

Data access pattern of a Tier2 in ATLAS

Günter Duckeck, Sergey Kalinin
Bergische Universität Wuppertal

Contents

- Short description of T2 at Wuppertal
- Access patterns at Wuppertal and LRZ
- Outlook

T2 at Wuppertal

- Hardware is from HP:
 - WNs: BL220c, 460c, CentOS
 - Storage: DL380G5, 2xRAID6(40TB) pro server, 10 Gb ports
- ~1000 cores(Torque/MAUI), 1 PB storage(dCache)
- New hardware:
 - WN: BL2x220 G7, 386 cores
- T3 at Dortmund use our storage

Sources of information

- Panda(distributed production system) website to get the list of dynamically distributed datasets to Wuppertal.
- dCache logs, no DB interface yet. The logs have detailed transfer description(file id, token, etc).A tool has been developed to parse logs and to extract interesting statistics from them. The results for the two last months are presented in this talk.

The most and the least used datasets

```
length($1)-5,5); print a"\t\t"$2}'
```

Dataset name	Accesses
<u>ddo.000001...01010</u>	50019
<u>mc11_7TeV...945_0</u>	38960
<u>ddo.000001...08010</u>	36595
<u>data11_7TeV...470_0</u>	32385
<u>mc10_7TeV...044_0</u>	31386
<u>data11_7TeV...510_0</u>	28849
<u>mc11_7TeV...665_0</u>	18872
<u>data11_7TeV...712_0</u>	17811
<u>data11_7TeV...698_0</u>	16030
<u>data11_7TeV...511_0</u>	14993
<u>data11_7TeV...629_0</u>	14689
<u>data11_7TeV...264_0</u>	14659
<u>data11_7TeV...4_m93</u>	14405
<u>user.abart...41900</u>	14296
<u>mc11_7TeV...387_0</u>	13964
<u>mc10_7TeV...062_0</u>	13635
<u>20110810_11081</u>	13503
<u>ddo.000001...05010</u>	13482
<u>user.mneum...15393</u>	11677
<u>data11_7TeV...189_0</u>	11568

← alignment, trigger, cond DB

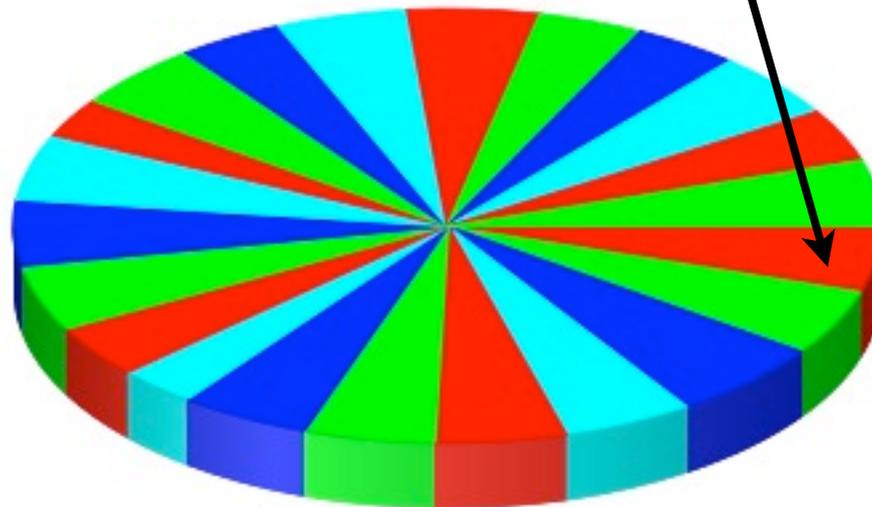
The tool extracts the information about the number of file accesses. 19 datasets out of 50 the most accessed were dynamically replicated to Wuppertal!

← HammerCloud

It also returns, of course, the list of datasets which were not used at all and the owner can be asked.

