# ATLAS Analysis Data Distribution and Panda PD2P

Alden Stradling, Kaushik De, Tadashi Maeno, Torre Wenaus, Paul Nilsson

12 August 2011
DPF, Brown University

# Once upon a time...

- There was the concept of an ideal Grid. These were indeed the best of times.

  - Bandwidth was infinite, authentication reliable, and execution was swift!

  - Data moved swiftly to wherever the battle was thickest, keeping efficiency high.
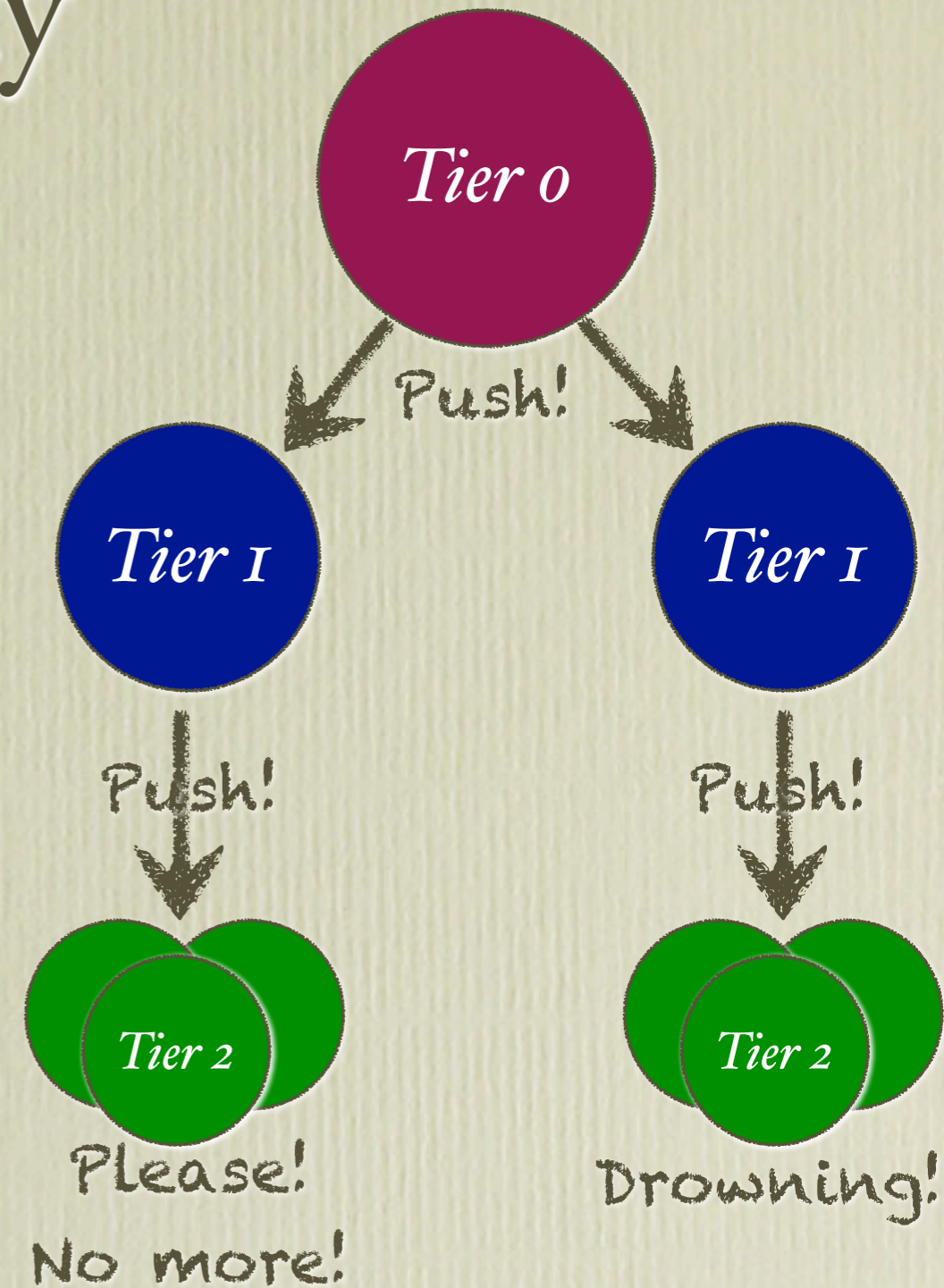
- Then came implementation.

# It so happened

- That the data started to get... big.

- This was, of course, anticipated. For production jobs, data movement time was (correctly) deemed relatively short, and data could easily be sent to the jobs.

- Trying the same thing with user analysis didn't work.

  - Users (ideally) need *instant* start, and quick results. Lots of time wasted in slow debug cycles, otherwise.

  - Moving a terabyte or more is not "instant" at all.

# Monarchy

- Decreeing by policy that $n$ copies of certain data types will be replicated based on the site's tier and etc...

  - Again, great for centrally-run jobs!

- *Creates terrible bottlenecks in user analysis.*



Tier 0

Push!

Tier 1    Tier 1

Push!    Push!

Tier 2    Tier 2

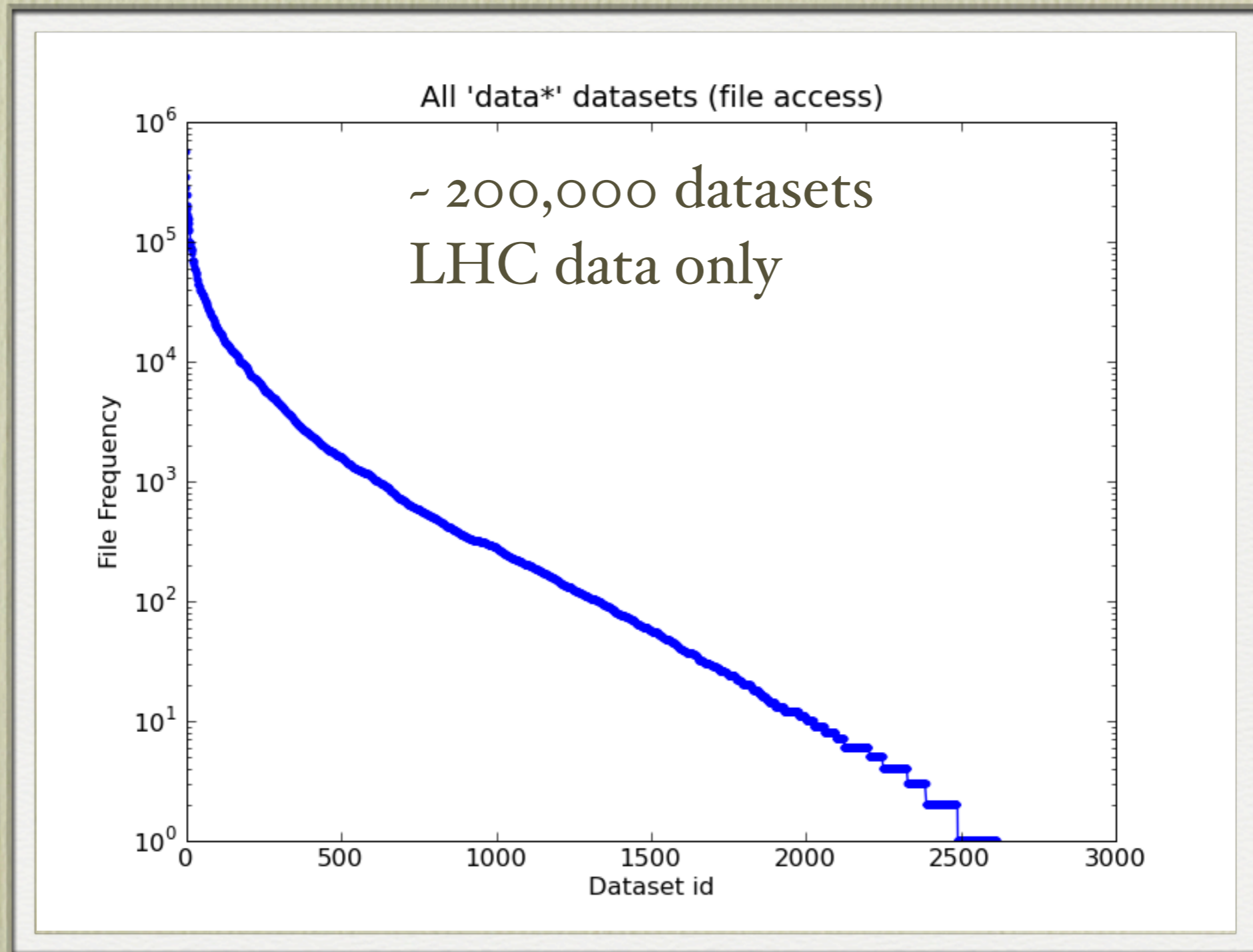Please! No more!    Drowning!

