# Wide Area Network Data Access Requirements in Analysis

Doug Benjamin

Duke University

# Setting the stage

- Since the beginning of the year, how much analysis was done?

- Why Jan 1 to May 1?
  - has analysis ramp up ahead of major Winter conferences (Moriond series)
  - Post conference activity perhaps indicative of steady state activity well into LS1 period
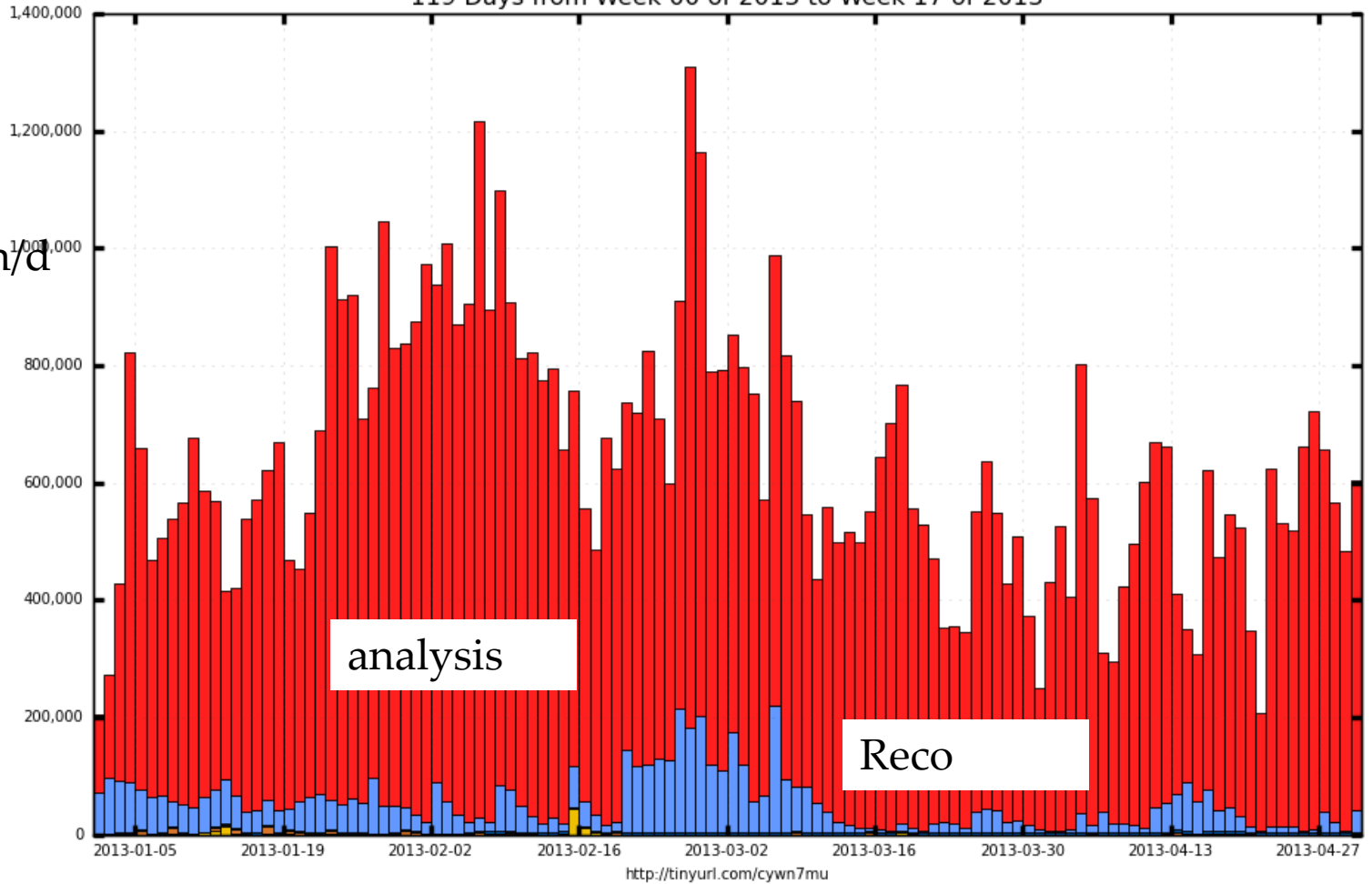
- What data do analyzers use?