Transfer statistics for pools

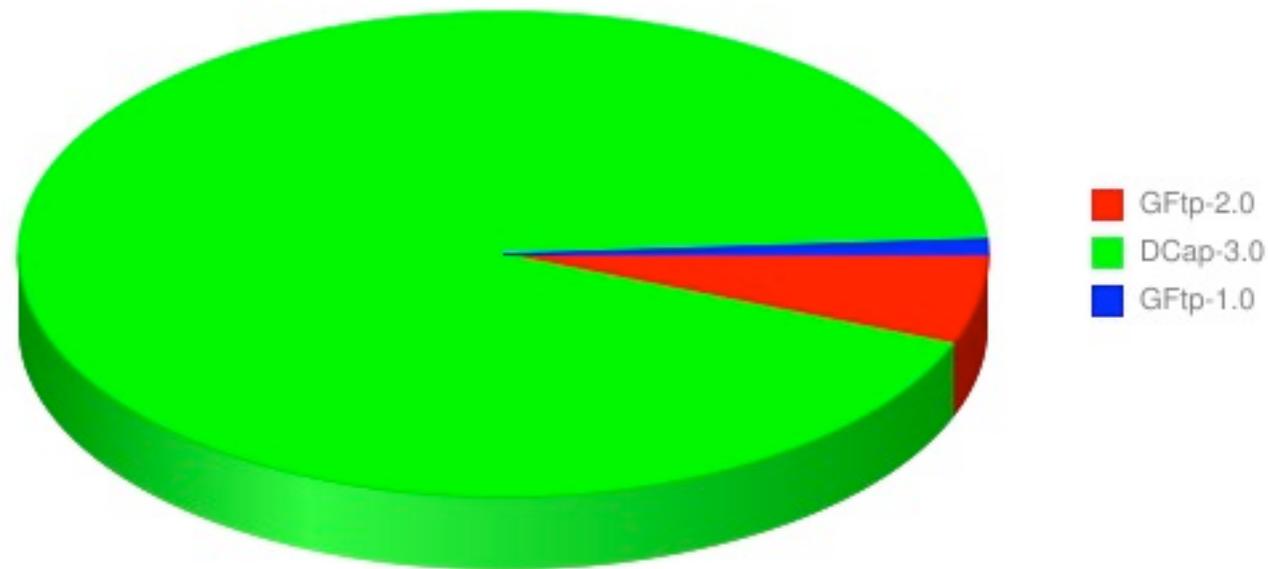
It starts at 3h00 →



- sn01_2@sn01_2Domain
- sn18_1@sn18_1Domain
- sn12_1@sn12_1Domain
- sn05_2@sn05_2Domain
- sn17_2@sn17_2Domain
- sn03_2@sn03_2Domain
- sn14_1@sn14_1Domain
- sn16_2@sn16_2Domain
- sn15_2@sn15_2Domain
- sn03_1@sn03_1Domain
- sn13_1@sn13_1Domain
- sn11_1@sn11_1Domain
- sn16_1@sn16_1Domain
- sn18_2@sn18_2Domain

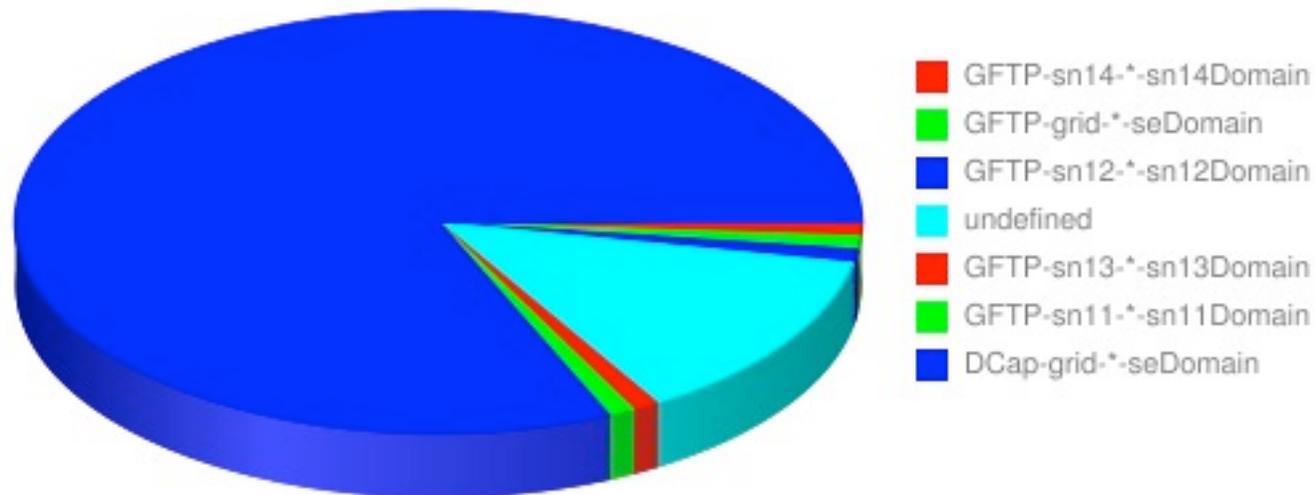
Thanks to our similar hardware,
the distribution is almost perfect.

