# Managing data movement between OSG sites in ATLAS experiment

**Alexei Klimentov**

Brookhaven National Laboratory

Jun 16, 2008

*Annual OSG Users Meeting'08*

# Outline

- ATLAS Event Model
- Distributed Data Management (DDM) software and components
- DDM Operation Activity
  - Data replication tests
    - CCRC08, FDR-II, Functional Tests
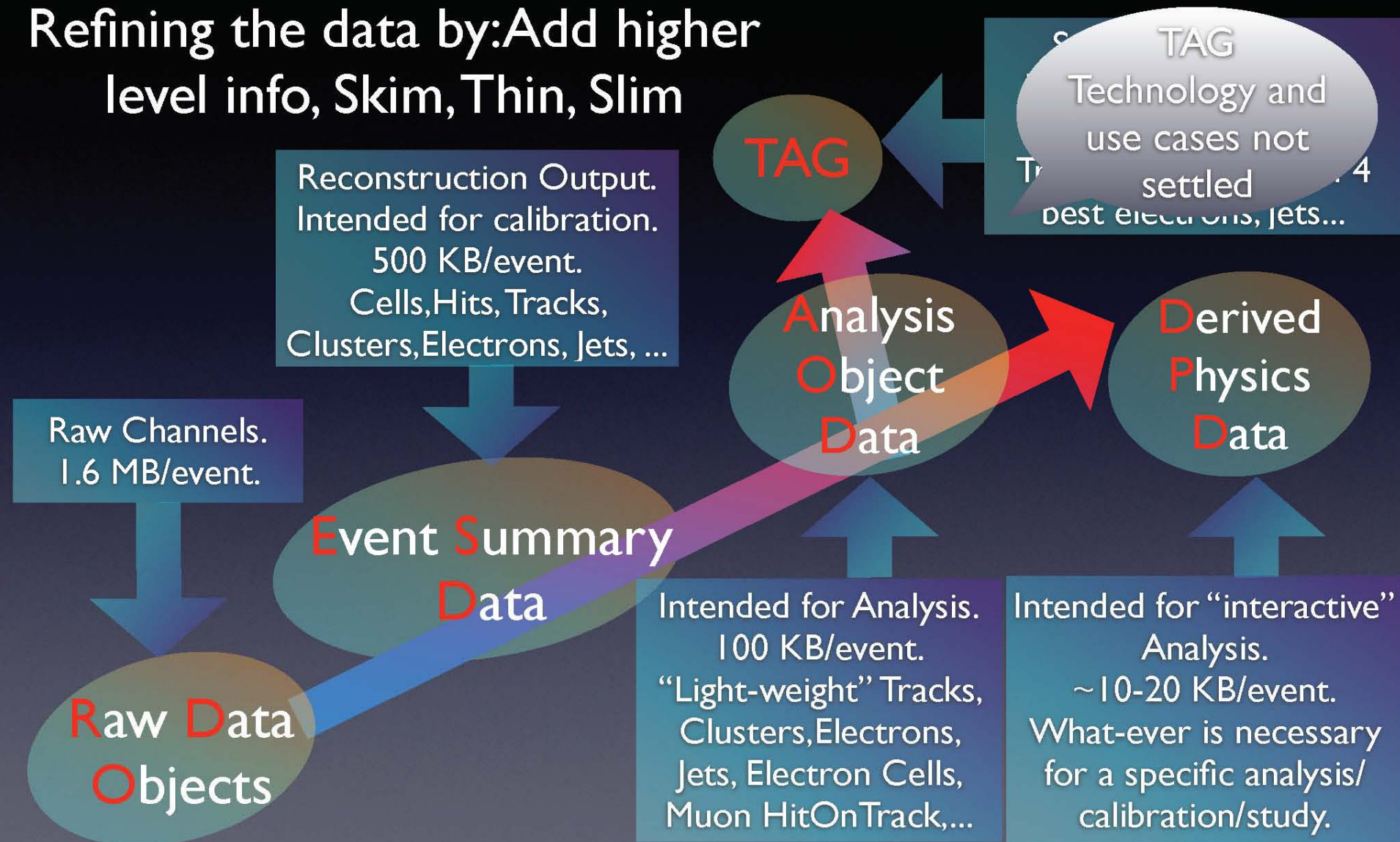- Conclusions and summary

# Event data flow from online to offline

- Events are written in "ByteStream" format by the Event Filter farm in ~2 GB files
  - ~1000 events/file (nominal size is 1.6 MB/event)
  - 200 Hz trigger rate (independent of luminosity)
  - Currently several streams are foreseen:
    - Express stream with "most interesting" events to be processed immediately
    - ~5 event streams, separated by trigger signature
      - e.g. muons, electromagnetic, hadronic jets, taus, minimum bias
    - Calibration events
    - "Trouble maker" events (for debugging)
  - One 2-GB file every 5 seconds will be available from the Event Filter
  - Data will be transferred from the pit to the Tier-0 input buffer at 320 MB/s (average)
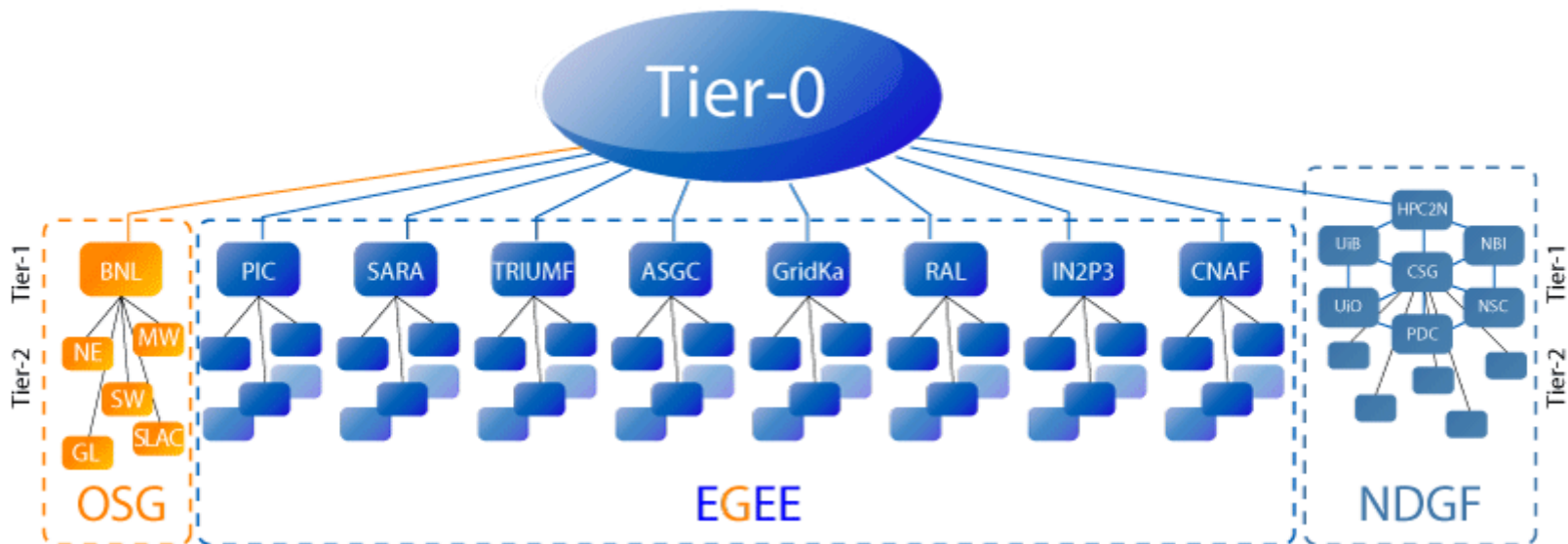
*from D.Barberis*