UT ARLINGTON

# Growing Pains

- So we broker the job to the data

  - Unless the sites that have the data are swamped! What then?

- Make the data as widespread as possible

  - But most of the data are unused (shown on next slide)

  - Sites clog quickly, and user jobs are still delayed

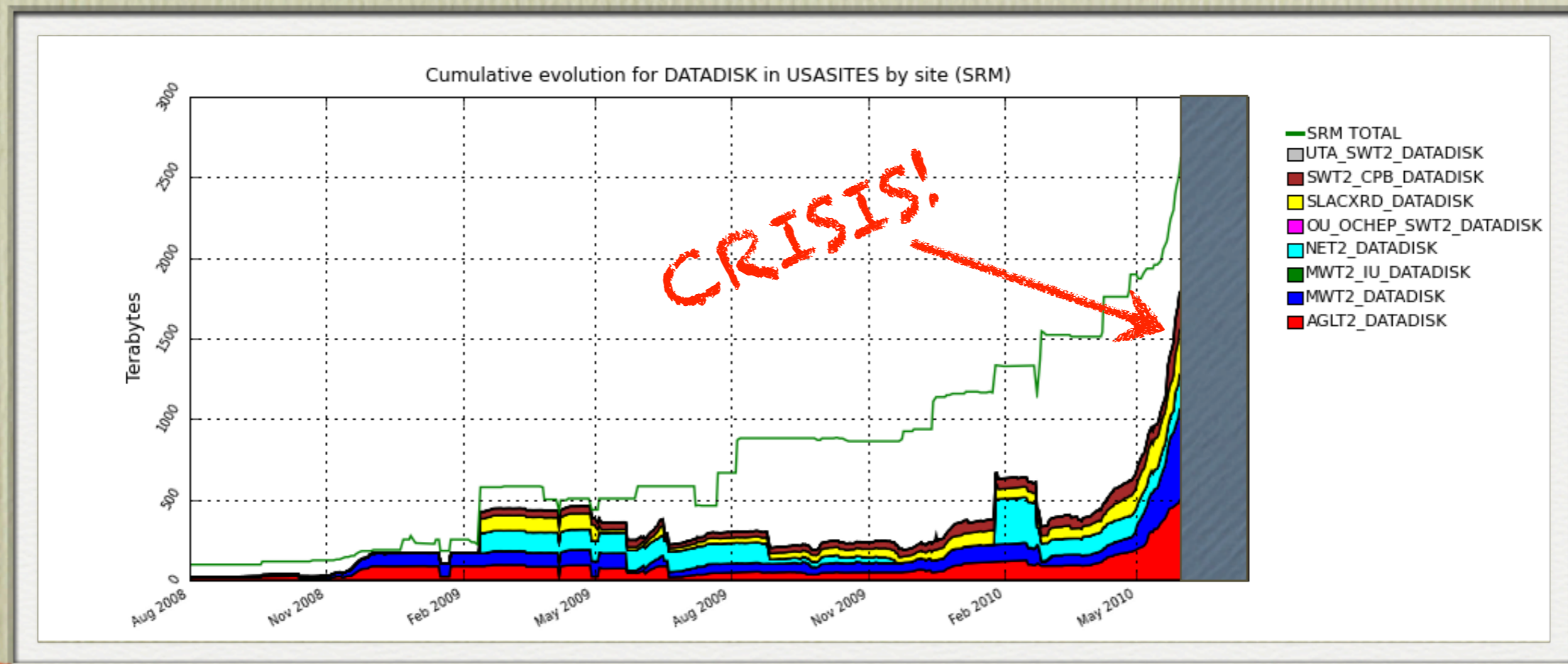  - And this is just the very beginning of LHC data

# Small Fraction Used



All 'data*' datasets (file access)