Transfer statistics for protocols



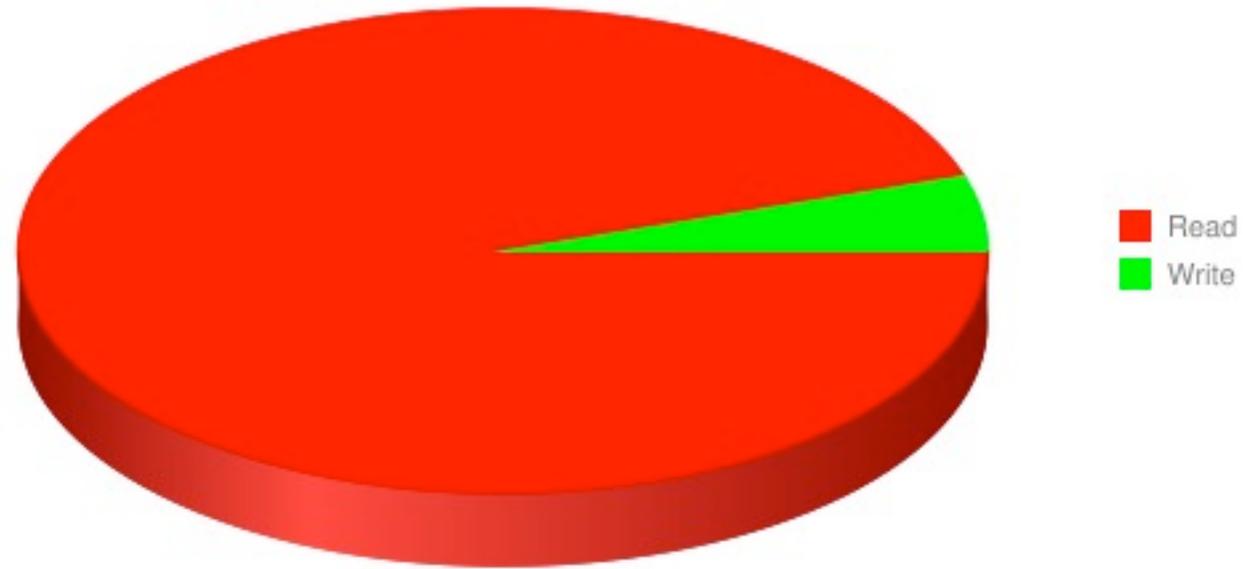
GFtp-1.0 is used by Auger and IceCube. The other ones are by ATLAS.

Transfer statistics for doors



We have 4 GFTP doors and 1 dcap door

Transfer Read/Writes



95% by volume(79% by number of transfers) are reads!