# Completed Jobs (1-Jan to 1 May 2013)



1 million/day

Job Type - Group Analysis, Group production, User analysis

# User Analysis jobs 1-Jan – 1-May 2013



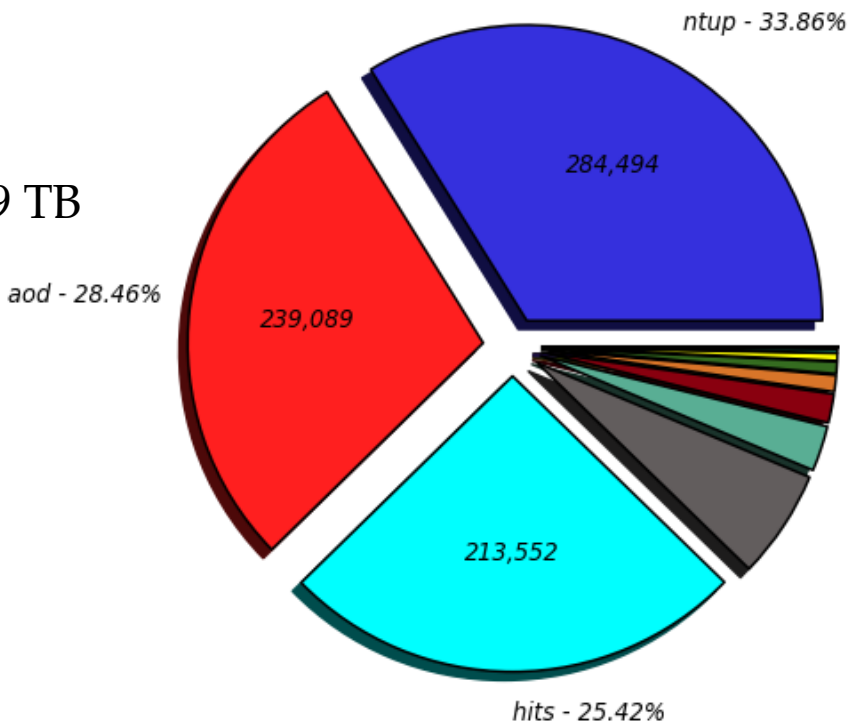Job Type  - Group Analysis, Group production, User analysis

# Data popularity
## (used datasets 2012-09-22 to 2013-03-21)



AOD
239,089 TB

NTUP
284,494 TB

AOD's used
to make
NTUP's

# Data unpopularity
## (Unused datasets)



Relative share of number of TB in unused datasets (2012-09-22 - 2013-03-21) (Sum: 44,306)

esd - 38.29% — 16,960
raw - 20.75% — 9,192
hits - 18.49% — 8,189
ntup - 10.20% — 4,516

AOD (1,919 TB)

NTUP (4,516 TB)

- esd - 38.29% (16,960)
- aod - 4.33% (1,919)
- log - 1.19% (528.00)
- daod - 0.26% (116.00)
- d2aod - 0.03% (14.00)
- d2esdm - 0.01% (5.00)
- pt20 - 0.00% (0.00)
- o1 - 0.00% (0.00)
- v15 - 0.00% (0.00)
- v15000000 - 0.00% (0.00)

- raw - 20.75% (9,192)
- rdo - 1.98% (876.00)
- cbnt - 0.68% (302.00)
- evnt - 0.12% (52.00)
- tag - 0.02% (8.00)
- d2esd - 0.00% (2.00)
- atlasproduction - 0.00% (0.00)
- mbts - 0.00% (0.00)
- 15 - 0.00% (0.00)
- filt1iet - 0.00% (0.00)

- hits - 18.49% (8,189)
- desd - 1.52% (675.00)
- hist - 0.39% (175.00)
- dpd - 0.09% (40.00)
- gen - 0.02% (7.00)
- rpcwbeam - 0.00% (2.00)
- pool - 0.00% (0.00)
- v140500 - 0.00% (0.00)
- l1caloem - 0.00% (0.00)
- v2 - 0.00% (0.00)

