.</b>

MTBI: mean time between interrupts; MTBF: mean time between failures; MTTR: mean time to restore

Source: Daniel A. Reed, UNC (via Chung-Hsing Hsu, LANL)



# What is “RAS”???

---

“MTBI”

“MTBF”

“MTTR”

“UP”

“DOWN”

**Problem:**

**We DO NOT agree on terms!**

**This prevents accurate comparisons,  
obscures meaningful discussion,  
and delays significant improvements.**

**Solution:**

**Let’s agree on definitions,  
And develop a reference implementation.**

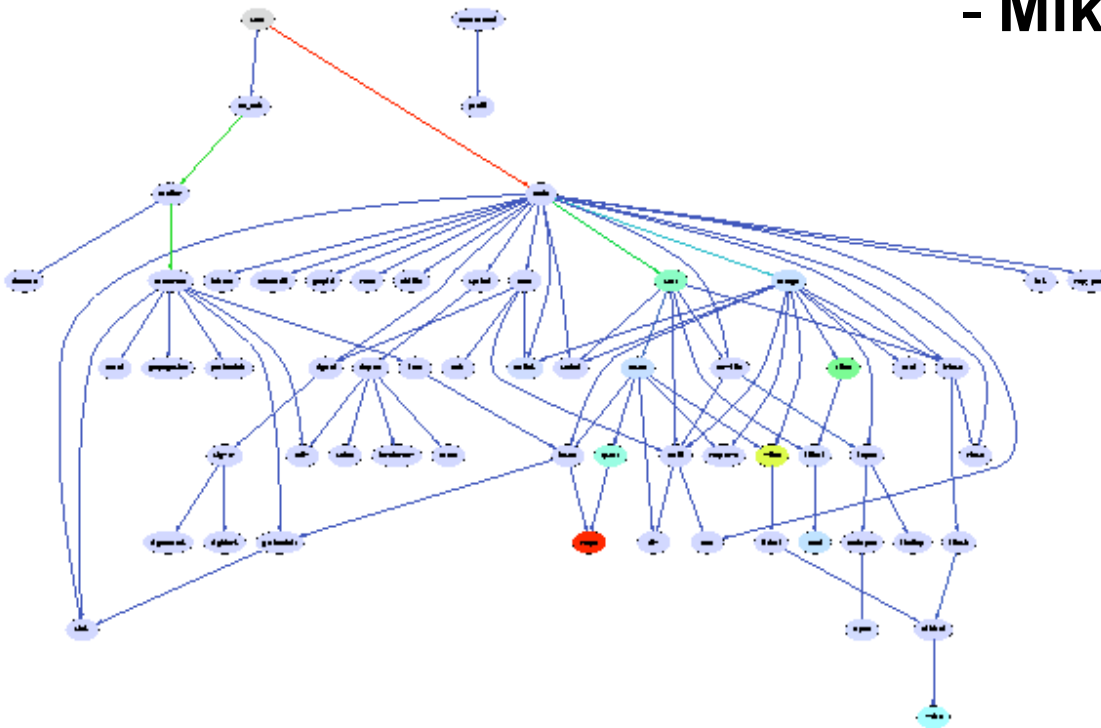
(Efforts are underway to establish a SNL/LLNL/LANL-endorsed “specification for defining and measuring high performance computing reliability, availability, and serviceability”.)



# Up or Down?

“A computer is in one of two situations. It is either known to be bad or it is in an unknown state.”

- Mike Levine (PSC)

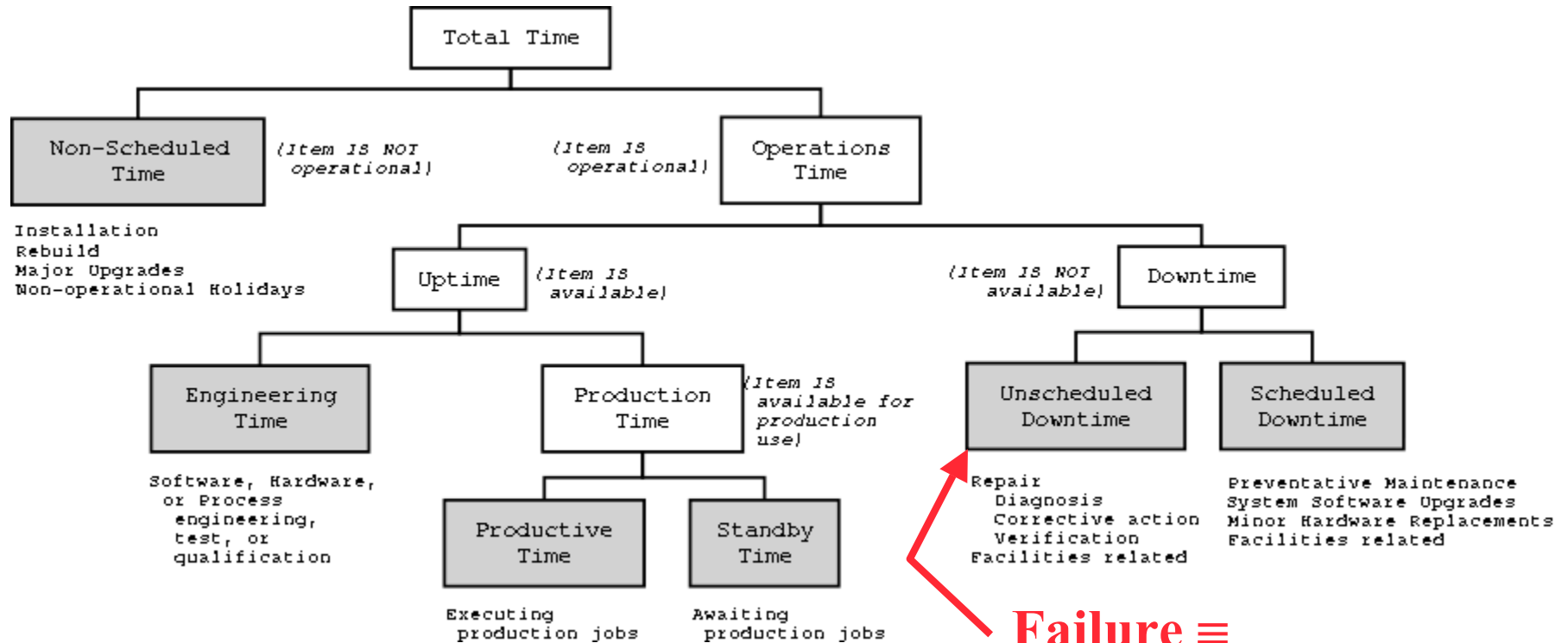


Is the  
“system”  
“up”  
(yet)?