# The Event Data Model

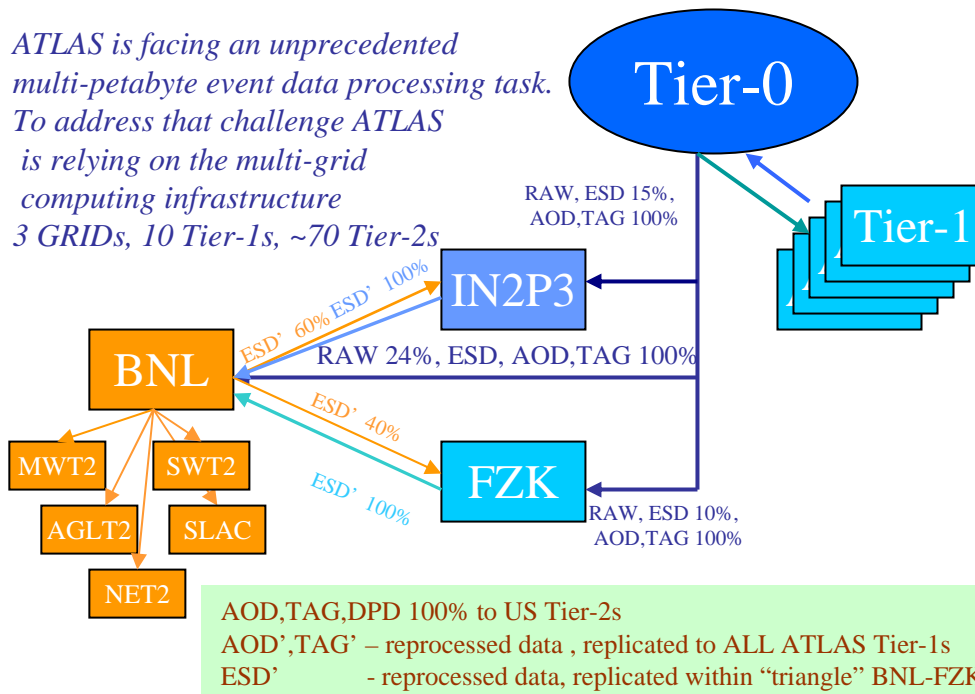Refining the data by: Add higher level info, Skim, Thin, Slim

TAG

Reconstruction Output. Intended for calibration. 500 KB/event. Cells, Hits, Tracks, Clusters, Electrons, Jets, ...

TAG Technology and use cases not settled

best electrons, jets...

**A**nalysis **O**bject **D**ata

**D**erived **P**hysics **D**ata

Raw Channels. 1.6 MB/event.

**E**vent **S**ummary **D**ata

Intended for Analysis. 100 KB/event. "Light-weight" Tracks, Clusters, Electrons, Jets, Electron Cells, Muon HitOnTrack,...

Intended for "interactive" Analysis. ~10-20 KB/event. What-ever is necessary for a specific analysis/ calibration/study.

**R**aw **D**ata **O**bjects

A.Farbin/UTA

# Event data model

- RAW:
  - "ByteStream" format, ~1.6 MB/event
- ESD (Event Summary Data):
  - Full output of reconstruction in object (POOL/ROOT) format:
    - Tracks (and their hits), Calo Clusters, Calo Cells, combined reconstruction objects etc.
  - Nominal size 1 MB/event initially, to decrease as the understanding of the detector improves
    - Compromise between "being able to do everything on the ESD" and "not enough disk space to store too large events"
- AOD (Analysis Object Data):
  - Summary of event reconstruction with "physics" (POOL/ROOT) objects:
    - electrons, muons, jets, etc.
  - Nominal size 100 kB/event (now 200 kB/event including MC truth)
- DPD (Derived Physics Data):
  - Skimmed/slimmed/thinned events + other useful "user" data derived from AODs and conditions data
  - Nominally 10 kB/event on average
    - Large variations depending on physics channels
- TAG:
  - Database (or ROOT files) used to quickly select events in AOD and/or ESD files

*from D.Barberis*

Tier-0

Tier-1

*ATLAS is facing an unprecedented*
*multi-petabyte event data processing task.*
*To address that challenge ATLAS*
*is relying on the multi-grid*
*computing infrastructure*
*3 GRIDs, 10 Tier-1s, ~70 Tier-2s*

RAW, ESD 15%,
AOD,TAG 100%

ESD' 60% ESD' 100%

IN2P3

RAW 24%, ESD, AOD,TAG 100%

BNL

ESD' 40%

ESD' 100%

FZK

RAW, ESD 10%,
AOD,TAG 100%

MWT2    SWT2

AGLT2    SLAC

NET2

AOD,TAG,DPD 100% to US Tier-2s
AOD',TAG' – reprocessed data , replicated to ALL ATLAS Tier-1s
ESD'          - reprocessed data, replicated within "triangle" BNL-FZK-LYON

ATLAS Tier-0 expected  average  rates
– Raw 320 MB/s
– ESD 200 MB/s
– AOD  20 MB/s
– TAG   2 MB/s

Total : 542 MB/s from Tier-0 to Tier-1s
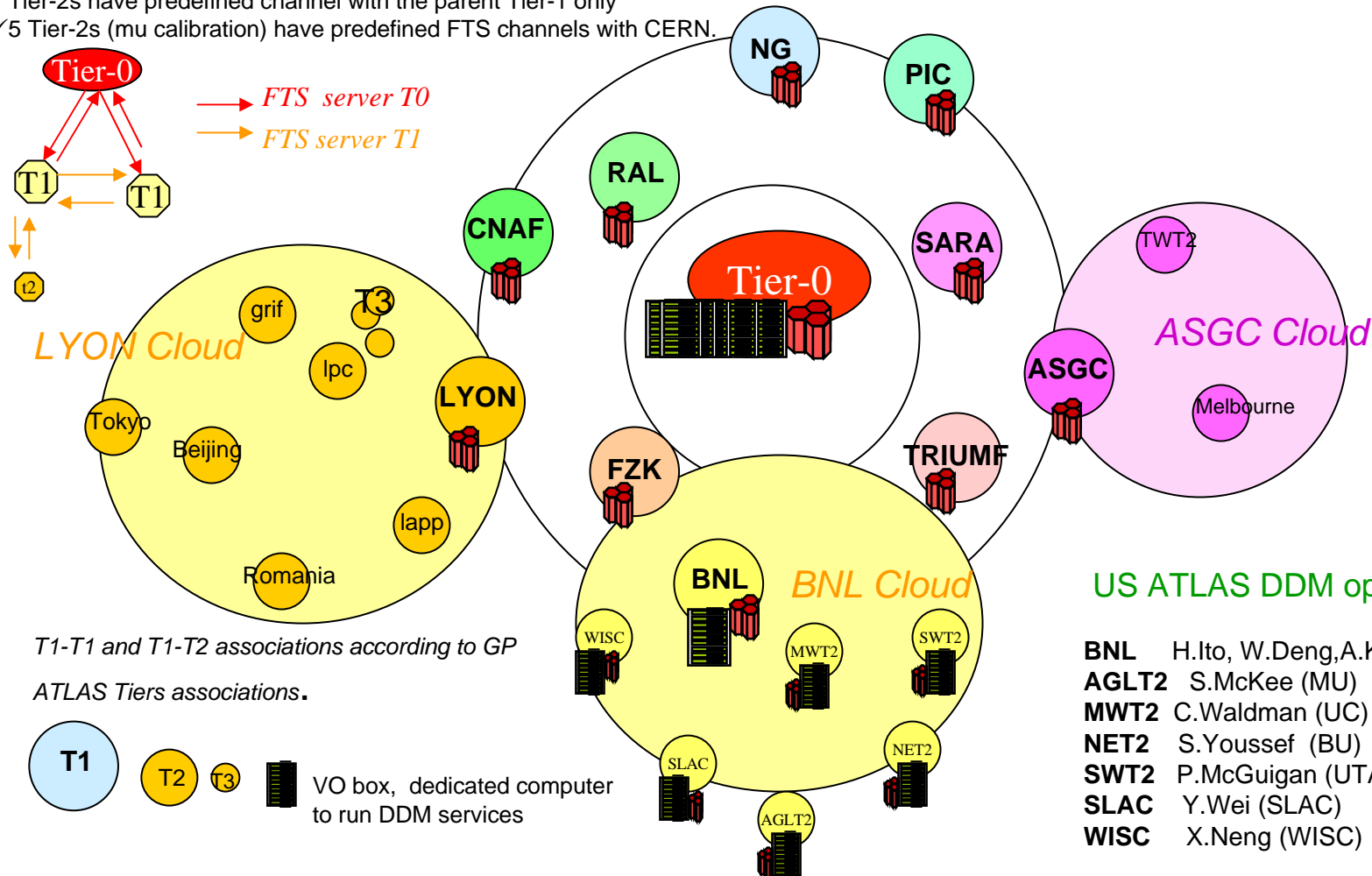
• BNL expected average rates (w/o re-processed data)
– Tape
  • RAW 80 MB/s
– Disk
  • ESD  50+150 MB/s
  • AOD  20 MB/s
  • TAG   2 MB/s

Total BNL  : 302 MB/s

A.Klimentov                          Annual OSG User's meeting '08                                    6

# DDM Deployment and Operations Model

✓EGEE and NDGF clouds have 1 File Catalog (LFC) per cloud
✓US cloud has 1 file catalog (LRC) per site
✓All Tier-1s have predefined FTS channel with CERN and with each other.
✓Tier-2s are associated with one Tier-1 and form the cloud
✓Tier-2s have predefined channel with the parent Tier-1 only
✓5 Tier-2s (mu calibration) have predefined FTS channels with CERN.