- ntup - 10.20% (4,516)
- desdm - 1.29% (569.00)
- d2aodm - 0.29% (129.00)
- draw - 0.05% (23.00)
- other - 0.01% (5.00)
- v1451 - 0.00% (0.00)
- txt - 0.00% (0.00)
- tgcwbeam - 0.00% (0.00)
- mcut - 0.00% (0.00)
- ... plus 86 more

# Interim summary

- 1-Jan to 1-May 2013 – peak of over 1M analysis jobs per day
- Prior to Moriond ~800 K jobs/day
- Post Moriond ~ 400 – 600 K jobs/day
- Jobs are typically short (most < 4 hours)
- In 6 month period End – Sep 2012 to end March 2013, NTUP's most popular data type 284 PB used (~ 67 PB/month)
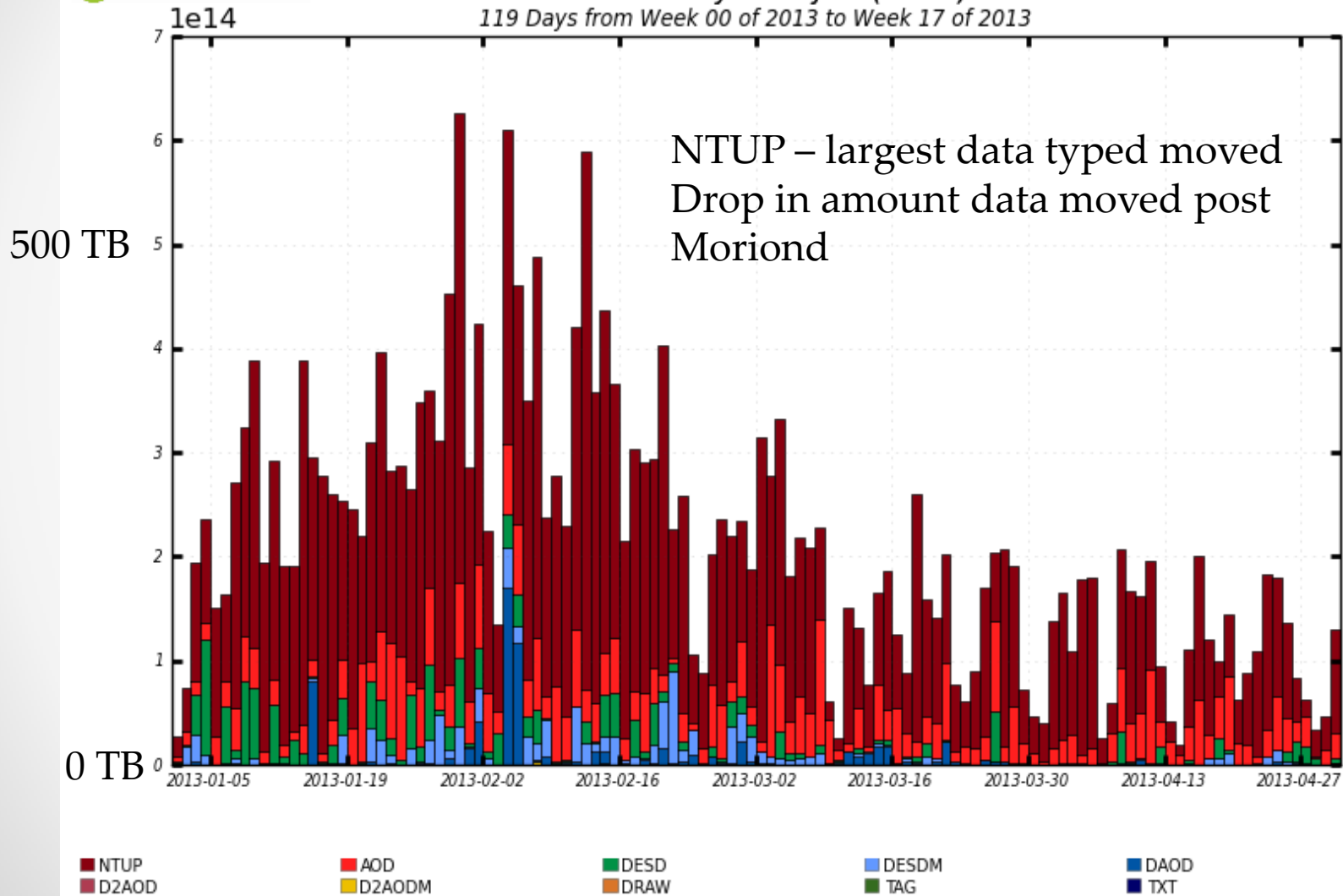- In same period - 4.5 PB NTUP's unused (~ 1.5%)