# State Model

(adapted from SEMI-E10)



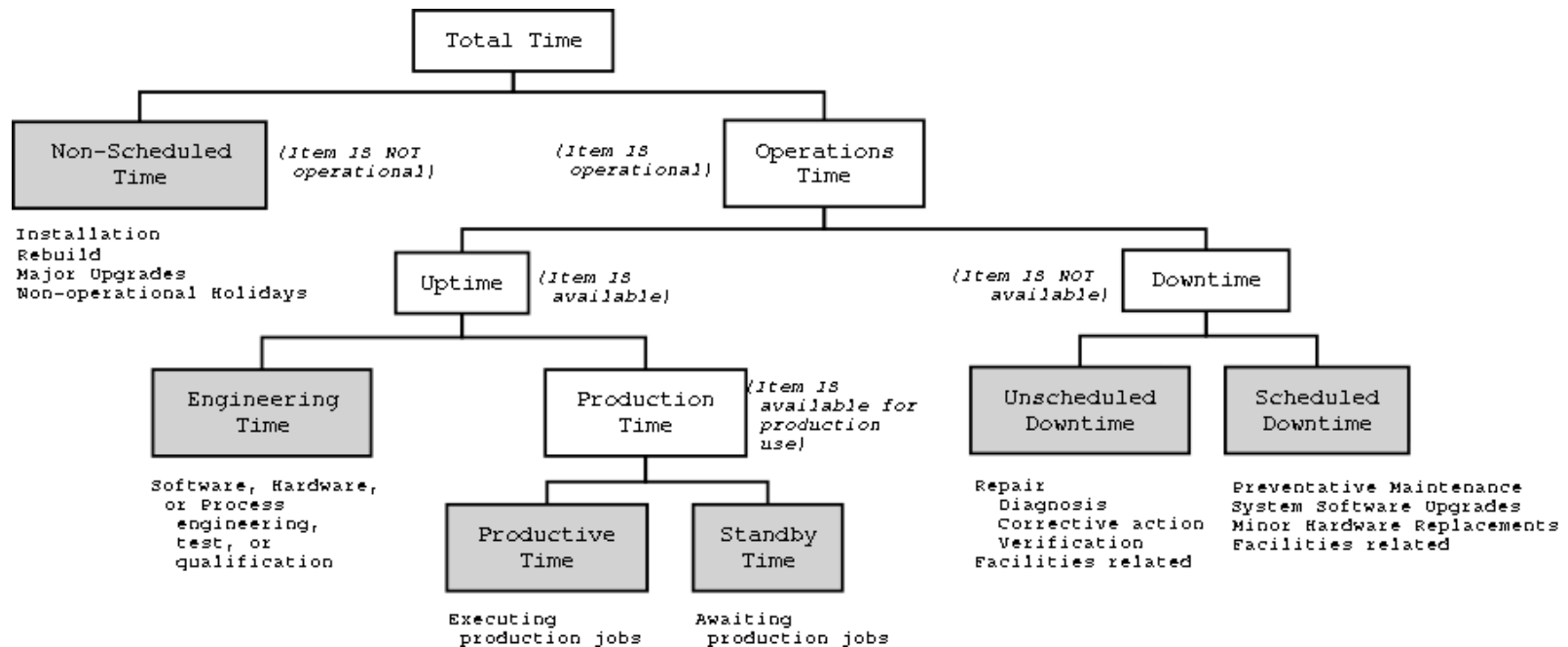
**Failure ≡  
unscheduled  
downtime**

Items are always in one of the six basic states (shaded).



# Availability

The fraction of a time period that an item is in a condition to perform its intended function (“available”). [IEEE]

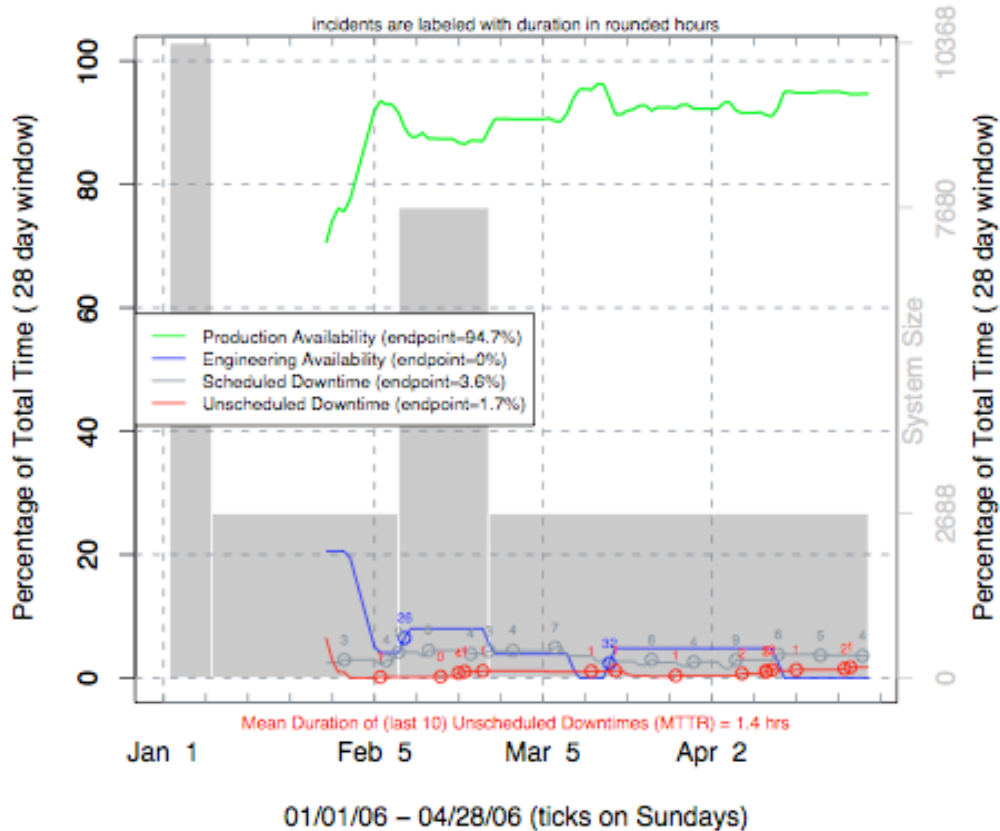


$$\text{Production Availability (\%)} = \frac{\text{Production Time}}{\text{Total Time}} * 100 \quad [\text{SEMI-E10}]$$