*DDM Deployment Model since Jun 2007*

→ *FTS server T0*
→ *FTS server T1*

*LYON Cloud*

*ASGC Cloud*

*BNL Cloud*

*T1-T1 and T1-T2 associations according to GP*

*ATLAS Tiers associations.*

VO box, dedicated computer to run DDM services

US ATLAS DDM operations team :

**BNL**     H.Ito, W.Deng,A.Klimentov,P.Nevski
**AGLT2**  S.McKee (MU)
**MWT2**  C.Waldman (UC)
**NET2**   S.Youssef  (BU)
**SWT2**  P.McGuigan (UTA)
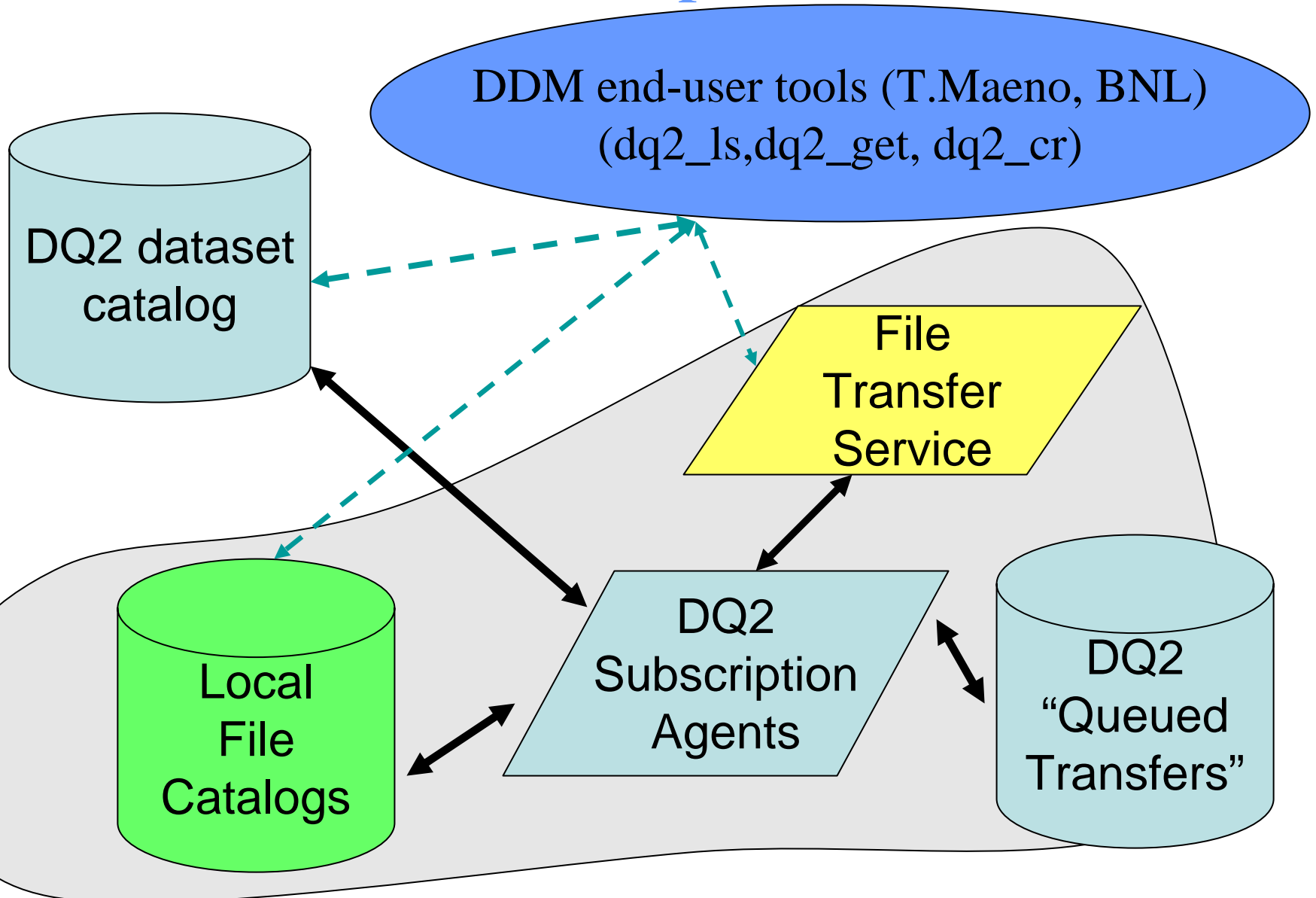**SLAC**   Y.Wei (SLAC)
**WISC**    X.Neng (WISC)

# ATLAS Data Management Software - Don Quijote

- The second generation of the ATLAS DDM system (DQ2)
  - DQ2 developers  M.Branco, D.Cameron, V.Garonne, T.Maeno(BNL), P.Salgado, T.Wenaus(BNL),…
  - Initial idea and architecture were proposed by M.Branco and T.Wenaus (BNL)
- DQ2  is built on top of Grid data transfer tools
  - Moved to *dataset* based approach
    - Datasets : an aggregation of files plus associated DDM metadata
    - Datasets is a unit of storage and replication
    - Automatic data transfer mechanisms using distributed site services
      - Subscription system
      - Notification system
  - SW releases deployed in 2007 /08
    - Central services : DQ2  0.3 (Jun 07)
    - Site Services : DQ2 0.4 (Oct 07)
    - Site Services : DQ2 0.5 (Dec 07)
    - Site Services : DQ2 0.6 (Mar-Apr 08)
    - Central catalogs, containers, site services : DQ2 1.0, 1.1 (May-Jun 08)

# Major Changes and Improvements in DDM/DQ2 Software (2007/08)

- 6 Software release since Jun 2007 to address Operations needs and to improve system performance
    - Switch to ORACLE backend for central datasets catalog
        - Performance and stability are improved
    - [ARDA] Monitoring
        - Data navigation
        - Cloud views
        - Errors summary
    - Site services
        - Dataset replicas look up is improved
        - Dataset subscription processing is improved
    - New packaging and deployment procedures
    - Fair-share implementation
        - Used for critical data replication (conditions data, database releases, etc)
    - New central catalogs schema
    - Containers implementation

# DDM components



DDM end-user tools (T.Maeno, BNL)
(dq2_ls,dq2_get, dq2_cr)

DQ2 dataset catalog

File Transfer Service

Local File Catalogs

DQ2 Subscription Agents

DQ2 "Queued Transfers"

# DDM Activities.

- Cosmic data replication
- Critical data (conditions datasets, database releases) replication
- Analysis Object Data (AOD) replication for physics analysis
- Support MC production
- DDM monitoring and control
- Data integrity check
- DDM Functional and throughput tests
- User and group support

# Conditions data replication monitoring page

**Conditions Datasets Distribution**

- Datasets are automatically subscribed to Tier-1s from CERN.
- *darkgreen* - site has a complete dataset replicas (data transfer is done)
- *green* - site has the same number of files as at CERN
- *lightgreen* - site has 90% of files at CERN
- *orange* - site has an incomplete dataset replicas. It also means that subscription is
- *red* - the subscription is not processed
- UTC time is used within the page