~ 200,000 datasets
LHC data only

# Rising Tide

- Exponential rise between February and June 2010
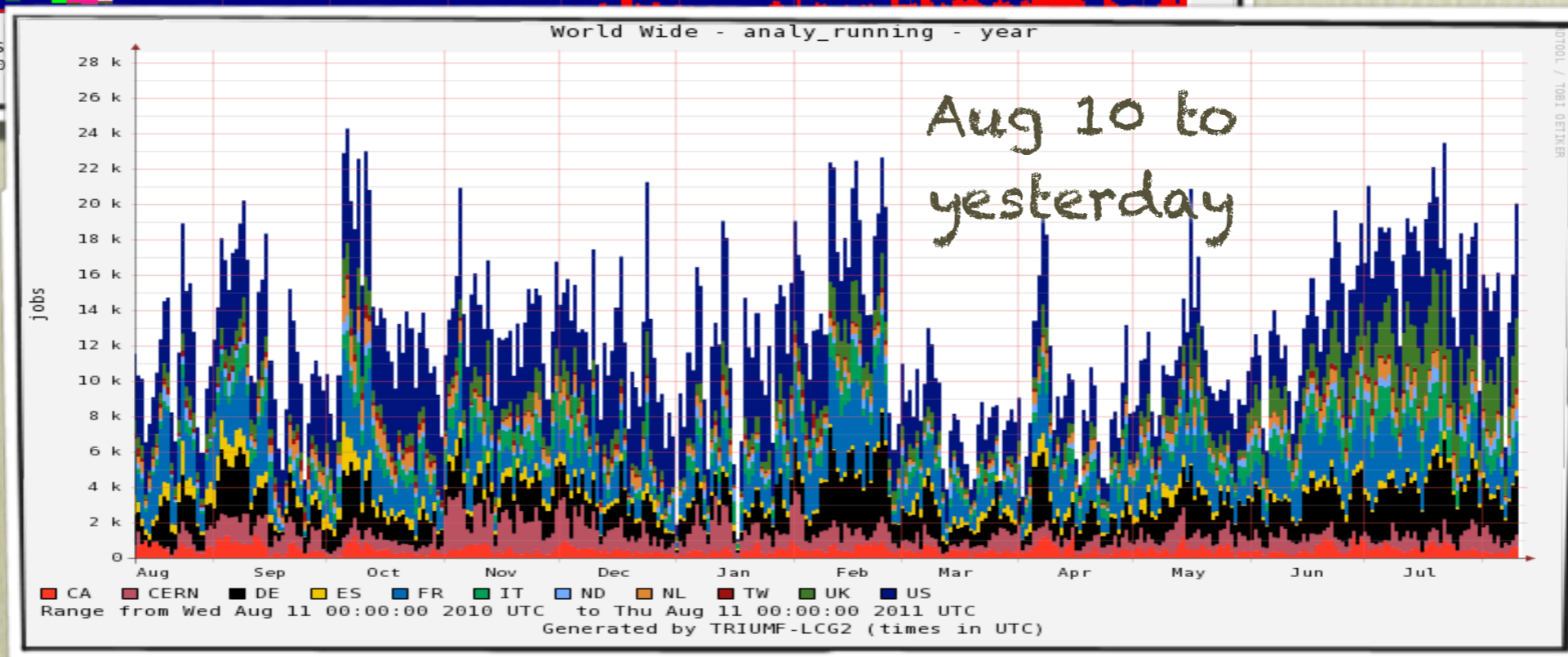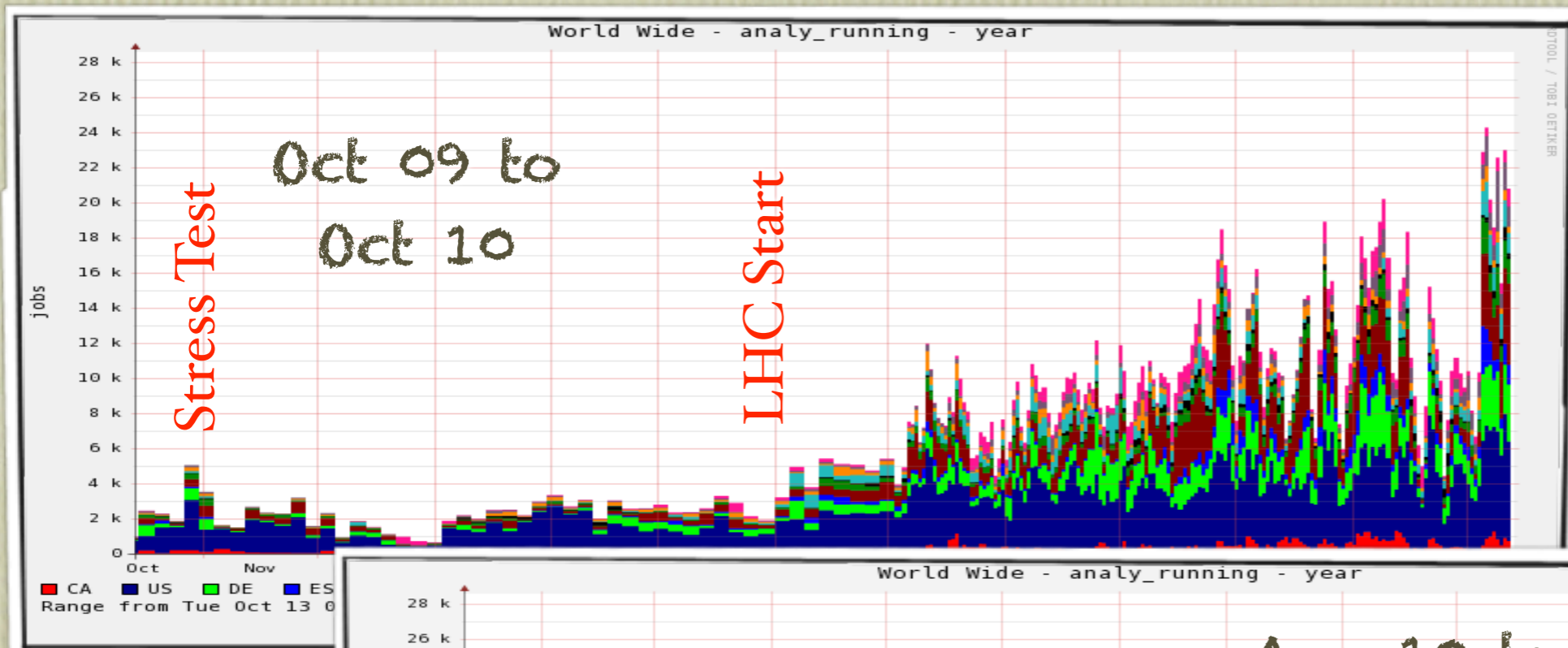- Unsustainable! (just the *beginning* of data collection)

# PD²P to the Rescue!

- **Panda Dynamic Data Placement (PD2P)**

- June 2010, put initial algorithm into play in the US cloud as a test as the *reverse* of the Monarch Model

  - *When a dataset is used* (even once), it is subscribed to a Tier2 site

  - Greater demand, more subscriptions

  - Unpopularity determines cleanup

  - Exclude heavier data in favor of user analysis types ($A_{nalysis}O_{bject}D_{ata}$, ntuple, even $E_{vent}S_{ummary}D_{ata}$)