# Analysis triggered data motion (PD2P)



**Number of Physical Bytes (in TBs)**
119 Days from Week 00 of 2013 to Week 17 of 2013

NTUP – largest data typed moved
Drop in amount data moved post Moriond

Legend: NTUP, AOD, DESD, DESDM, DAOD, D2AOD, D2AODM, DRAW, TAG, TXT

Maximum: 625,444,452,253,292 , Minimum: 18,567,501,102,876 , Average: 213,946,526,963,683 , Current: 130,511,950,942,669

# PD2P by Algorithm

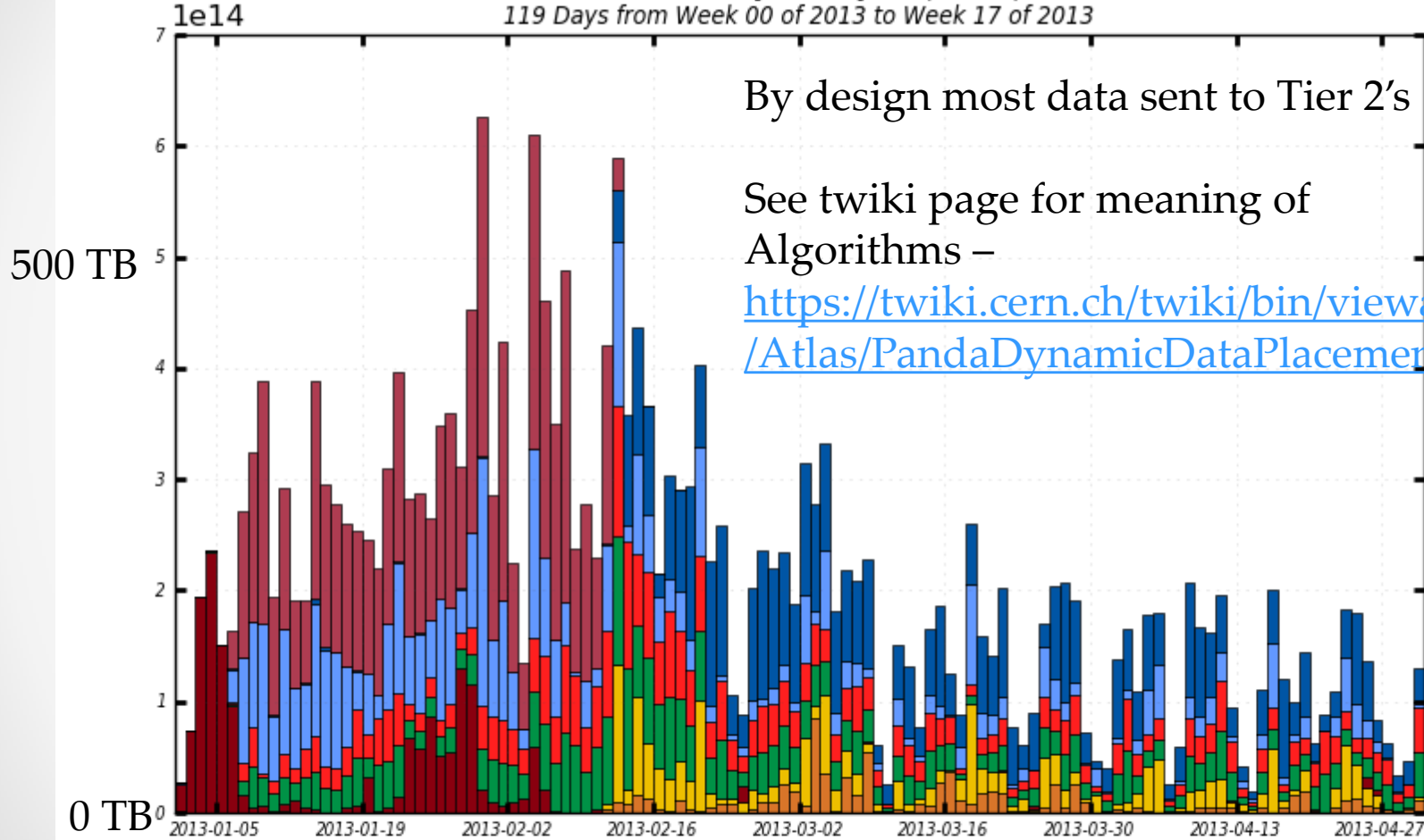By design most data sent to Tier 2's

See twiki page for meaning of Algorithms –