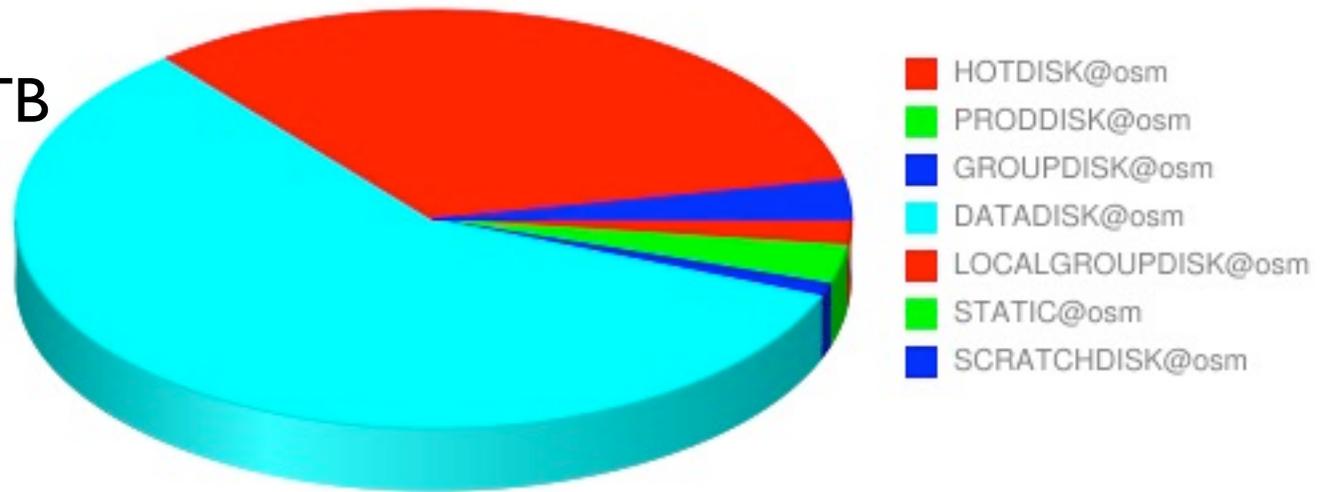
Transfer statistics for space tokens

Used space

DATADISK:471 TB

LOCALGROUPDISK: 118 TB

GROUPDISK: 61 TB

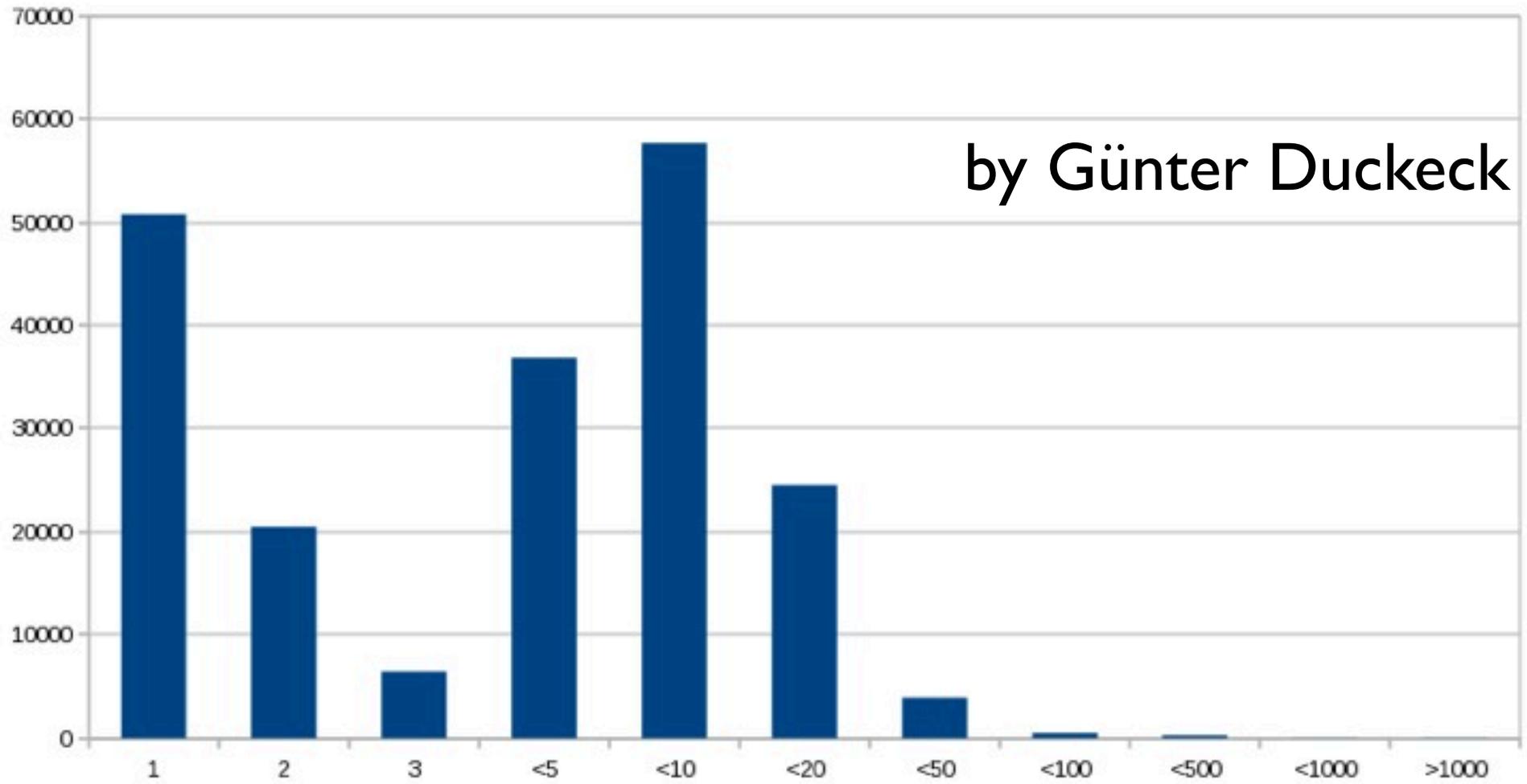


DATADISK is a general purpose token for data and mc.
LOCALGROUPDISK is for local non pledged storage.

File popularity at LRZ

N-read per file
dcap - 7 day

by Günter Duckeck



... nicely pronounced after 1 week ...

Data volume transferred vs storage accessed