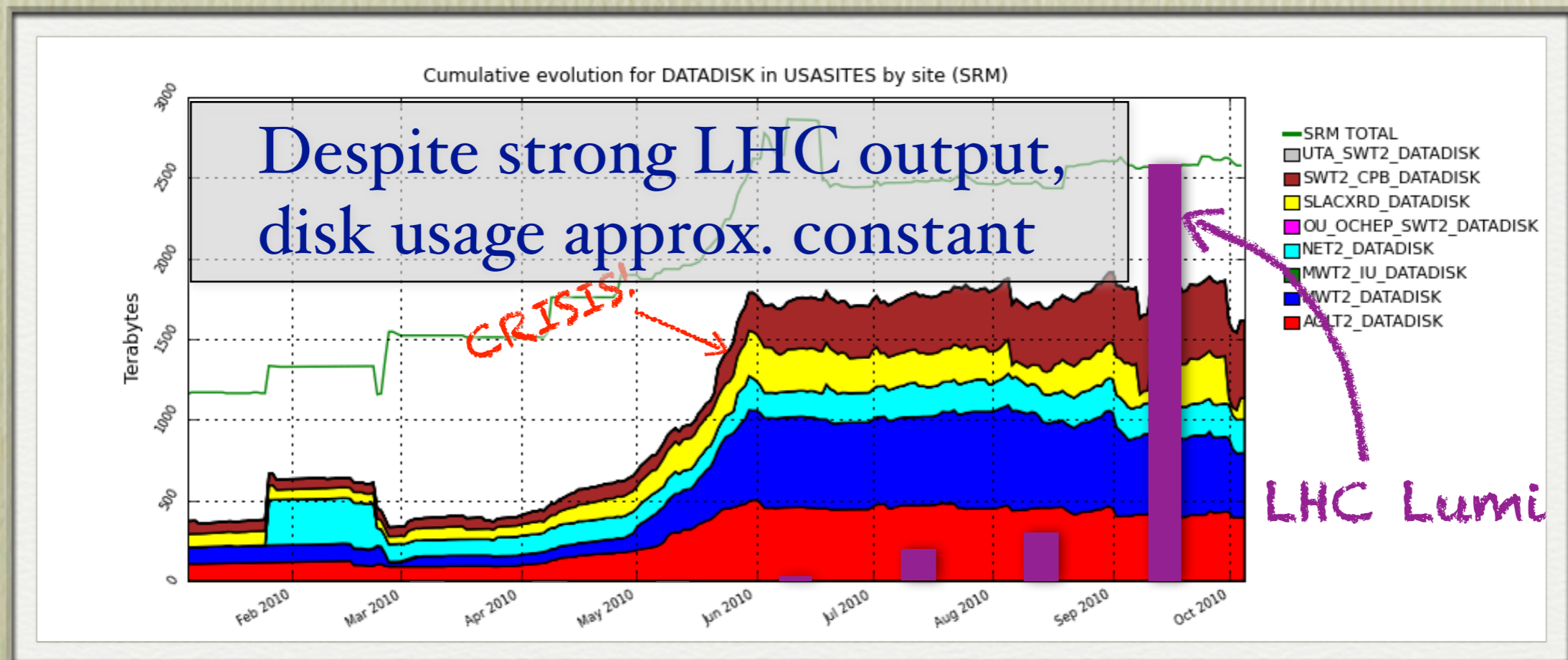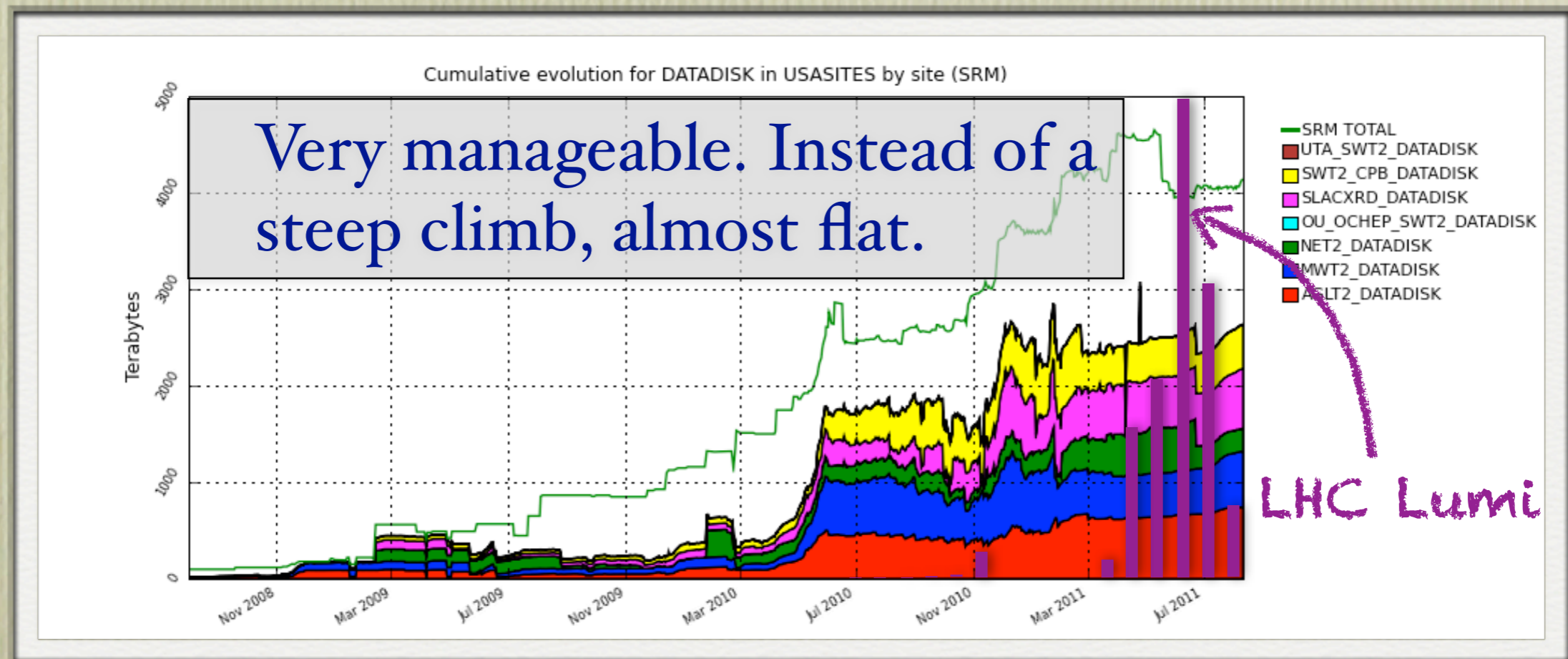  - User-created datasets are left out of the algorithm.

# Distributed Analysis

# Initial Results (Oct 2010)

- Very encouraging. Plateau in data growth.

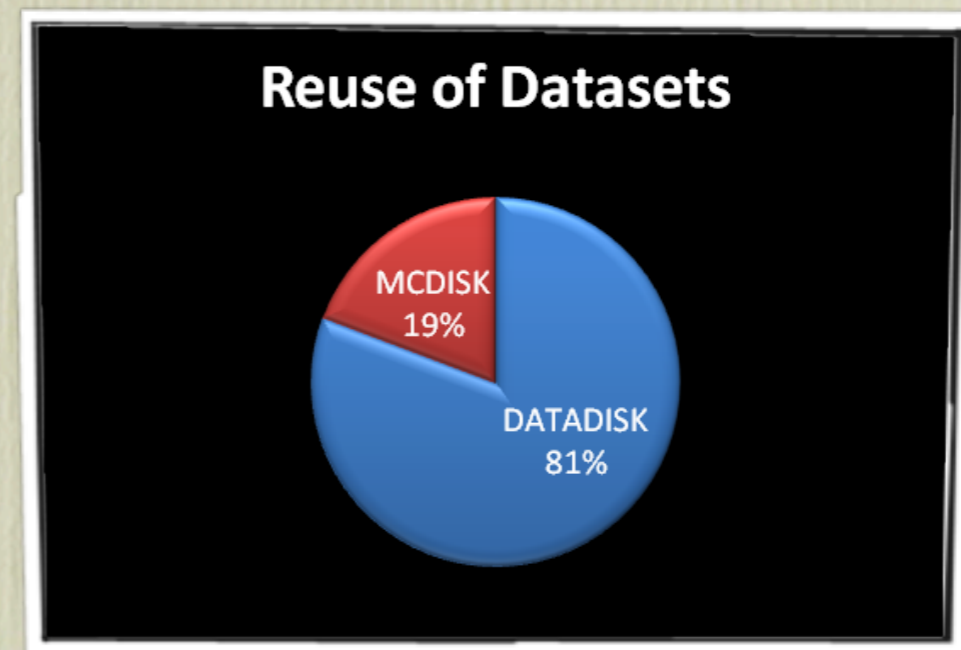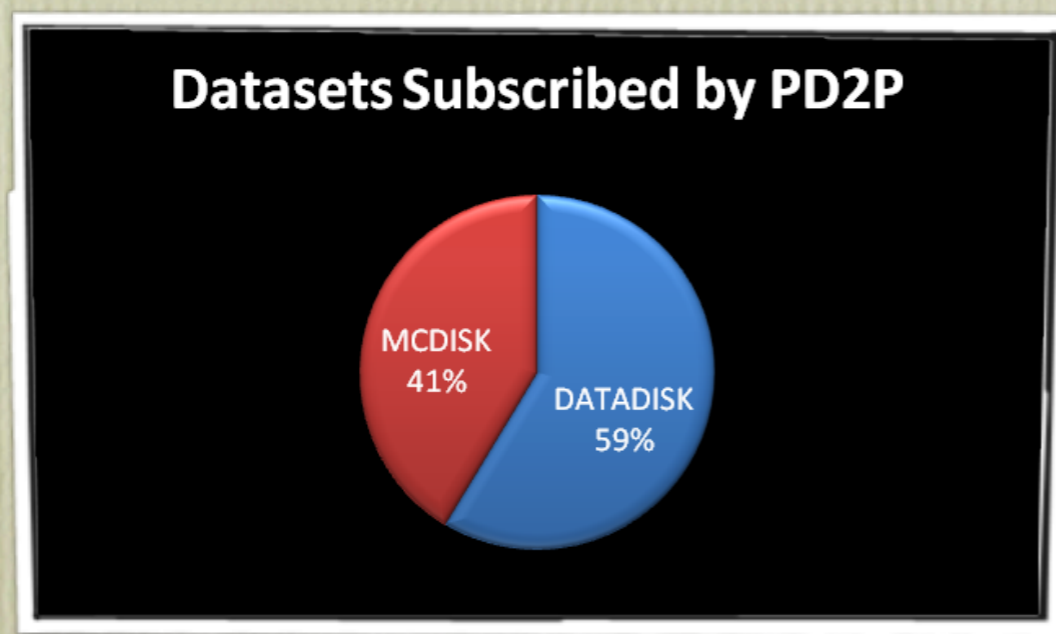- In spite of *rapid* lumi growth and constant user analysis

Cumulative evolution for DATADISK in USASITES by site (SRM)

Despite strong LHC output, disk usage approx. constant

CRISIS!