*Comments*

*Critical Data Replication (conditions data) BNL and US Tier-2s have a complete replicas of DB releases and conditions data. Monitoring is integrated with Panda*

| Datasets | Total Files in datasets | Last Subscription | LFC Checked | Last Transfer |
|---|---|---|---|---|
| 23 | 1994 | Jun 11 15:49:41 | Jun 11 19:45:19 | Jun 11 19:45:19 |

| Dataset | ASGC | BNL | CNAF | FZK | LYON | NDGF | PIC | RAL | SARA | TRIUMF |
|---|---|---|---|---|---|---|---|---|---|---|
| cmccond.000001.conditions.recon.pool.v0000 (Subscription Time: Jul 20 2007 09:16) | (10/10) | (10/10) | (10/10) | (10/10) | (10/10) | (10/10) | (10/10) | (10/10) | (10/10) | (10/10) |
| cmccond.000001.conditions.simul.pool.v0000 (Subscription Time: Sep 14 2007 11:56) | (24/24) | (24/24) | (24/24) | (24/24) | (24/24) | (24/24) | (24/24) | (24/24) | (24/24) | (24/24) |
| comcond.000001.conditions.recon.pool.v0000 (Subscription Time: Sep 14 2007 11:56) | (26/26) | (26/26) | (26/26) | (26/26) | (26/26) | (26/26) | (26/26) | (28/28) | (30/30) | (26/26) |
| comcond.000001.lar_conditions.recon.pool.v0000 (Subscription Time: Mar 04 2008 10:49) | (352/352) | (352/352) | (352/352) | (352/352) | (352/352) | (352/352) | (352/352) | (352/352) | (352/352) | (352/352) |
| comcond.000002.lar_conditions.recon.pool.v0000 (Subscription Time: Jun 11 2008 15:49) | (142/142) | (142/142) | (142/142) | (142/142) | (142/142) | (142/142) | (142/142) | (142/129) | (142/142) | (142/142) |
| comcond.000002.lar_conditions.recon.pool.v000001 (Subscription Time: Mar 04 2008 10:49) | (1/1) | (1/1) | (1/1) | (1/1) | (1/1) | (1/1) | (1/1) | (1/1) | (1/1) | (1/1) |
| comcond.000002.lar_conditions.recon.pool.v000002 (Subscription Time: Mar 04 2008 10:49) | (1/1) | (1/1) | (1/1) | (1/1) | (1/1) | (1/1) | (1/1) | (1/1) | (1/1) | (1/1) |
| comcond.000002.lar_conditions.recon.pool.v000003 (Subscription Time: Mar 04 2008 10:49) | (1/1) | (1/1) | (1/1) | (1/1) | (1/1) | (1/1) | (1/1) | (1/1) | (1/1) | (1/1) |
| comcond.000003.lar_conditions.recon.pool.v0000 (Subscription Time: Mar 04 2008 10:49) | (129/129) | (20/20) | (129/129) | (20/20) | (20/20) | (129/129) | (129/129) | (129/129) | (347/906) | (129/129) |
| comcond.000004.lar_conditions.recon.pool.v0000 (Subscription Time: Mar 04 2008 10:49) | (382/382) | (166/166) | (166/166) | (166/166) | (166/166) | (166/166) | (382/382) | (166/166) | (383/2403) | (166/166) |
| oflcond.000001.bfield_conditions.simul.pool.v0000 (Subscription Time: Mar 04 2008 10:49) | (3/3) | (3/3) | (3/3) | (3/3) | (3/3) | (3/3) | (3/3) | (3/3) | (3/3) | (3/3) |

**Panda monitor**
Now in UTC
Shift log   Wiki

**Jobs** - search
Recent running, activated, waiting, assigned, defined, finished, failed jobs
Select analysis, prod, install, test jobs
**Quick search**
Job
Dataset
Task request
Task status
File

**Summaries**
Blocks: ___ days
Errors: ___ days
Nodes: ___ days
Daily usage

**Tasks** - search
Generic Task Req
EvGen Task Req
CTBsim Task Req
Task list
New Tag
Bug Report
Task browser

**Datasets** - search
Dataset browser
Aborted MC datasets
Panda subscriptions

**Datasets Distribution**
DDM Req
Req list
AODs
EVNTs
RDOs
Conditions DS
DB Releases
Validation Samples
Functional Tests
CosmicRuns
FDR_Datasets

Not logged in. List users

# DDM Replication Metrica

- Test Software Components
- Combined Computing Readiness Challenge of 2008
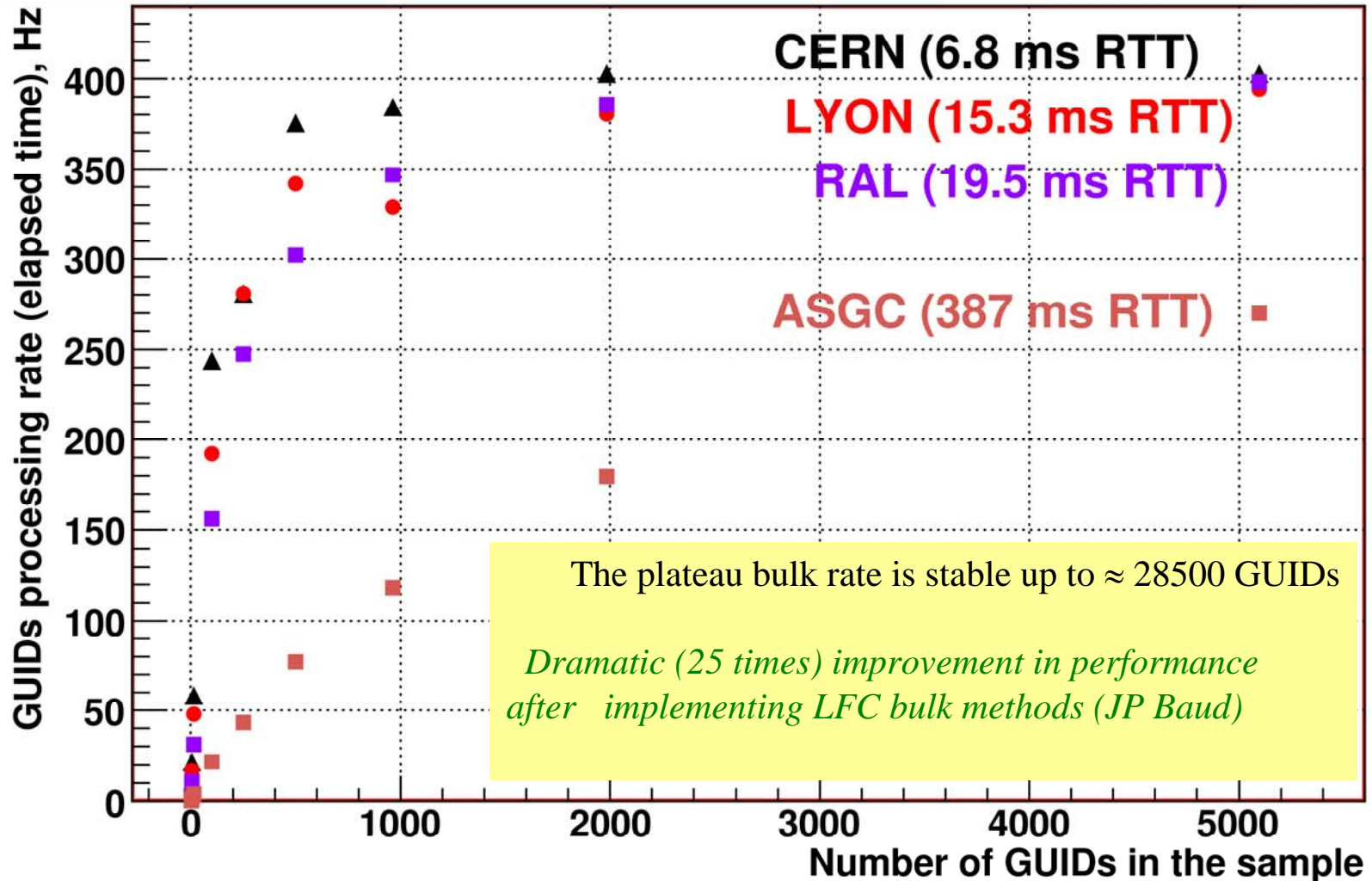- Full Dress Rehearsal
- Functional and Throughput Tests