https://twiki.cern.ch/twiki/bin/viewauth/Atlas/PandaDynamicDataPlacement



**Number of Physical Bytes (in TBs)**
119 Days from Week 00 of 2013 to Week 17 of 2013

Legend:
- SELECTEDT2
- SELECTEDT2_NOREP
- SELECTEDT2_T2MOU
- SELECTEDT1
- SELECTEDT2_T1MOU
- unknown
- SELECTEDT2_JOB
- SELECTEDT2_WAIT

Maximum: 625,444,452,253,292 , Minimum: 18,567,501,102,876 , Average: 213,946,526,963,683 , Current: 130,511,950,942,669
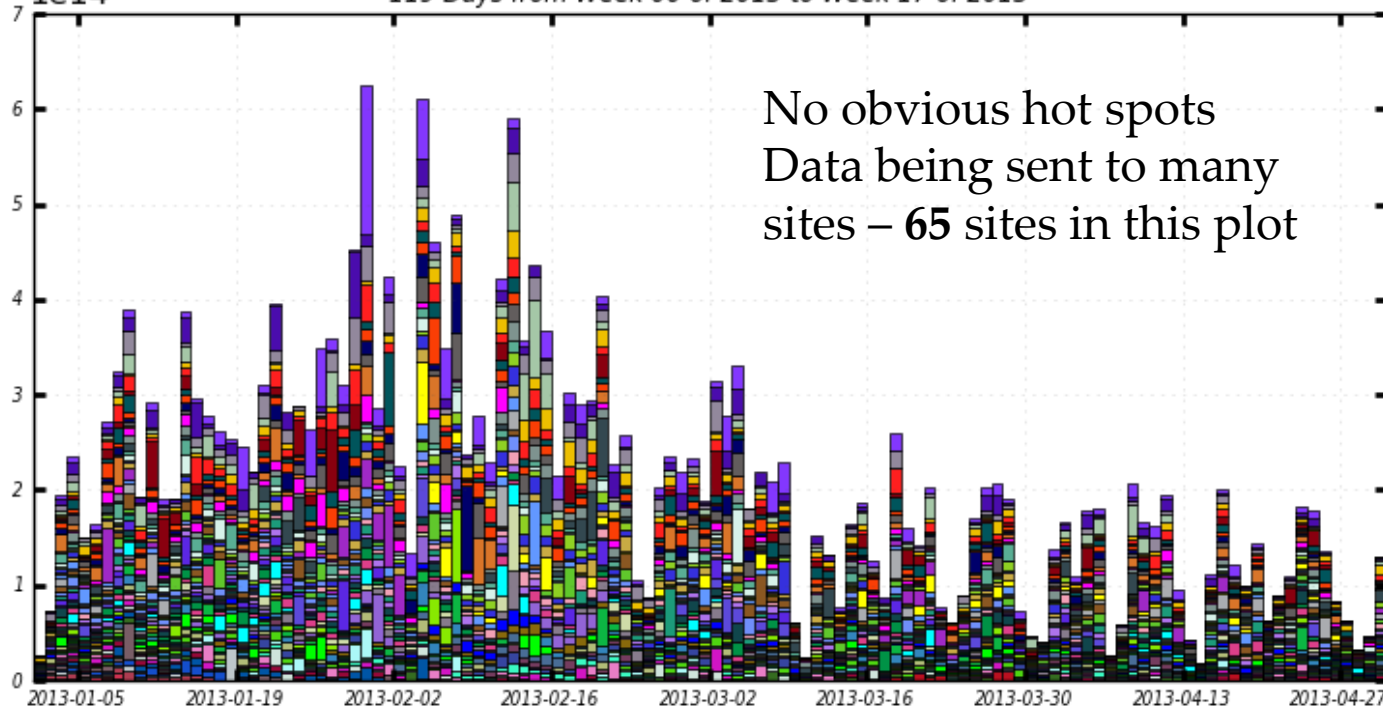
# PD2P sends data everywhere