LHC Lumi

- SRM TOTAL
- UTA_SWT2_DATADISK
- SWT2_CPB_DATADISK
- SLACXRD_DATADISK
- OU_OCHEP_SWT2_DATADISK
- NET2_DATADISK
- MWT2_IU_DATADISK
- MWT2_DATADISK
- AGLT2_DATADISK

# Up Until Today...

- The plateau continues, with dips in some places

- User analysis, as seen above, is huge. So is **lumi**.



Very manageable. Instead of a steep climb, almost flat.
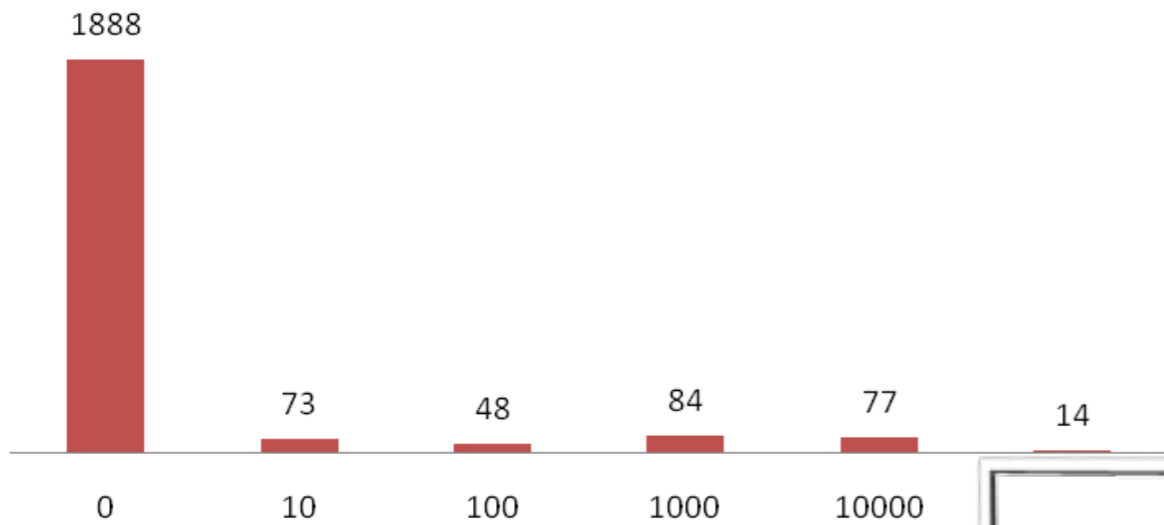
LHC Lumi

# Initial Behavior

- Steady rate of subscriptions (no big spikes)

- Site distribution is fairly even

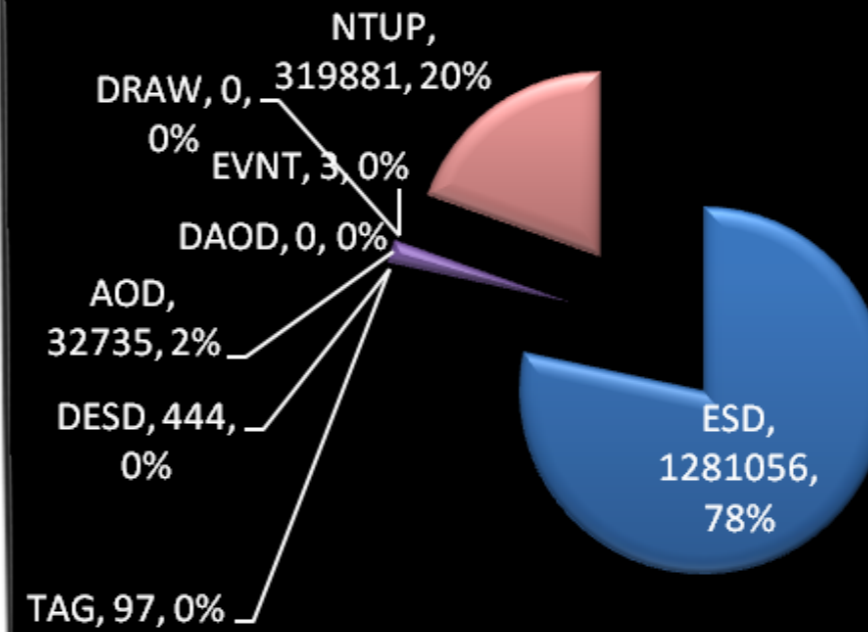- LHC data are (of course) more often subscribed and reused than Monte Carlo
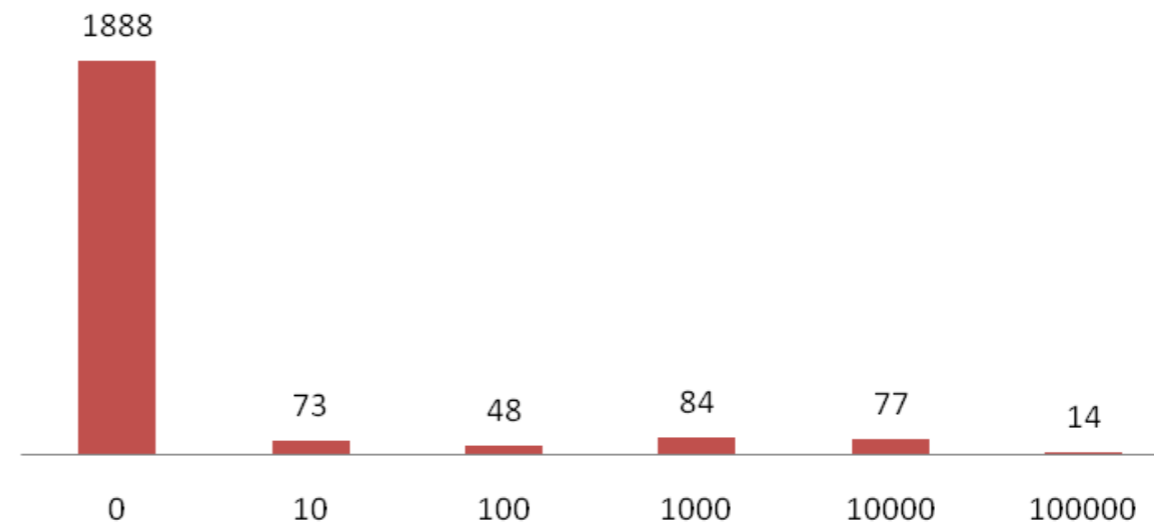
# How Much Reuse?