# Test DDM SW Components. LFC

- LFC – one of vital components, out of DDM
- It is important the two work together as expected
- We organized systematic measurements to understand performance to spot and fix problems...
- Initial LFC performance was found poor : 12 Hz
- It was joint effort of ATLAS DDM Operations team, DQ2 developers, CERN ARDA and LFC Author to understand and to improve the catalog's performance.
- ATLAS performance requirement was driven by our computing and event model

# Test DDM SW Components

## Results on  Performance Testingof the LFC @ CERN



### GUIDs processing rate vs number of GUIDs

GUIDs processing rate (elapsed time), Hz

CERN (6.8 ms RTT)
LYON (15.3 ms RTT)
RAL (19.5 ms RTT)

ASGC (387 ms RTT)

The plateau bulk rate is stable up to ≈ 28500 GUIDs

*Dramatic (25 times) improvement in performance after   implementing LFC bulk methods (JP Baud)*
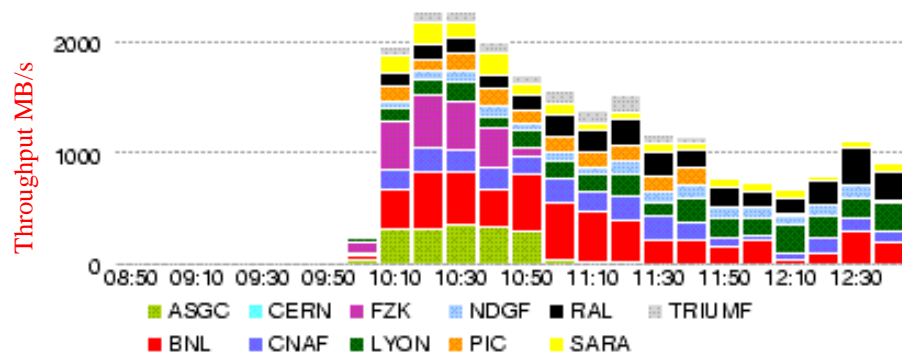
Number of GUIDs in the sample

# CCRC08

- Combined Computing Readiness Challenge of 2008
  - LHC-wide computing challenge,
    - a preparatory Phase 1 series of tests conducted in February,
    - main Phase 2 conducted in May, 2008.
  - ATLAS in CCRC08
    - Tests carried along for the all month
      - **CCRC08 ONLY during week days**
      - Cosmic data during the weekend (commissioning and M7)
    - Focused on data distribution according to Computing Model
    - Tier0->Tier1's, Tier1->Tier1's, Tier1->Tier2's
    - Very demanding metrics
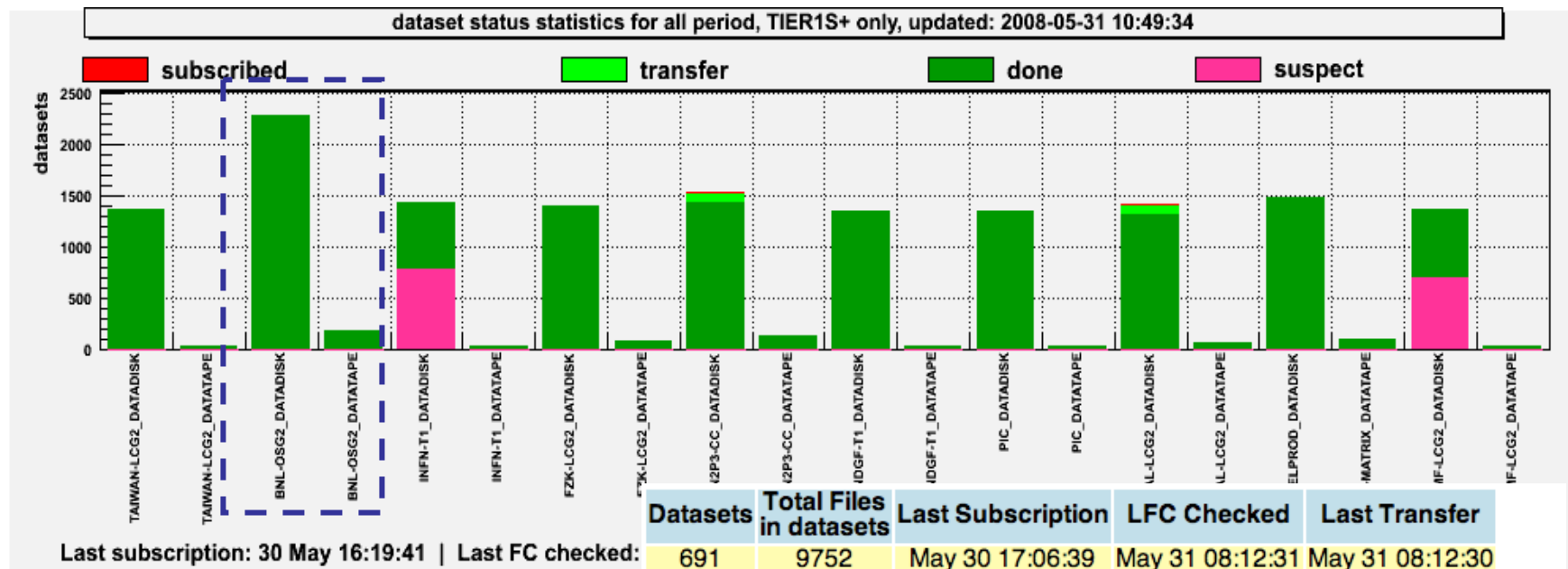      - More than you will need to do during 2008 data taking

# CCRC08 (Tier0-Tier1's)

- T0->T1: sites should demonstrate to be capable to import 100% of the subscribed datasets (complete datasets) within 6 hours from the end of the exercise



| | | Transfers | | | Registrations | | Errors | |
|---|---|---|---|---|---|---|---|---|
| Cloud | Efficiency | Throughput | Successes | Datasets | Files | Transfer | Registration |
| ASGC | 100% | 219 MB/s | 300 | 46 | 300 | 0 | 0 |
| BNL | 100% | 471 MB/s | 597 | 10 | 597 | 0 | 0 |
| CERN | 0% | 0 MB/s | 0 | 0 | 0 | 0 | 0 |
| CNAF | 100% | 195 MB/s | 196 | 17 | 196 | 0 | 0 |
| FZK | 100% | 229 MB/s | 331 | 40 | 329 | 0 | 0 |
| LYON | 99% | 147 MB/s | 155 | 9 | 156 | 2 | 0 |
| NDGF | 100% | 83 MB/s | 98 | 22 | 98 | 0 | 0 |
| PIC | 100% | 132 MB/s | 156 | 19 | 156 | 0 | 0 |
| RAL | 99% | 154 MB/s | 152 | 17 | 152 | 1 | 0 |
| SARA | 100% | 132 MB/s | 207 | 16 | 208 | 0 | 0 |
| TRIUMF | 100% | 105 MB/s | 94 | 26 | 92 | 0 | 0 |

# CCRC08 (Tier0-Tier1's)



dataset status statistics for all period, TIER1S+ only, updated: 2008-05-31 10:49:34