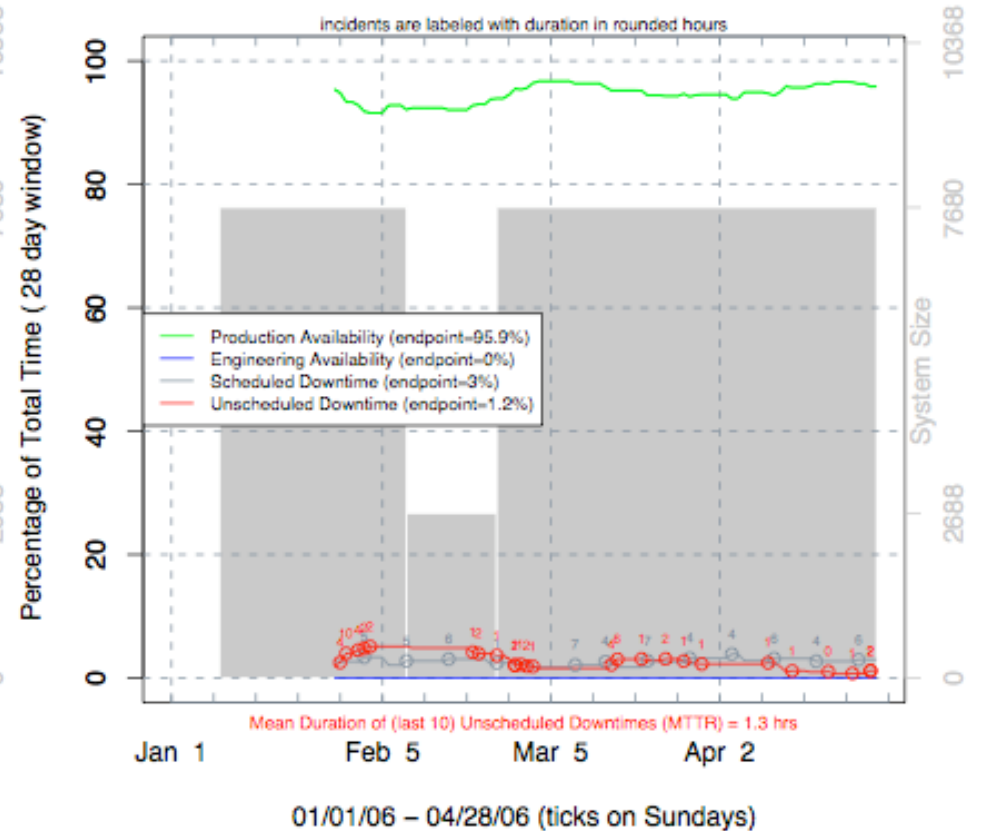


# System Availability

### RedStorm/Unclassified : Availability



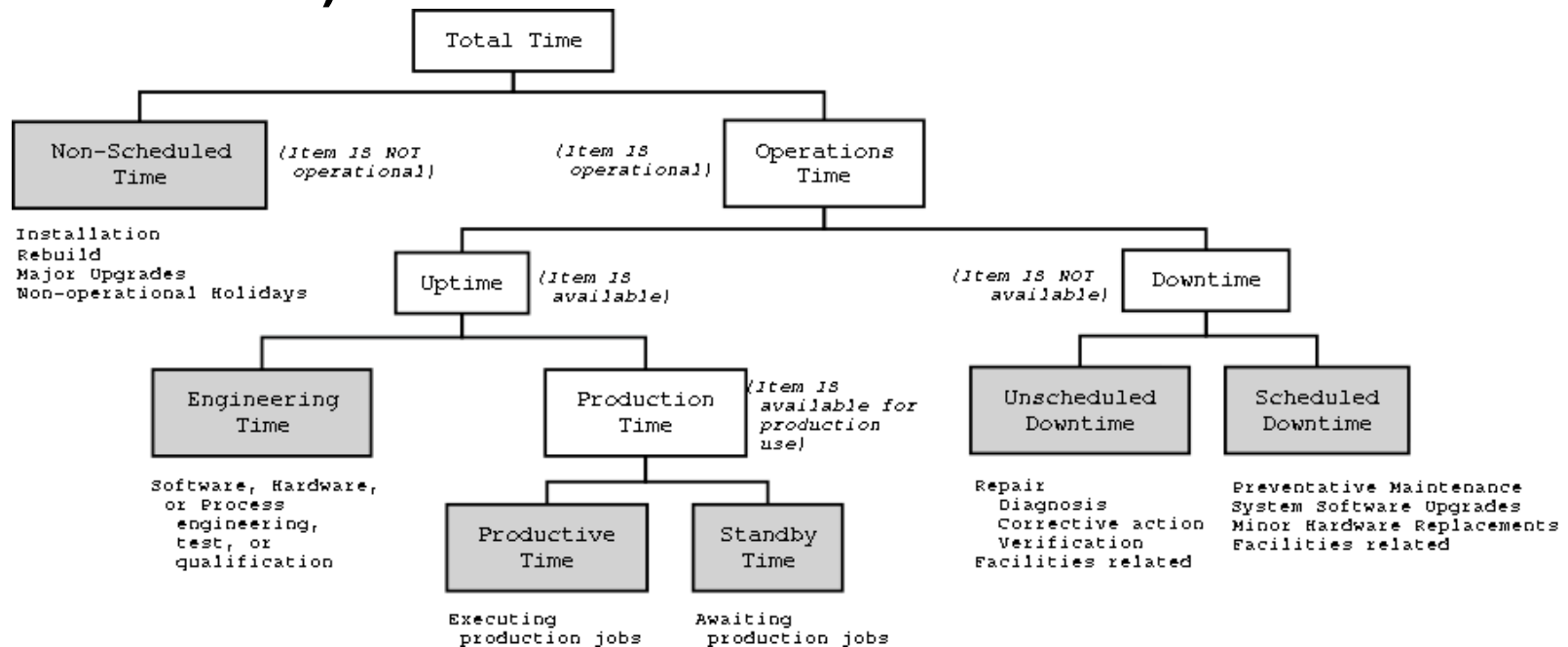
### RedStorm/Classified : Availability





# MTBIService

A measure of reliability describing how long the system stays in a production state (regardless of why an interruption in production service has occurred).