Number of Physical Bytes (in TBs)
119 Days from Week 00 of 2013 to Week 17 of 2013

No obvious hot spots
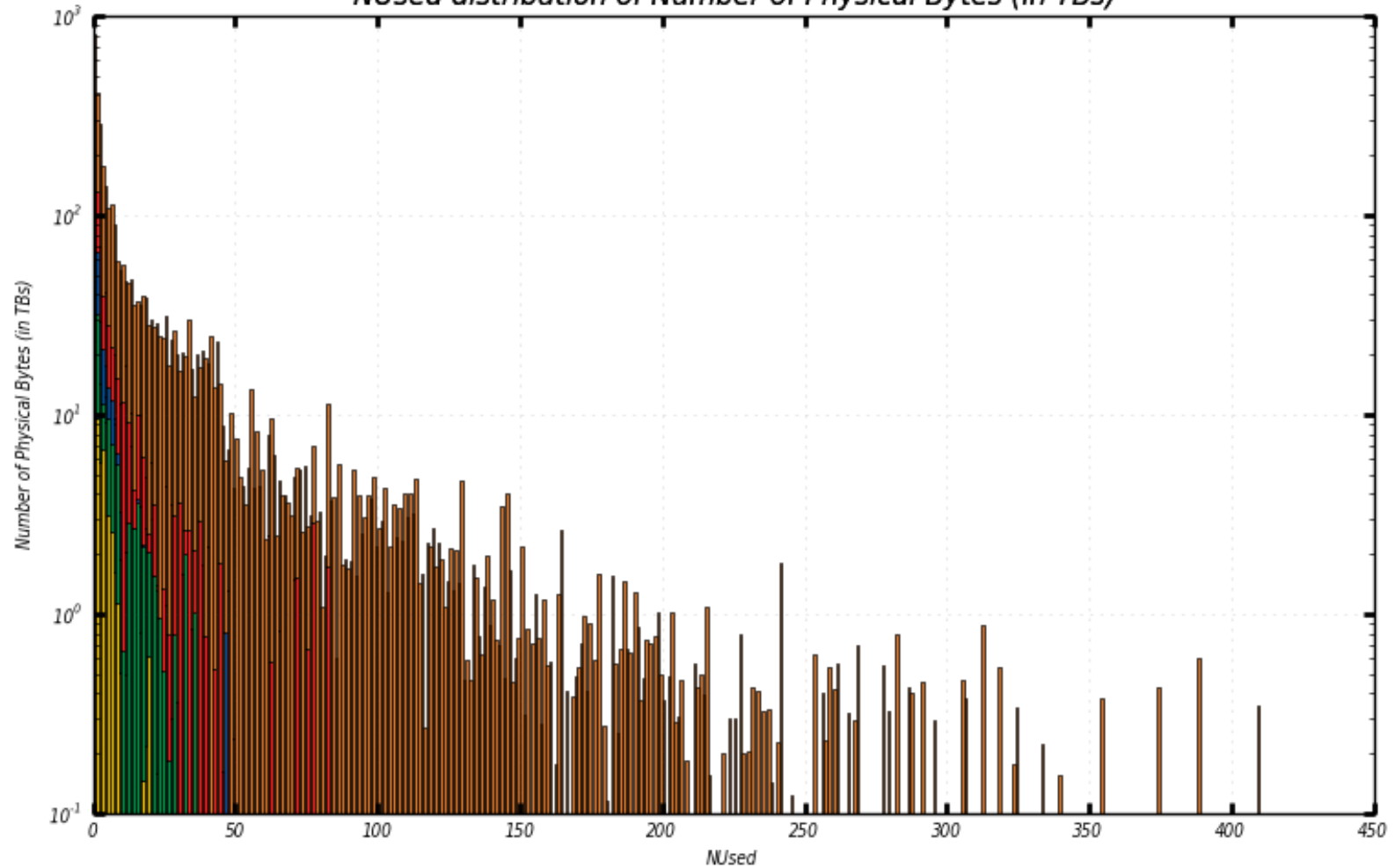Data being sent to many
sites – **65** sites in this plot