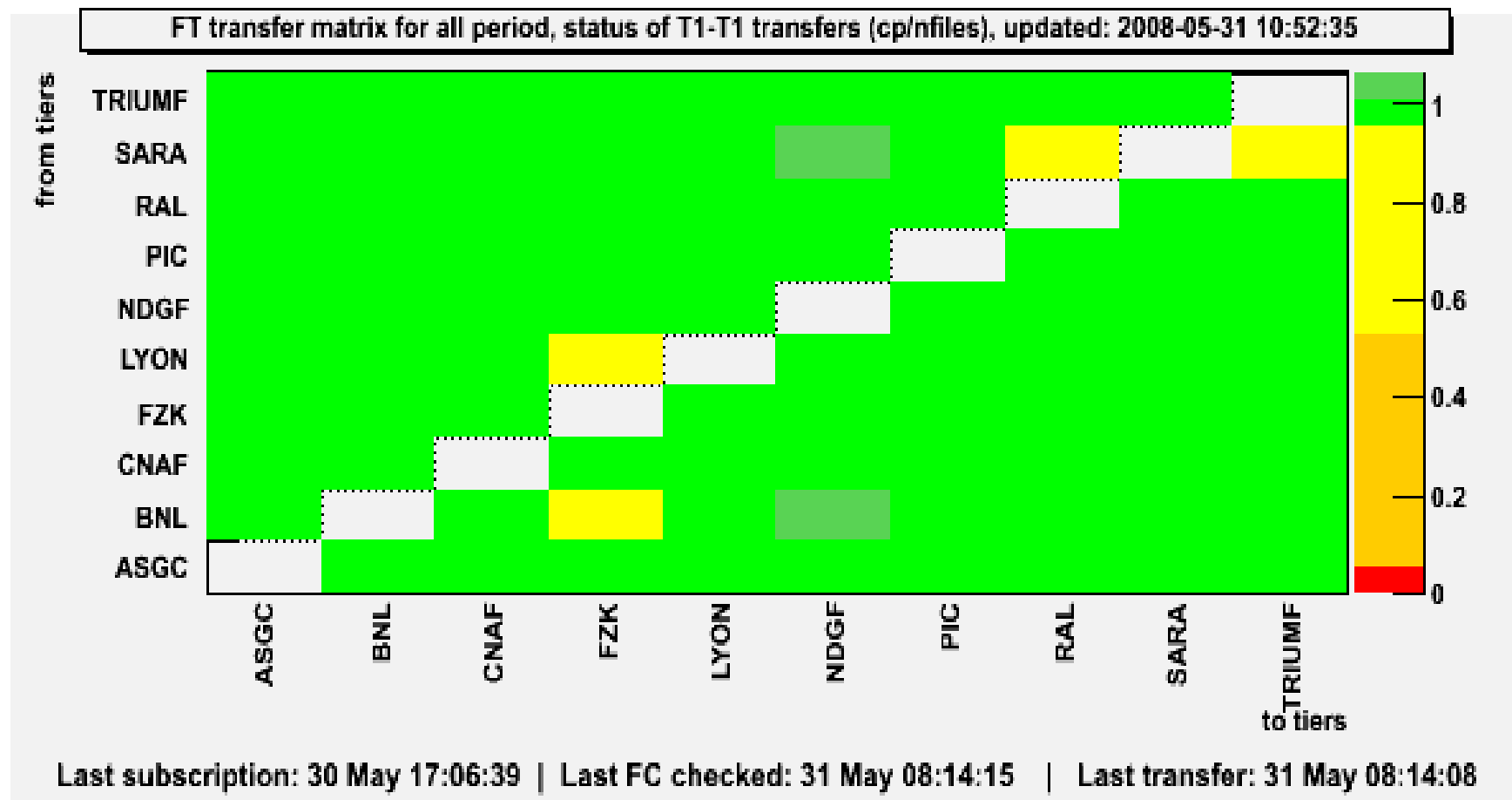| Datasets | Total Files in datasets | Last Subscription | LFC Checked | Last Transfer |
|----------|------------------------|-------------------|-------------|---------------|
| 691 | 9752 | May 30 17:06:39 | May 31 08:12:31 | May 31 08:12:30 |

BNL managed to get all data with average throughput 470+MB/s (1.55 times more than nominal)

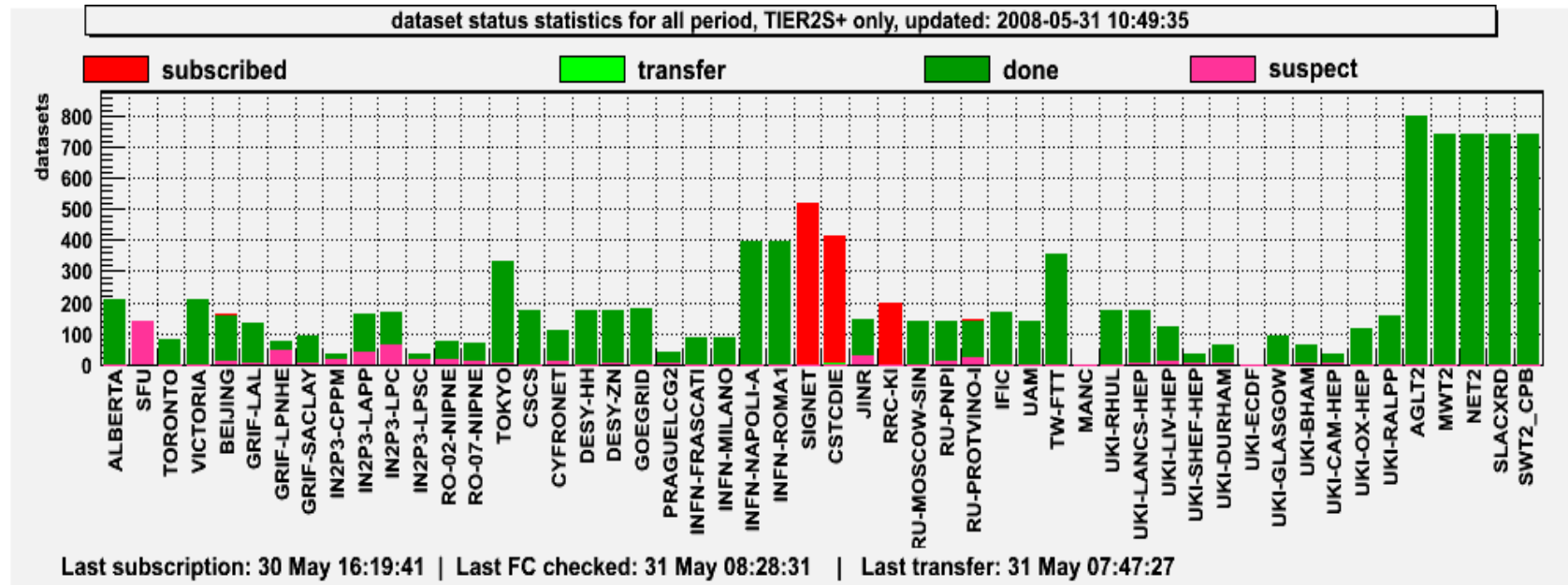All transfers to BNL instances (DISK and TAPE) were completed

Backlog was recovered within 24h (after CERN power cut)

| Tier1 | Datasets | Total Files in datasets | Total CpFiles in datasets | Completed | Transfer | Subscribed |
|-------|----------|------------------------|---------------------------|-----------|----------|------------|
| BNL | 549 | 8170 | 8170 | 549 | 0 | 0 |
| FZK | 442 | 3400 | 3097 | 422 | 9 | 11 |
| IN2P3 | 432 | 3528 | 3432 | 426 | 0 | 6 |
| INFN | 464 | 3530 | 3530 | 464 | 0 | 0 |
| NDGF | 477 | 4033 | 4137 | 472 | 0 | 0 |
| PIC | 483 | 4046 | 4044 | 482 | 0 | 1 |
| RAL | 505 | 5013 | 4900 | 485 | 18 | 2 |
| SARA | 421 | 3137 | 3136 | 420 | 0 | 1 |
| TAIWAN | 470 | 4050 | 4036 | 464 | 5 | 1 |
| TRIUMF | 488 | 4221 | 4120 | 477 | 10 | 1 |

# CCRC08 : Tier1-Tier1's



FT transfer matrix for all period, status of T1-T1 transfers (cp/nfiles), updated: 2008-05-31 10:52:35

Last subscription: 30 May 17:06:39 | Last FC checked: 31 May 08:14:15 | Last transfer: 31 May 08:14:08

# CCRC08 : Tier1-Tier2's



dataset status statistics for all period, TIER2S+ only, updated: 2008-05-31 10:49:35

■ subscribed   ■ transfer   ■ done   ■ suspect

Last subscription: 30 May 16:19:41 | Last FC checked: 31 May 08:28:31 | Last transfer: 31 May 07:47:27

Datasets subscribed:
-upon completion at T1 -every 4 hours

All Clouds

AODs distribution in US cloud
CCRC08 May 16 – Jun 14

# Full Dress Rehearsal – II (FDR-II)

*ATLAS tried to practice everything — from data coming off the experiment right through to it being shipped around and analyzed — under conditions just as they will be when real data–taking begins.*

Test scope :

– Simulated data in RAW data format are pre-loaded on the output buffers of the online computing farm and transmitted to the Tier-0 farm at nominal rate (200 Hz, 320 MB/s), mimicking the LHC operation cycle

– Data are calibrated/aligned/reconstructed at Tier-0 and distributed to Tier-1 and Tier-2 centres, following the computing model