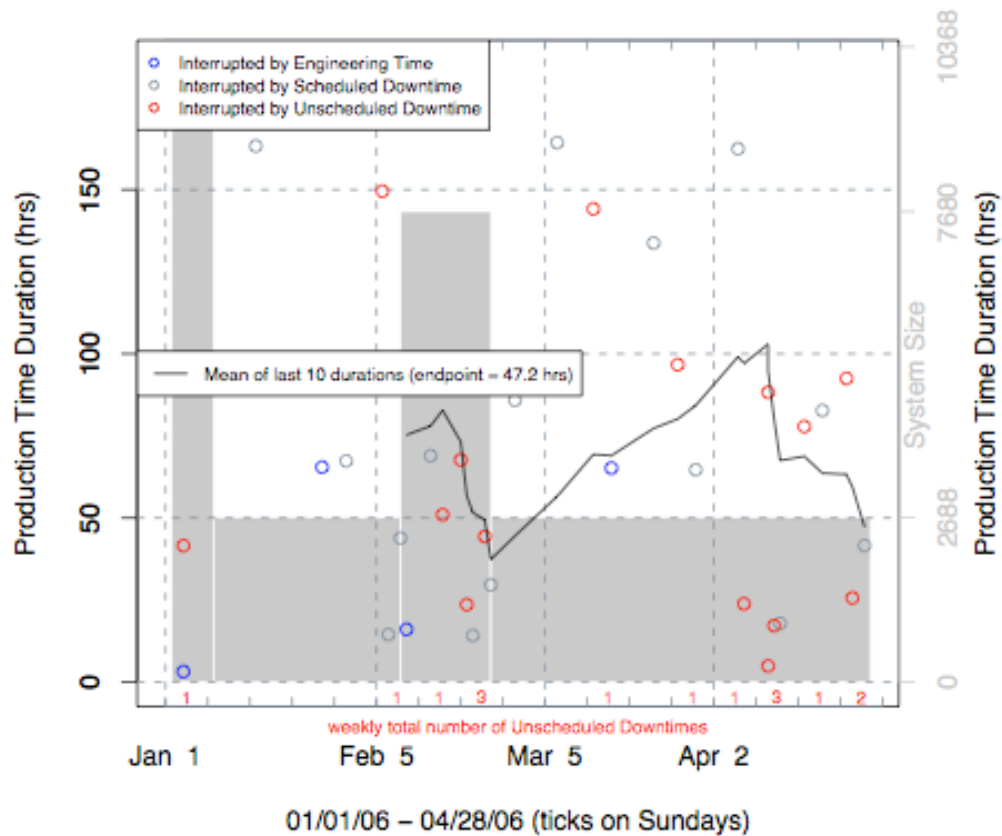
$$\text{MTBIService (hours)} = \sum \frac{\text{Production Time}}{\text{Service Interrupts}}$$