# Dataset reuse



NUsed distribution of Number of Physical Bytes (in TBs)

Maximum: 814.96 , Minimum: 0.00 , Average: 13.92

# Fraction of Data volume reused

**Fraction Data used vs PD2P NUsed**

— AOD   — NTUP

Using Nused quantity to determine reuse
- NTUP – 90 % of data volume 50 times or less
- AOD - 90% of data volume 24 times or less

**Fraction data used** (y-axis): 0, 0.2, 0.4, 0.6, 0.8, 1, 1.2

**Number of Time data is used (NUsed)** (x-axis): 0, 10, 20, 30, 40, 50, 60, 70, 80, 90, 100, 110, 120, 130, 140, 150, 160, 170, 180, 190, 200

# PD2P interim summary

- Jan-Apr 2013 – PD2P moved at peak ~ 600 TB/day
- NTUP most popular data type to move
  - Not a surprise given the popularity plots
- Most data moved to Tier 2 sites
- Data moved to a wide variety of sites (> 67 sites)
- … something about dataset reuse….

# Planned data replication



DDM – Daily Data Brokering data Volume (1-Jan to 1-May 2013)

# Data Brokering transfer Rate



Transfer Throughput (averaged over a day) 1-Jan to 1-May 2013

# Data Brokerage transfer rate
# 28-Jan – 31-Jan (48 hrs)



Transfer Throughput (one hour bins) 28-Jan to 1-Feb 2013

# Transfer rate zoom in further



Transfer Throughput In/Out
2013-01-30 08:00 to 2013-01-30 10:00 UTC

# Planned transfers
# Data Brokering, Group
# Subscriptions, User subscriptions

# Planned Data Brokerage summary

- Peak 400 TB/day and ~ 200 TB/day post Moriond
- Transfer rates
  - Peak 4.5 GB/s (day resolution)
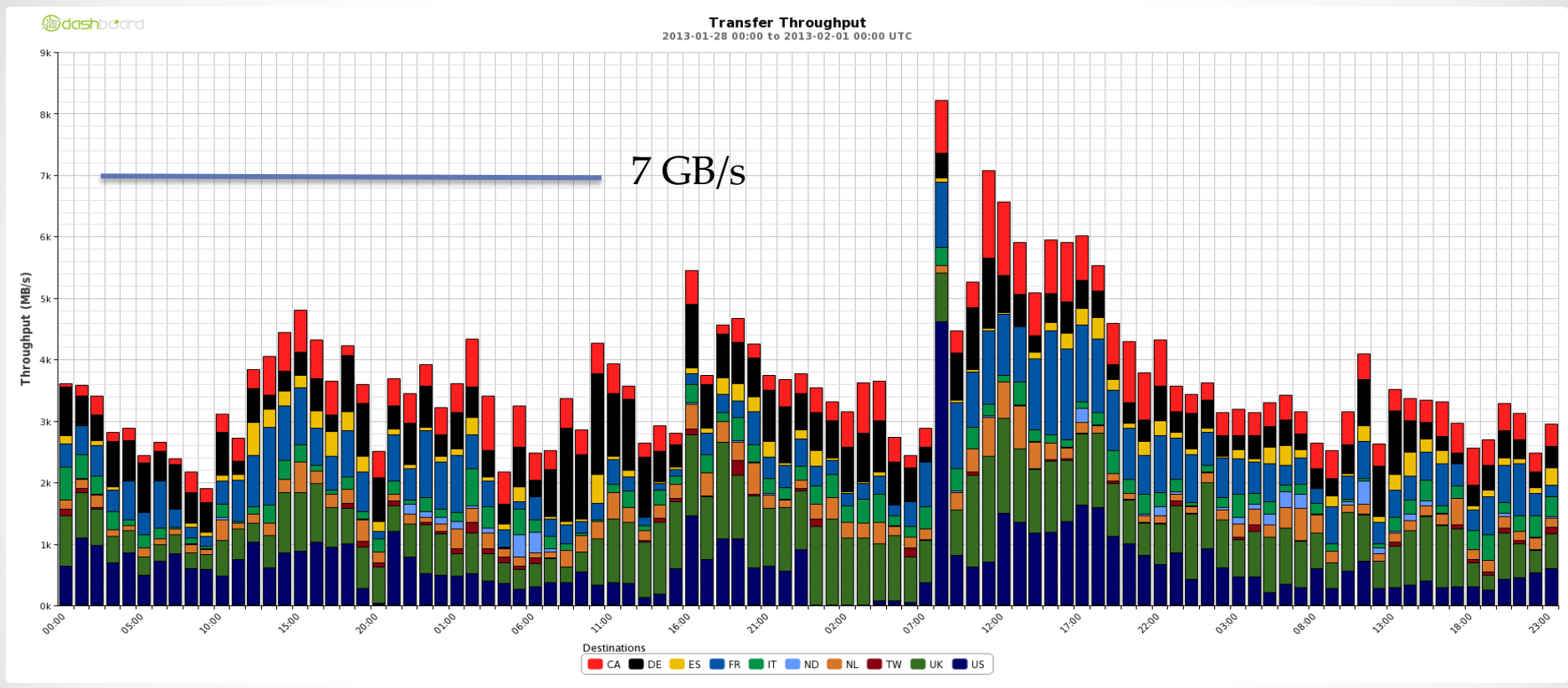  - During busiest few days 28-Jan through 31-Jan-13 (48 hrs)
  7 GB/s (one hour resolution)
  - Jan 30, 2013  (08:00 – 10:00) peak time, in Peak 10 minute period
  30 GB/s -  reading mostly from US Cloud to rest of ATLAS
  ~ 5 GB/s other 10 minute periods
- User data subscriptions