– At the same time, distributed simulation production and distributed analysis activities continue, providing a constant background load

– Reprocessing at Tier-1s is also tested in earnest for the first time

# FDR-II. Data sharing within clouds

| | | | | |
|---|---|---|---|---|
| **TAIWAN** | TW-FTT_DATADISK | 50% | | **registered: 2** |
| | AU-ATLAS | | | **ready: 1** |
| **UK** | UKI-LT2-RHUL_DATADISK | 50% | Egamma | |
| | UKI-NORTHGRID-LANCS-HEP_DATADISK | 50% | Bphys | |
| | UKI-NORTHGRID-LANCS-HEP_DATADISK | 50% | Muon | |
| | UKI-NORTHGRID-LIV-HEP_DATADISK | 50% | Egamma | |
| | UKI-NORTHGRID-SHEF-HEP_DATADISK | 50% | Minbias | |
| | UKI-SCOTGRID-GLASGOW_DATADISK | 50% | Bphys | |
| | UKI-SCOTGRID-GLASGOW_DATADISK | 50% | Minbias | **registered: 12** |
| | UKI-SOUTHGRID-BHAM_DATADISK | 50% | Jet | **ready: 9** |
| | UKI-SOUTHGRID-CAM-HEP_DATADISK | 50% | Muon | |
| | UKI-SOUTHGRID-OX-HEP_DATADISK | 50% | Jet | |
| | UKI-SOUTHGRID-RALPP_DATADISK | 50% | Egamma | |
| | MANC | | | |
| | UKI-SCOTGRID-DURHAM_DATADISK | | | |
| | UKI-SCOTGRID-ECDF_DATADISK | | | |
| **USA** | AGLT2_DATADISK | 100% | | |
| | MWT2_DATADISK | 100% | | |
| | NET2_DATADISK | 100% | | |
| | SLACXRD_DATADISK | 100% | | **registered: 8** |
| | SWT2_CPB_DATADISK | 100% | | **ready: 5** |
| | MWT2_IU | | | |
| | OU | | | |
| | WISC | | | |

Each cloud defines data sharing policy
US : 100% per T2
UK : share data stream between T2s

Done

# FDR-II (data replication)
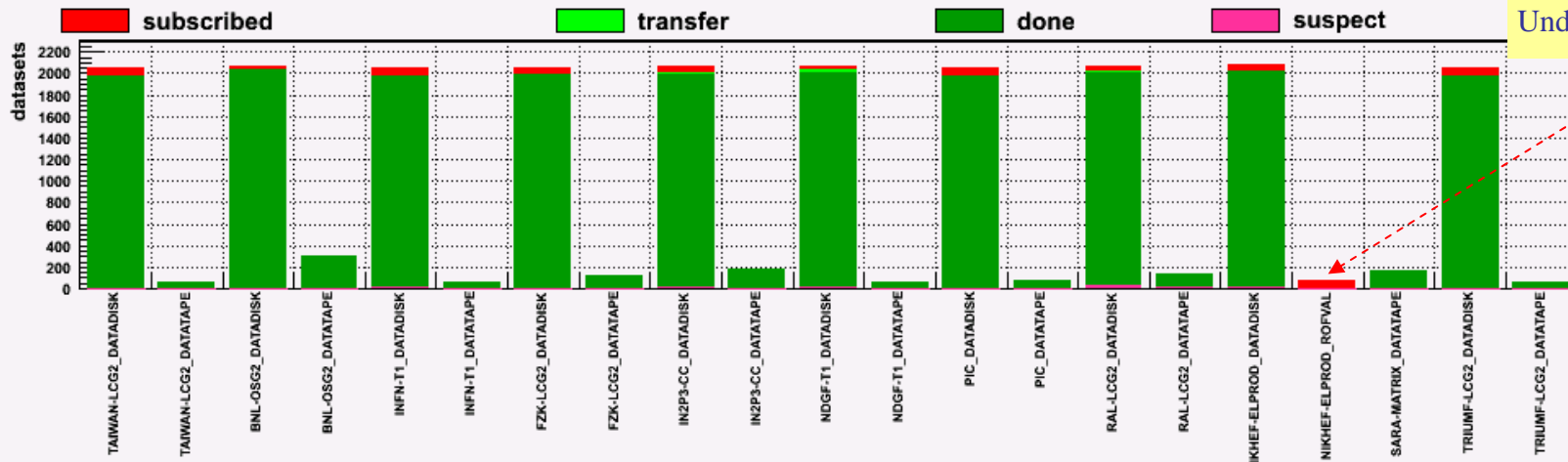
# FDR-II DDM issues (summary)

- From ATLAS Physics Coordination
  - Data transport - OK
    - DDM generally worked well - both in shipping files to CERN and shipping them out after FDR running Shipped 18h of FDR-2 RAW data out within <24h, successfully to all Tier-1s and Tier-2s
- From ATLAS Distributed Computing (FDR-II+CCRC08)
  - Known problems :
    - Double files registration
      - The problem never happened in US cloud
        » Working idea of DQ2 developers that it is in between DQ2 and LFC SW
      - Description
        » The file is transferred correctly to site and registered in LFC
        » "Something" goes wrong and the file is replicated again
        » Another entry in LFC, same GUID, different SURL
    - Delay up to 4h with central catalog update (site replicas info)
      - The problem is understood and fixed
    - Delay up to 2h with start of data replication

# Functional and Throughput Tests

– Starting from 2007 ATLAS Operations runs FT and TT

    *FT was proposed by P.Nevski and A.Klimentov to test system functionality and performance*

- Initially 5 times per year after major SW releases
- From Jun 2008 : continuously to exercise data transfer

– Test scope

    » Run data generator at 10% of nominal rate:

    » Distribute RAW, ESD, AOD according to Computing Model

    » Tier-0 – Tier-1's data transfer

    » Tier-1-Tier-1's data transfer

    » Tier-1-Tier-2's data transfer

    » Subscribe 'calibration' datasets from CERN to 5 Tier-2s + CNAF and BNL

    » 10 Tier-1s and 60 Tier-2s are participating

    » Statistics is generated automatically

    » Exercise Operations shifts

# Functional Tests (cont)



dataset status statistics for all period, TIER1S+ only, updated: 2008-06-15 21:15:20