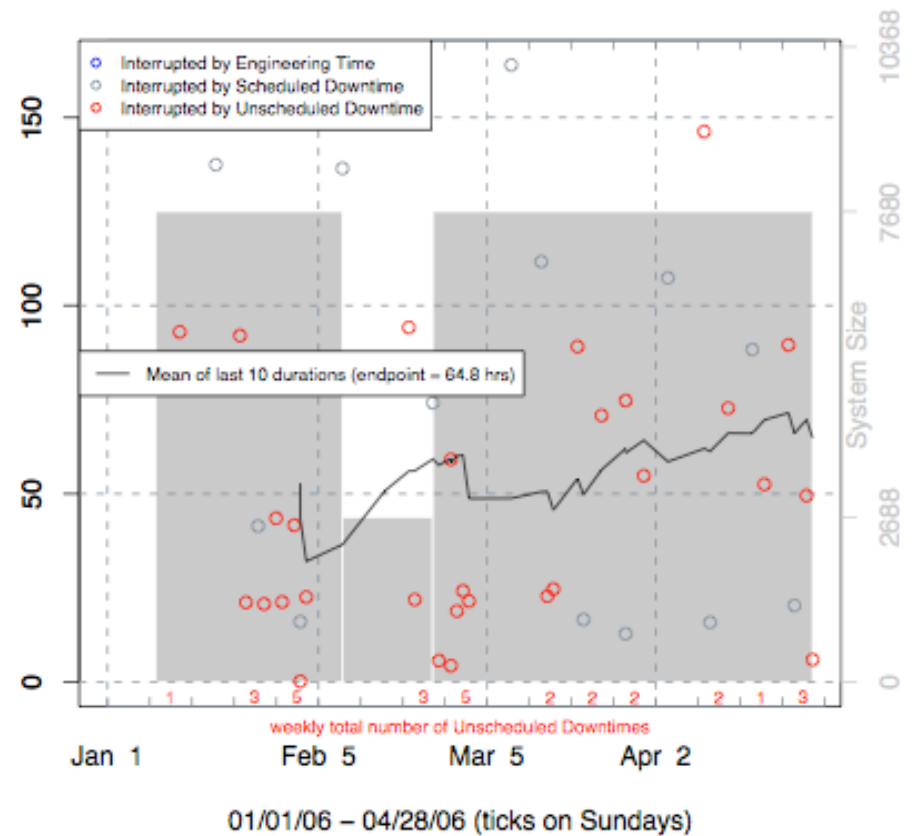


# MTBIService

### RedStorm/Unclassified : MTBIService



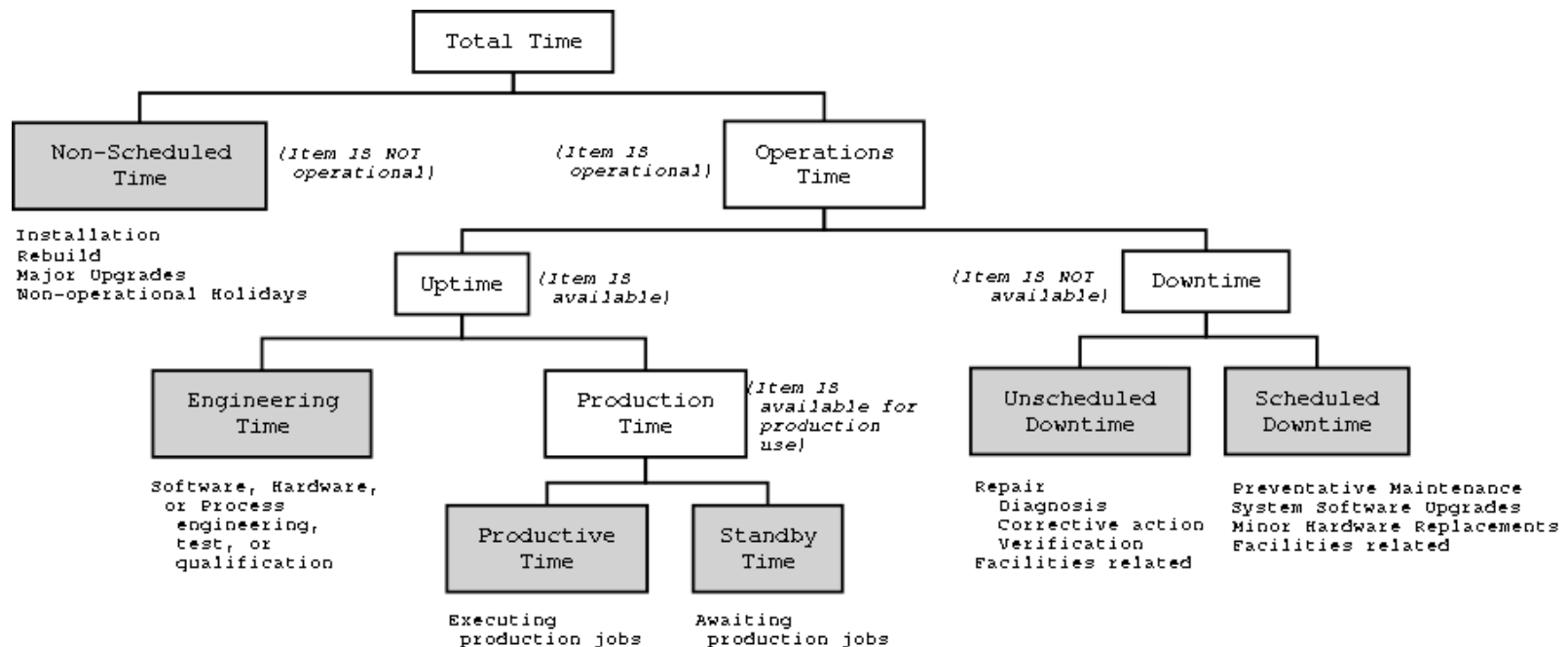
### RedStorm/Classified : MTBIService



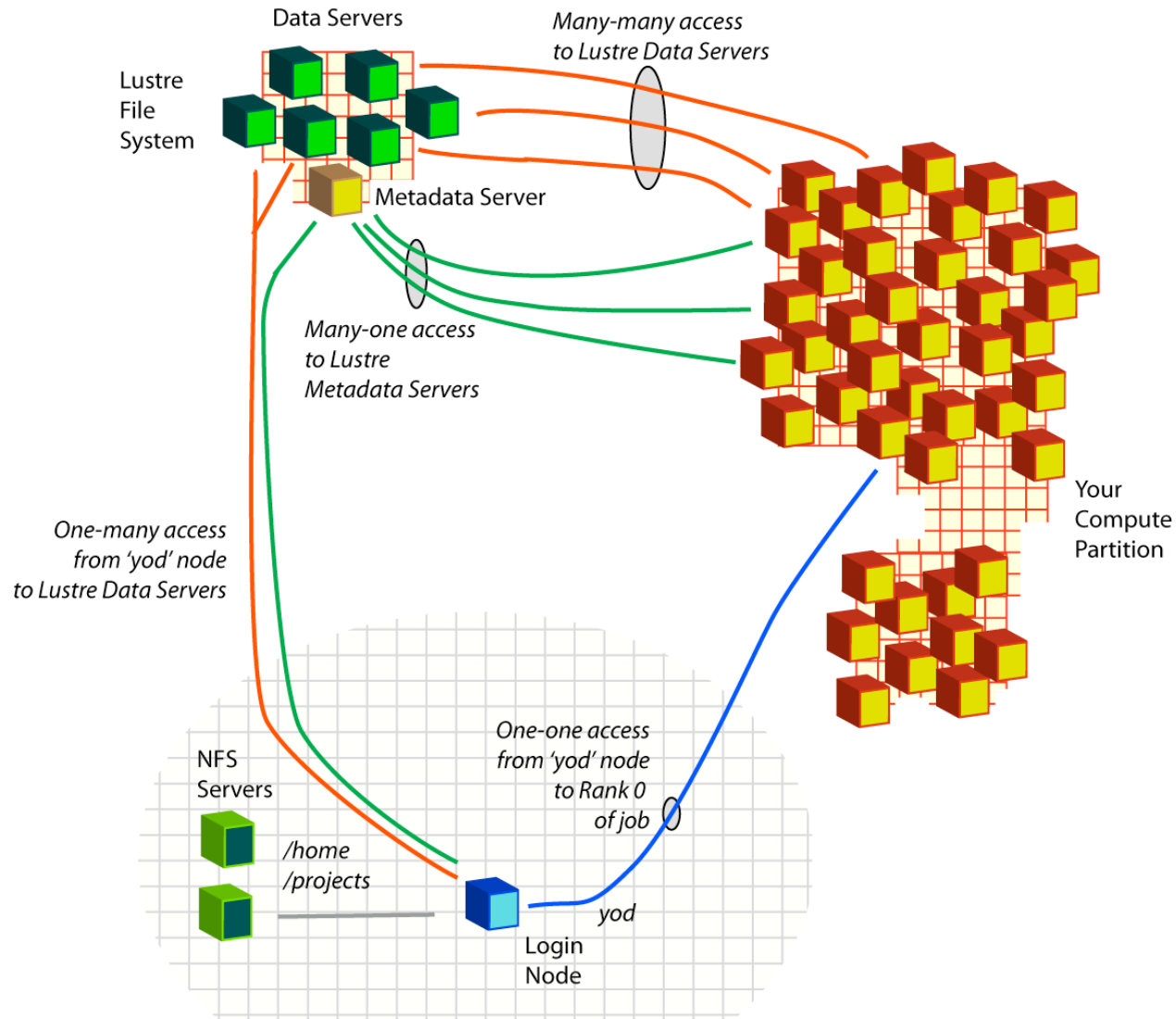