150-200 TB/day  pre Moriond

50-100 TB/day post Moriond

- All networking plans should include accounting for User Output

# Analysis Sites



Completed Jobs per site

Jan 1 to April 30 2013   -    MWT2 – direct read site

# How much of D3PD do users read?



Fraction of file read in User analysis job w/ group Ntuple

Many small reads of
Physics group D3PD files

# How much of D3PD do users read?



Fraction of file read in User analysis job w/ group Ntuple

Semi-log version of previous plot

# How much of D3PD do users read?

**Integral in User analysis job w/ group Ntuple**



- 90% of all access read 10% or less of file.
- Implies if users are reading mostly the same variables – we more a lot of data never read.

*(y-axis)* per 1 percent

*(x-axis)* Fraction of file read (in percent)

# Interim summary  user access of D3PD's

- Current D3PD's from the Physics groups typically very large
- They contain more Branches than the users typically use
- New analysis model is working to merge AOD and D3PD's
  - Implication – input files will contain much more information that the users actually read
  - Efficient Skim/Slim centralized service will be crucial
- Need some mechanism for capturing what variables the users are really using and then provide them  files with mostly those variables.

# Open questions needing further study

- How much PD2D data is read only once or twice after it has been replicated via PD2P?

- How long does a file stay popular (ie read frequently)?
  - We want to replicate the popular files and not the other ones?

- Can we reduce the amount of data in the Analysis files that is rarely read?

- Should we have caches for data files at the Tier 1 and Tier 2 sites ?
  - Do for analysis data what was done for database data and software releases (frontier/squid and cvmfs/squid).
  - What will it take to have partial event caches?

- Can we estimate the network growth needs when jobs and input data are in separate locations (WAN access from jobs)

# Conclusions

- During busy times – 1 Million user analysis and D3PD production jobs per day, user analysis jobs are short duration.

- PD2P moves at peak ~ 600 TB/day and planned replication of similar files ~ 400 TB/day, User subscriptions ~ 200 TB/day

- 90% of NTUP datasets (by volume) are read 50 times or less. (according to PD2P values)

- Users still read a small fraction of the centrally produced D3PD's (NTUP's).  Implies much of the data moved is not read

- What network issues will be seen w/ factor of 3 increase in Trigger Rate in 2015-2016?