time	Number of different	total Vol (TB)	unique Vol (TB)	unique Vol (files read >
1d	67752	25	15	1
2d	115943	56	24	5
7d	201543	186	52	22
40d	525620	929	120	61

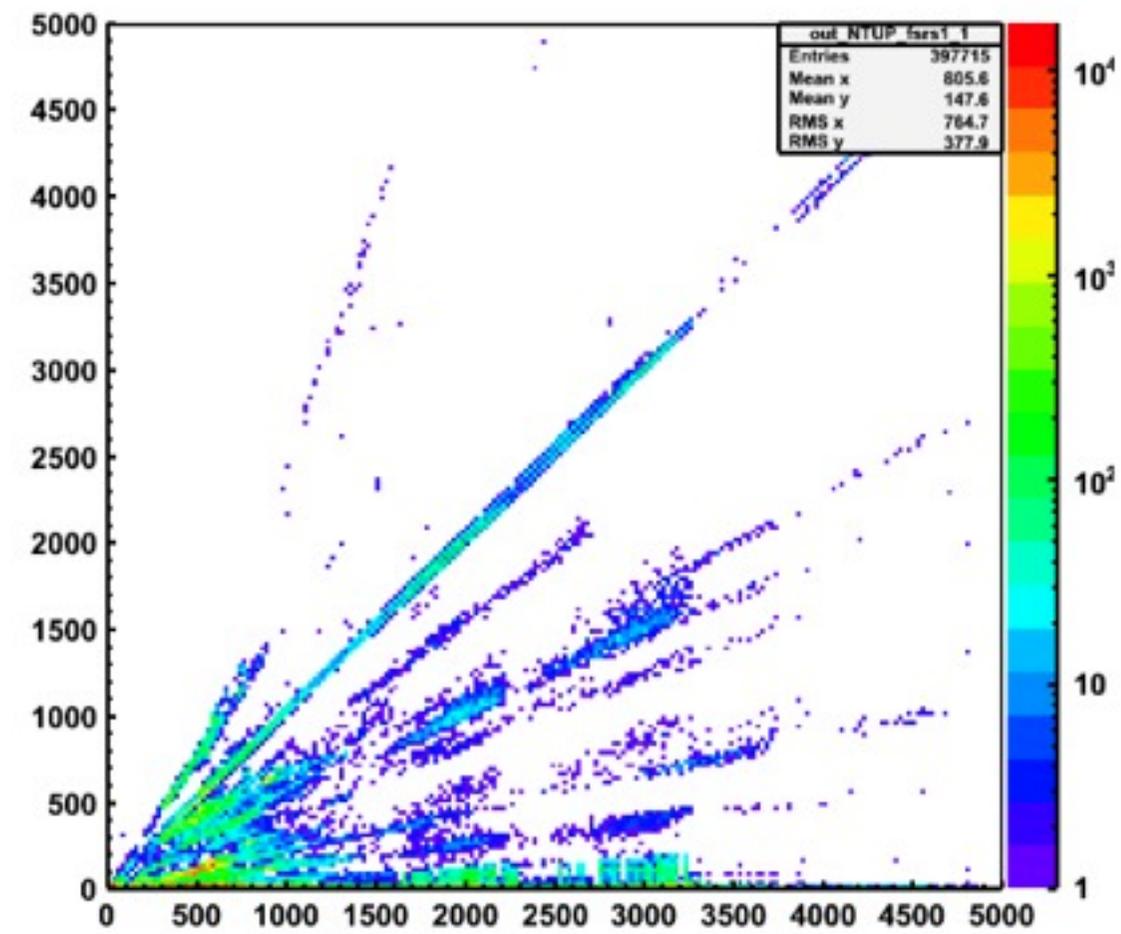
- Only ~20% of data read within 40days (120 TB vs 600 TB)
- For 7 days the ~20 TB most popular files yield ~150 TB of transfer volume
 - Replicating/caching those on fast storage might help a lot for performance