# MTBIjob

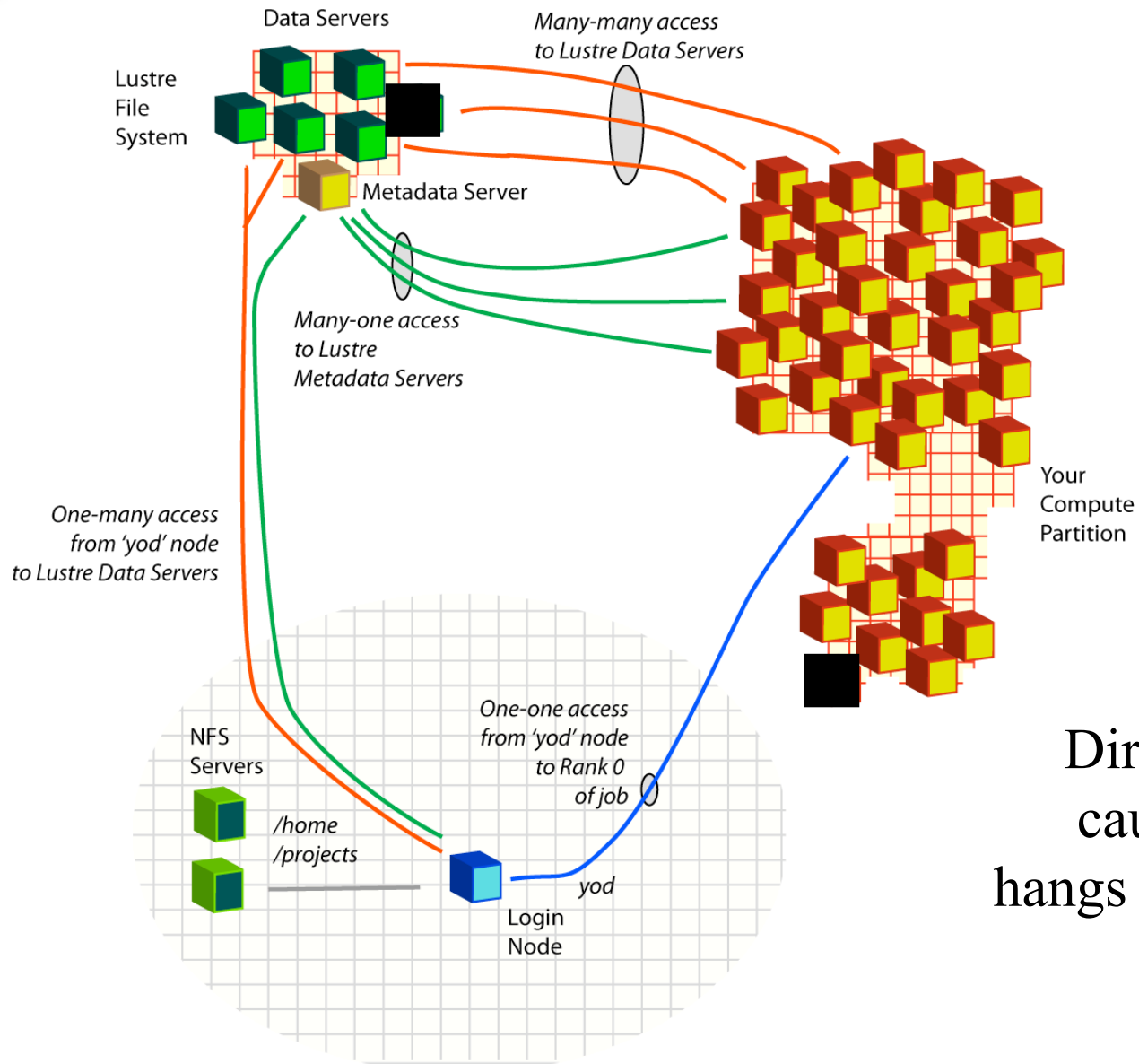
A measure of reliability describing how long the system is in a production state before any job is interrupted (by the system).