Reuse of PD2P Datasets - Month 4
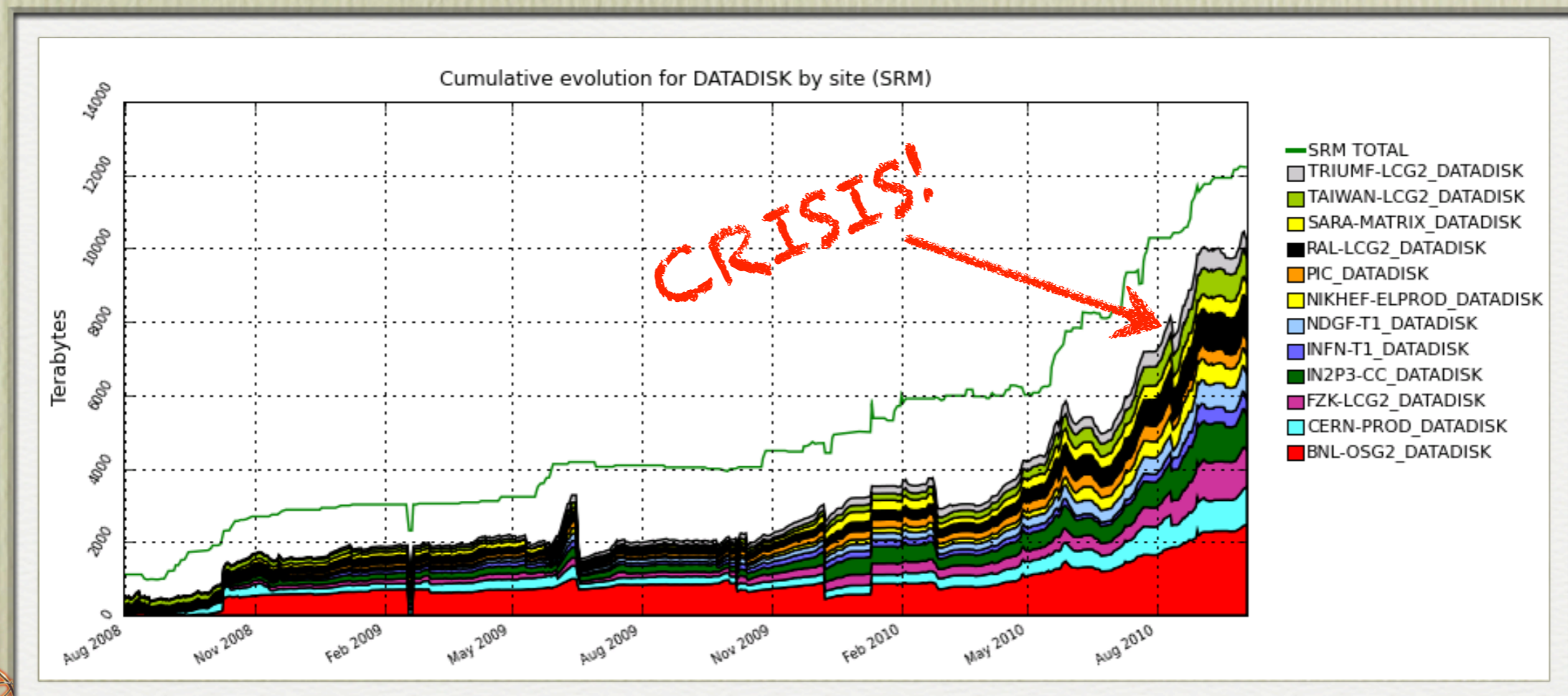(Sep 15 - Oct 14)



Reuse of Datasets by Type



Reuse of PD2P Datasets - Month 4
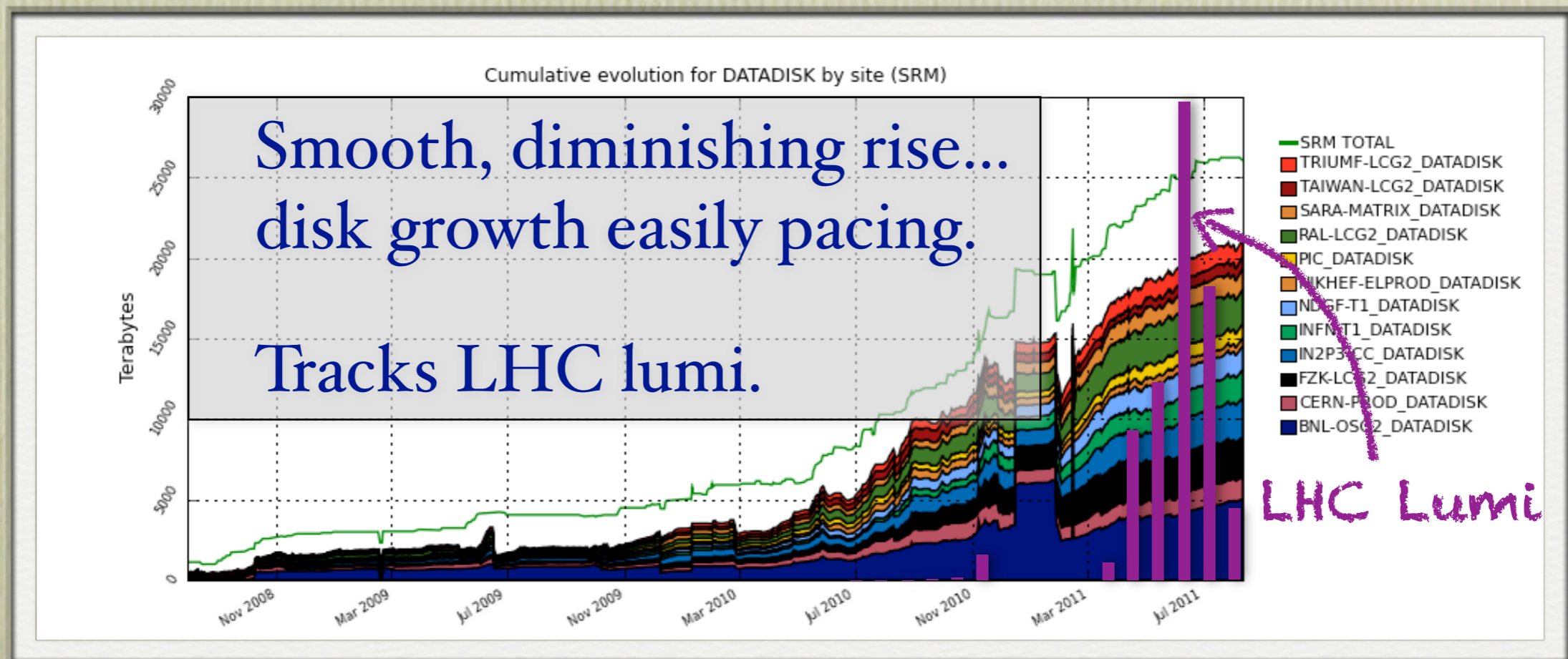(Sep 15 - Oct 14)