by Günter Duckeck

Transfer size vs File size

- Ntuple

File-size vs Read-size

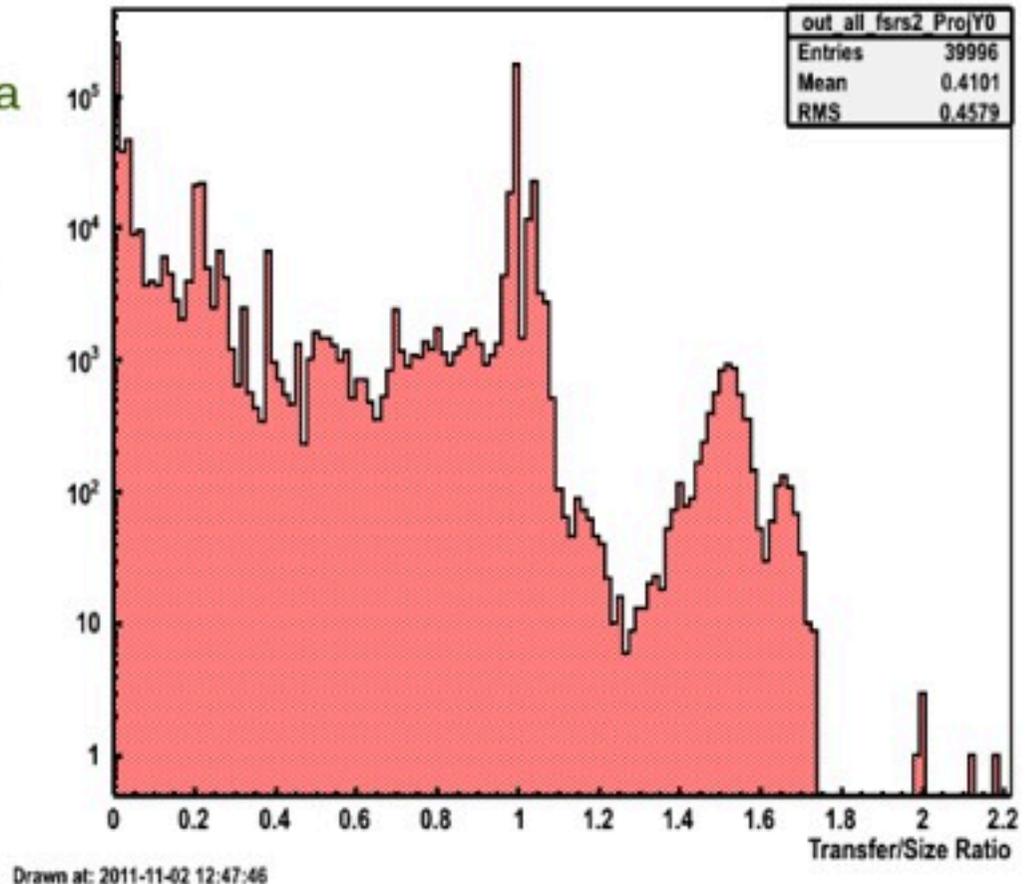


Drawn at: 2011-10-29 00:03:33

by Günter Duckeck

Transfer size/File size ratio

- Often for AOD/ntuple not all data is needed
 - Proven by billing stat when comparing transfer size vs file-size
 - Sometimes backward jumps cause repeated reading of blocks → overhead



by Günter Duckeck

Outlook

- A tool has been developed to show real usage of dCache based mass storage for all VOs and tested at Wuppertal
- ATLAS data replication is effective
- There is a room for further developments(AOS/ESD, transfer destination, etc). Can also be combined with other tools(e.g. developed at Munich)