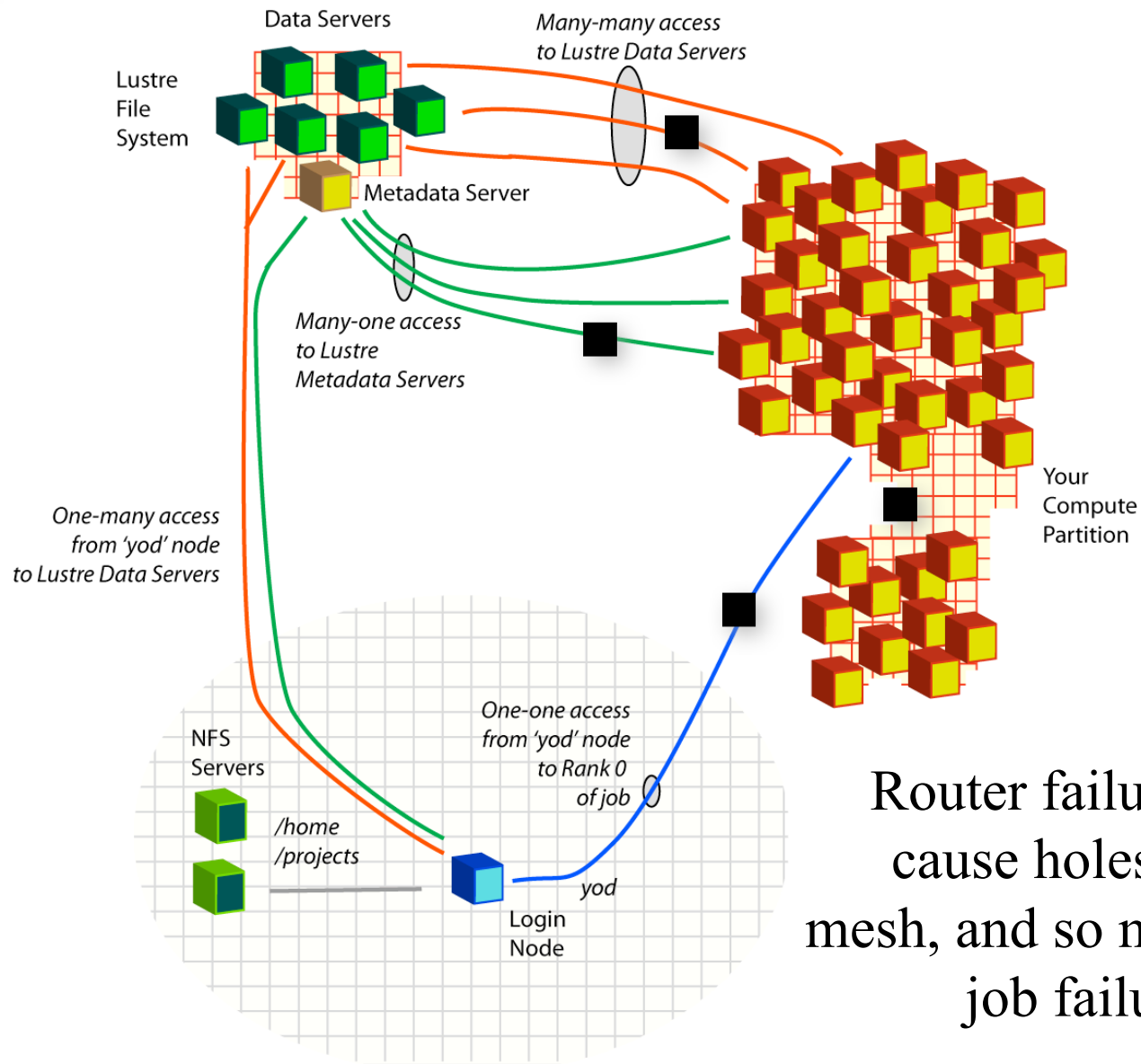
$$\text{MTBIjob (hours)} = \sum \frac{\text{Production Time}}{\text{Job Interrupts}}$$



How do you detect a job failure in an XT3?