# T1 Sites – Another Crisis

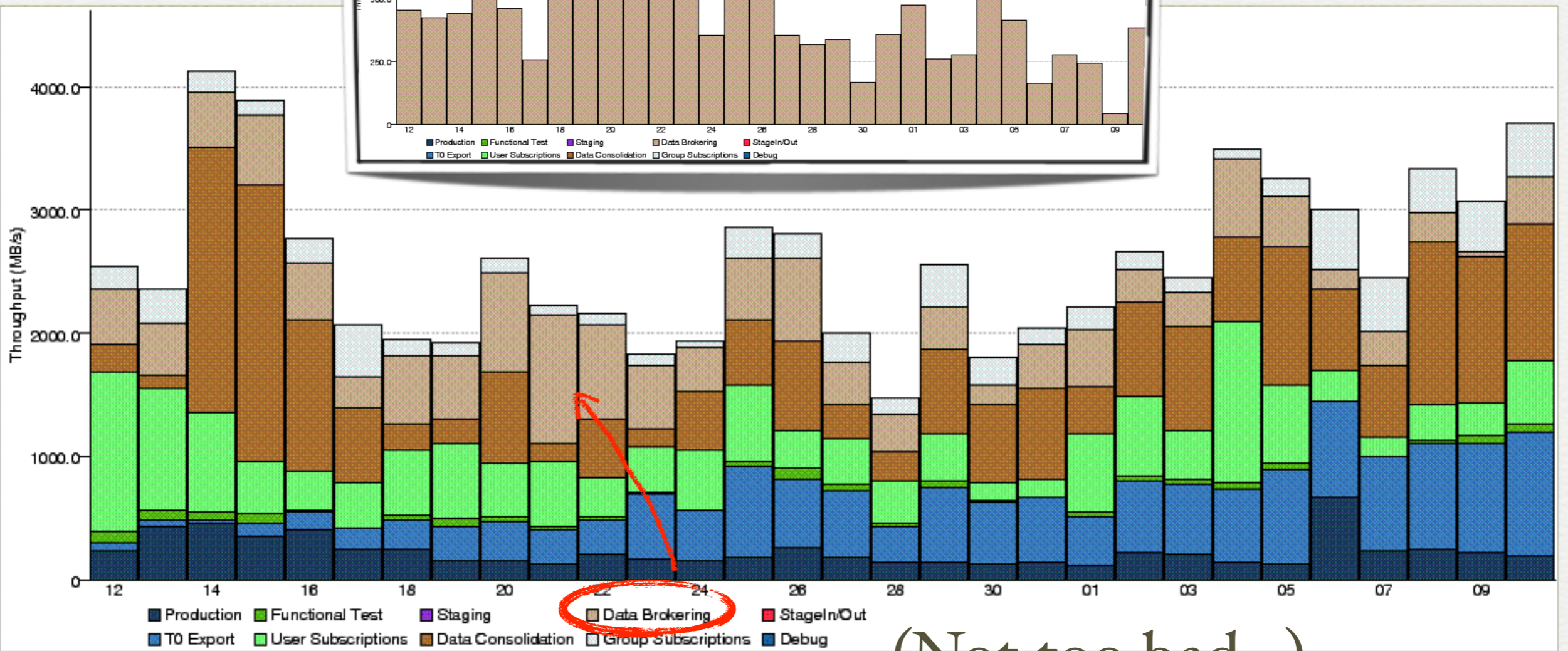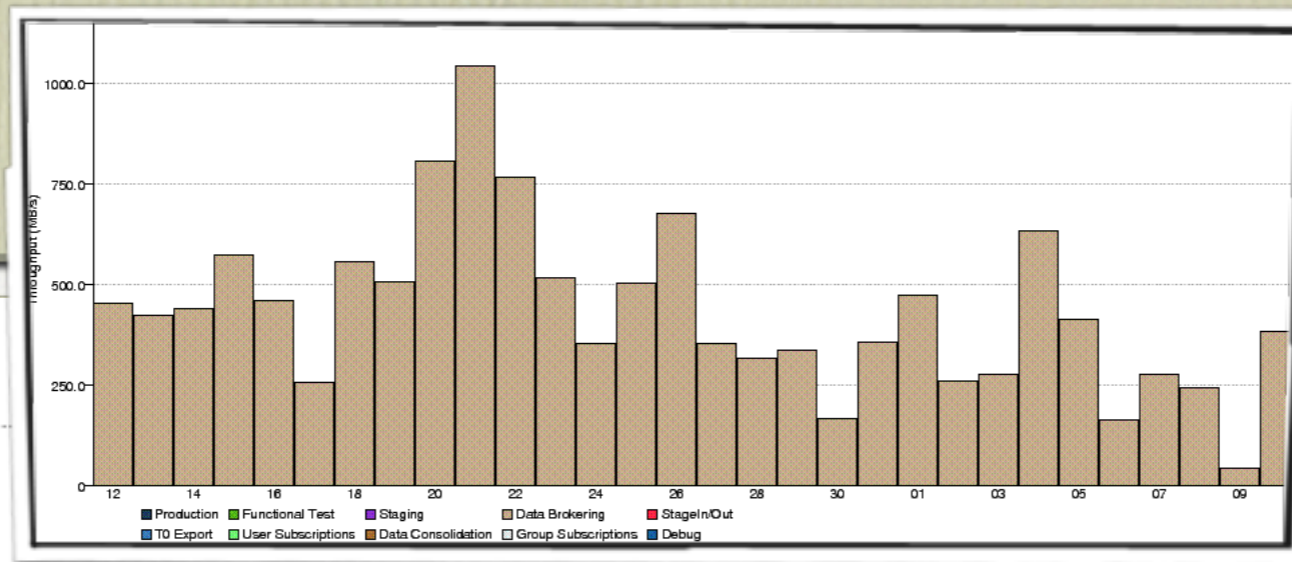- With the Tier2 peril newly overcome, the Tier1 data become... unruly

# PD2P To The Rescue!

- Distributes differently – according to MoU share, and then popularity

  - Logarithmic – new copies every 10th/100th/1,000th/10,000th usage.



Cumulative evolution for DATADISK by site (SRM)

Smooth, diminishing rise...
disk growth easily pacing.

Tracks LHC lumi.
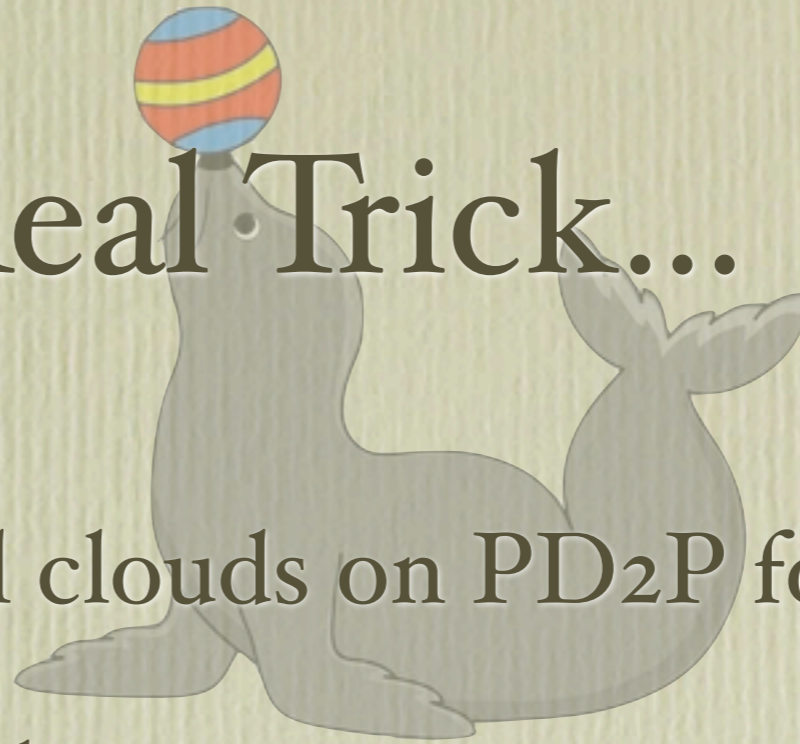
LHC Lumi

# How Much Load?

Past 30 Days



(Not too bad…)

# The Real Trick...

- We've moved to all clouds on PD2P for Tier2

- The transition has been **transparent** to users.