Direct hits  
cause job  
hangs or failures

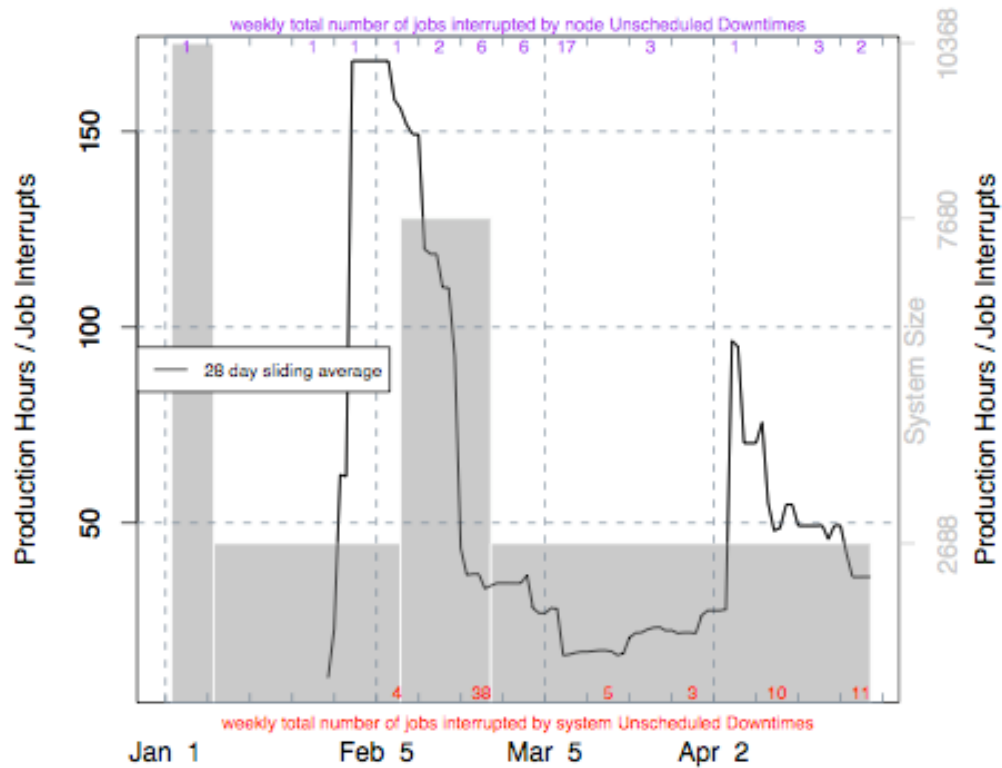


Router failures also cause holes in the mesh, and so may induce job failures



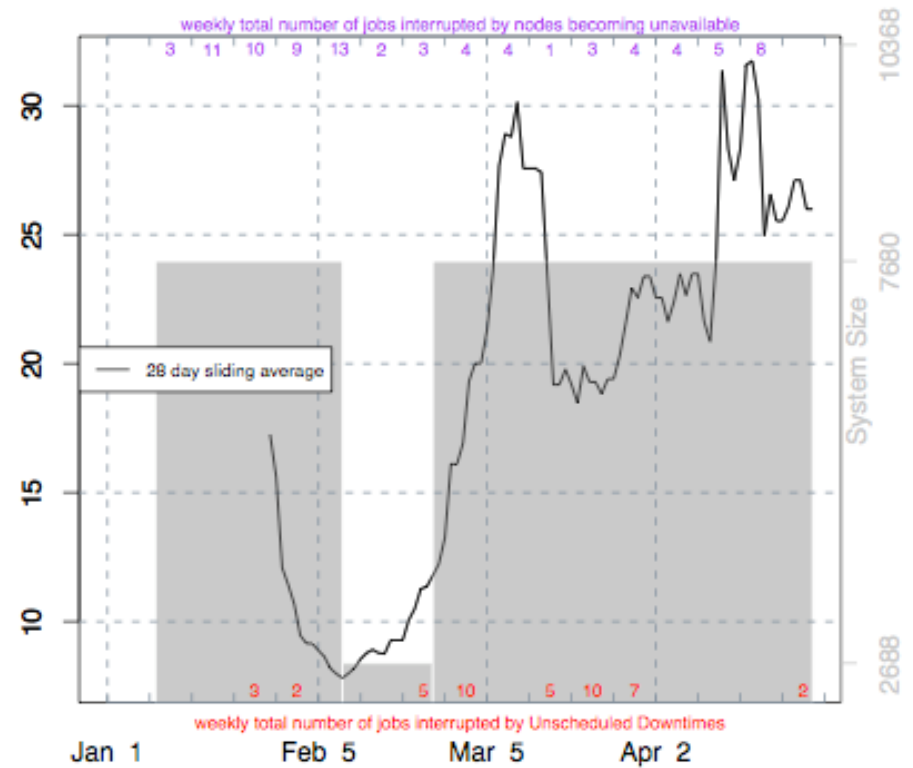
# MTBIjob

### RedStorm/Unclassified : MTBIjob



01/01/06 - 04/28/06 (ticks on Sundays)

### RedStorm/Classified : MTBIjob



01/01/06 - 04/28/06 (ticks on Sundays)



## Lessons Learned

---

- **Clear RAS definitions (states, failures, interrupts, ...) need to be established (the HPC community should do this)!**
- **Means to easily capture data regarding the above is needed (e.g. vendors should supply this)!**
- **Need for the above is increasing with system size and complexity!**
- **Distilling RAS information into clear plots is challenging.**
  - **Buy an XT3 today! ;)**



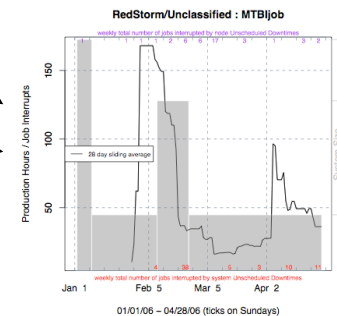
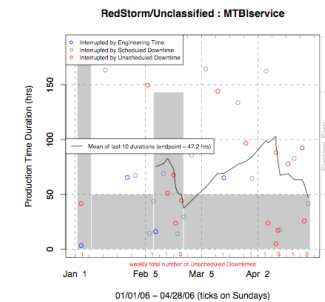
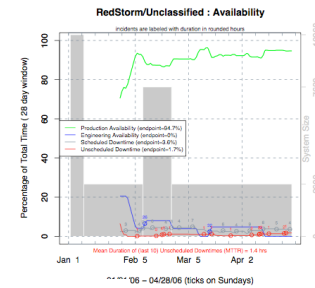
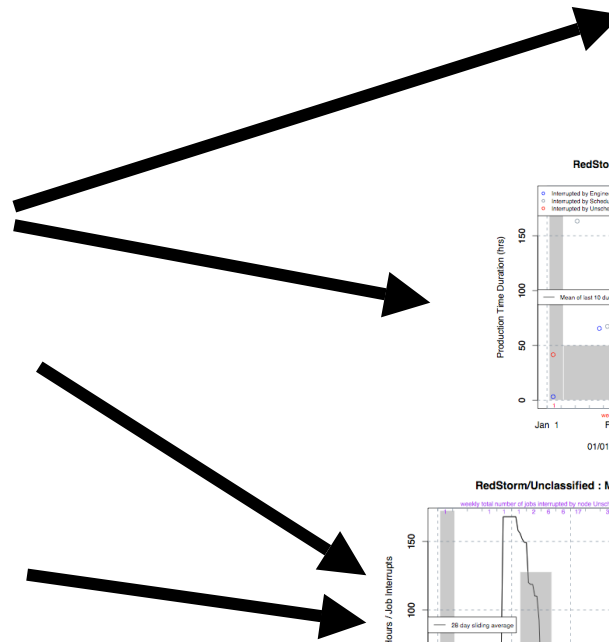
# RASM Toolkit

**Input**  
(ASCII files)

**Output**  
pdf|ps|png|...

**“state\_transitions.tab”**  
time state size

**“job\_interrupts.tab”**  
time interrupt\_type



Written in R, includes documentation and example scripts.





## Nodehour-Scaled Metrics (TBD...)

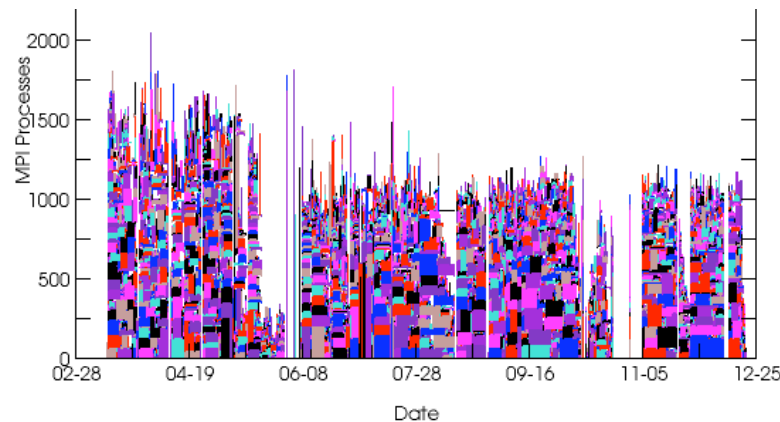
---

**Utilize workload information in order to measure work per interrupt rather than simply time per interrupt.**

$$\text{MNBI}_{\text{service}} (\text{hours}) = \sum \frac{\text{Productive Nodehours}}{\text{Service Interrupts}}$$

$$\text{MNBI}_{\text{job}} (\text{hours}) = \sum \frac{\text{Productive Nodehours}}{\text{Job Interrupts}}$$

...





## For More Info...

---

See <http://www.cs.sandia.gov/~jrstear/ras/>



# Workload information is vital!

---