  - No complaints on the help lists

  - No delays noticed (meaning that there are no delays that exceed other slowdowns)

- The transition has been a boon for the site admins and deletion services operators

# And So It Happened...

- That the situation once again became manageable!

    - Storage stayed reasonable

    - Small tweaks to the algorithm improved the situation incrementally even more

    - All the T1 sites were added to PD2P

- Jobs that had languished with excessive delays could be **rebrokered**, for the data would be moved as well (since there was already a delay)
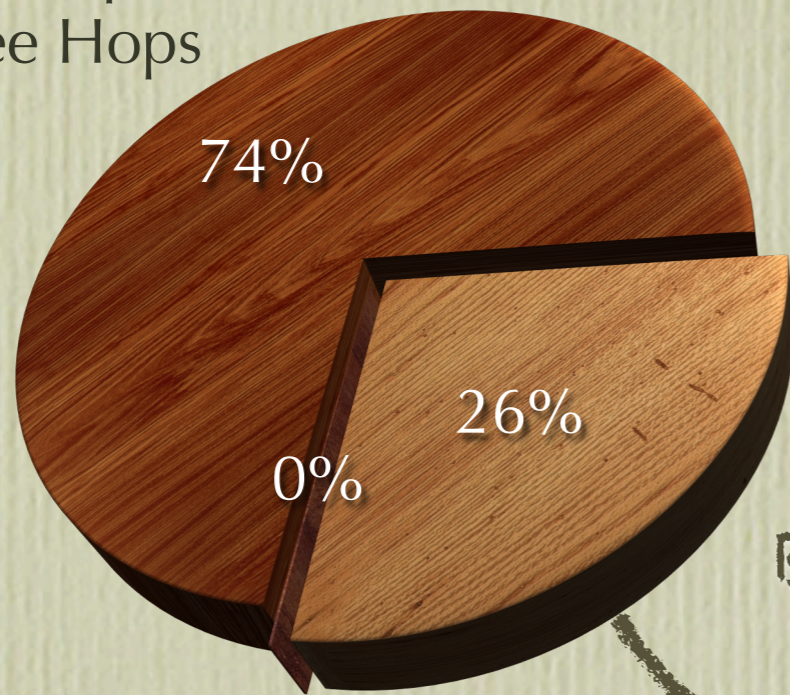
# Rebrokerage

- And so 300,000 jobs that waited to start for more than 72 hours were sent to a new site

- It was decreed that 72 hours was too long, and the delay was decreased to a single day

- And the rebrokerage rate leapt fivefold! (5x)

- But the dreaded "bouncing job" never came about, nor was PD2P flooded with data.

# Job Hopping
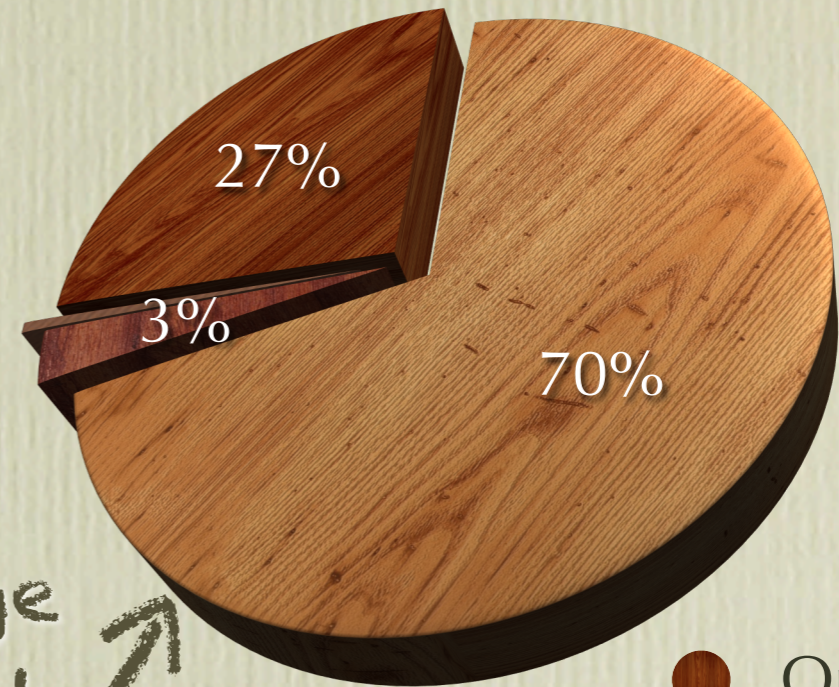


**New**

One Hop
Two Hops
Three Hops

74%
26%
0%

**Old**

Big change in 2-hops!

27%
3%
70%

One Hop
Two Hops
Three Hops
Four Hops

(But no runaway behavior, and two 24h hops are still better than one 72h hop!)

US ATLAS

UT ARLINGTON ™

# And They Lived Happily...

- Of course not. Much to do to make data management "smart"

  - Widen the gates and let more copies be made to T1 sites. Do some pre-placement.

  - Make more T2 copies of popular datasets

- Attempt to copy only the parts of the dataset that will be used in the job

- Possibly find patterns in use, and make predictive copies of datasets likely to be popular?

# Ever After…

- The Ideal Grid, where Data speed to where they're needed, allowing transfers of individual files (or even events!)

- Or where data are read directly over Federated Xrootd, from any site to any site.

- Final abolition of all "cloud" boundaries, and have access to all data from any site, from the least to the greatest.
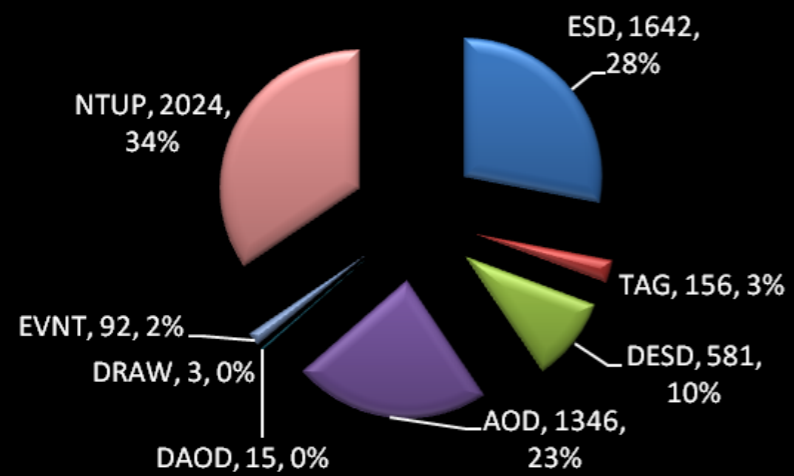
*The End*

# Backup Slides

# What Kinds of Data?



# of Datasets Subscribed by Type

ESD, 1642, 28%
NTUP, 2024, 34%
TAG, 156, 3%
EVNT, 92, 2%
DRAW, 3, 0%
DESD, 581, 10%
DAOD, 15, 0%
AOD, 1346, 23%

Past

Last Week



D2P replication by data type in ALL Cloud in last 7 days (2011-08-11.15.00.02 CEST)

NTUP 1188
ESD 0
DESD 272
AOD 434
DAOD 235