Site was set up Jun 14th
Under investigation

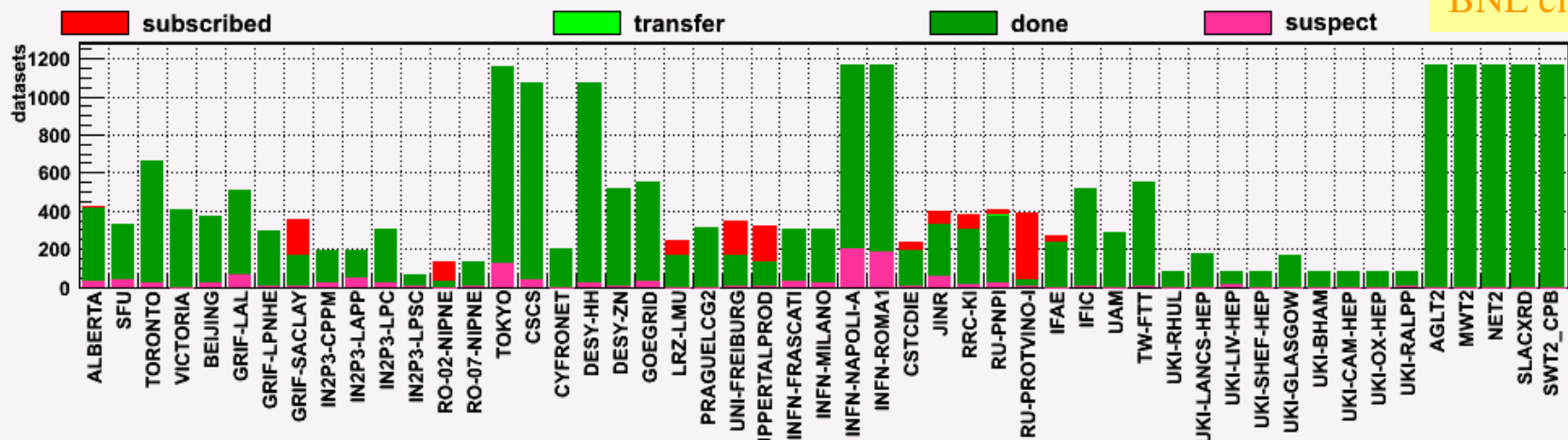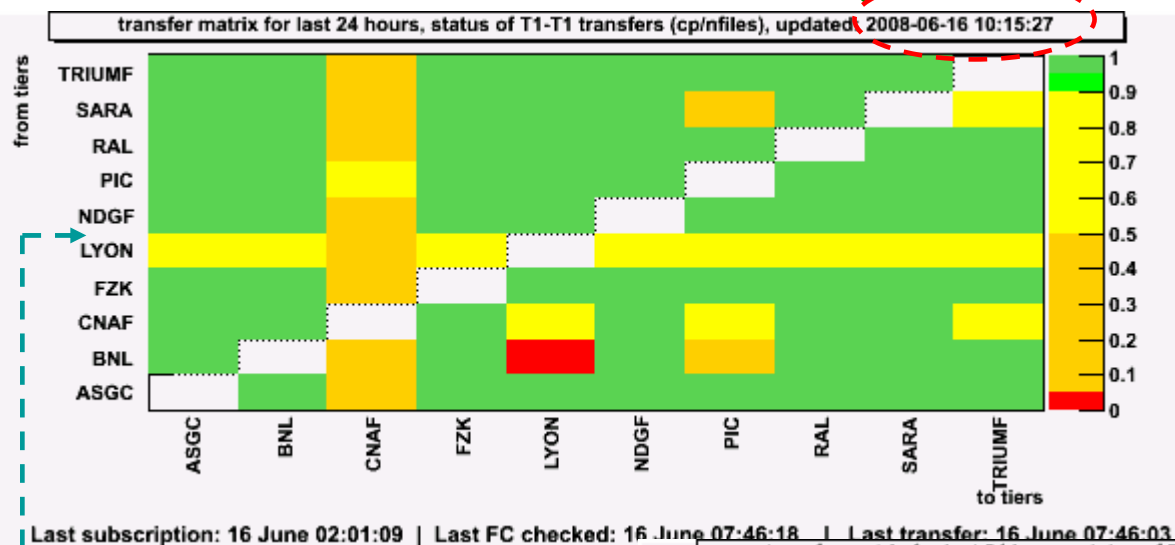Last subscription: 15 June 15:45:15  |  Last FC checked: 15 June 19:12:14  |  Last transfer: 15 June 19:11:46

dataset status statistics for all period, TIER2S+ only, updated: 2008-06-15 21:15:22

BNL cloud

Last subscription: 15 June 15:45:15  |  Last FC checked: 15 June 19:12:14  |  Last transfer: 15 June 19:11:46

# Functional Test ("current" status)



transfer matrix for last 24 hours, status of T1-T1 transfers (cp/nfiles), updated 2008-06-16 10:15:27

Last subscription: 16 June 02:01:09  |  Last FC checked: 16 June 07:46:18  |  Last transfer: 16 June 07:46:03

FT data replication status for the last 24h
LYON and CNAF experienced problems
Data replication performance between
other Tier-1s is close to 100%



transfer matrix for last 24 hours, status of T1-T1 transfers (cp/nfiles), updated 2008-06-16 11:22:15

Last subscription: 16 June 02:01:09  |  Last FC checked: 16 June 09:13:23  |  Last transfer: 16 June 09:13:18

*and 1h later LYON
recovered backlog
Problem with INFN persists*

# Conclusions

- The data distribution scenario has been tested well beyond the use case for 2008 data taking

- US ATLAS Computing and Networking infrastructure met the experiment's requirements.

- ATLAS DDM Software is stable and met data replication requirements

# Acknowledgements

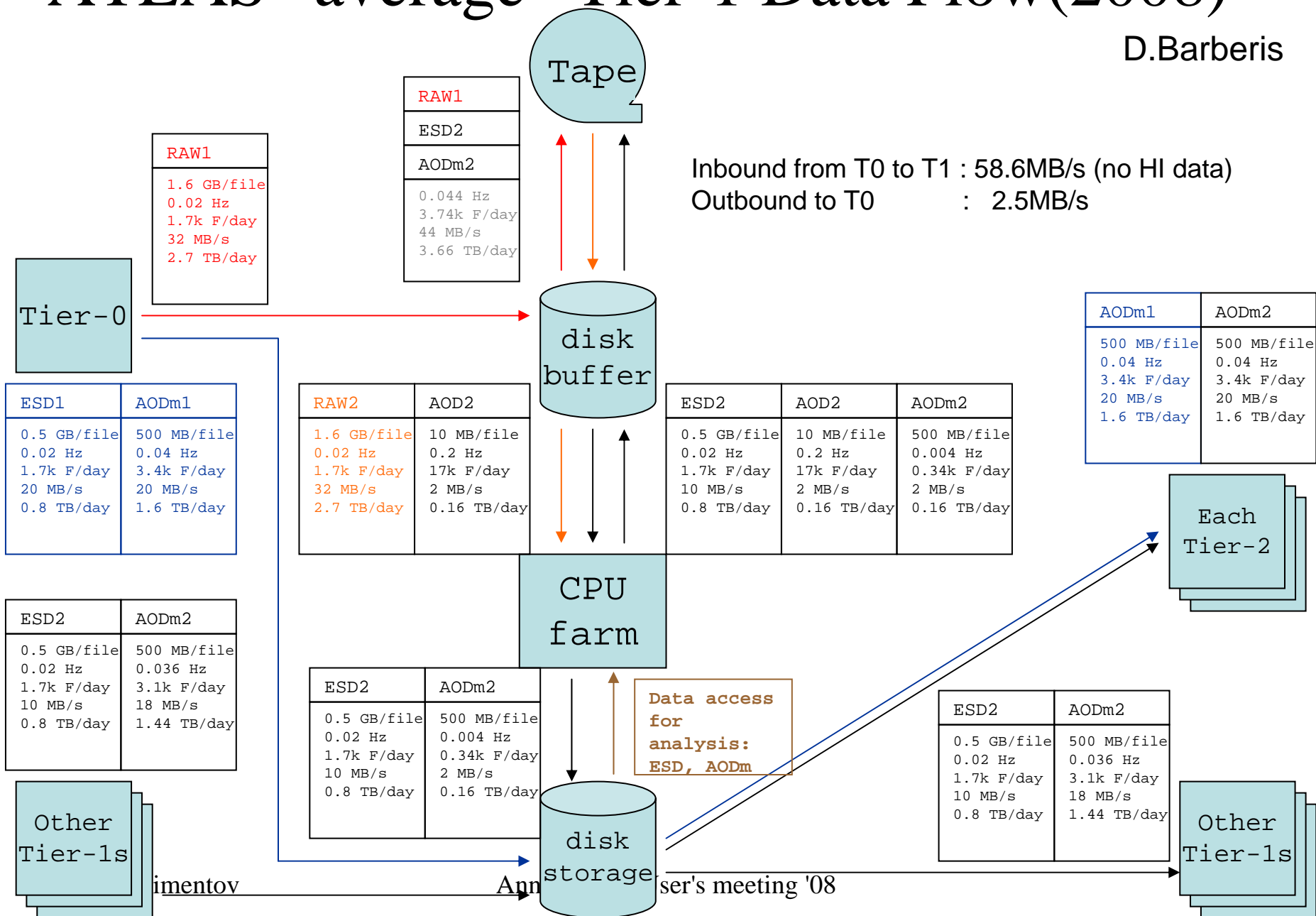- Thanks to my colleagues A.Anisenkov, D.Barberis, M.Branco, S.Campana, A.Farbin and A.Vanyashine for slides and courtesy pictures used in this presentation.

# BACKUP SLIDES

# ATLAS "average" Tier-1 Data Flow(2008)

D.Barberis

Tape

| RAW1 |
|---|
| ESD2 |
| AODm2 |
| 0.044 Hz |
| 3.74k F/day |
| 44 MB/s |
| 3.66 TB/day |

| RAW1 |
|---|
| 1.6 GB/file |
| 0.02 Hz |
| 1.7k F/day |
| 32 MB/s |
| 2.7 TB/day |

Inbound from T0 to T1 : 58.6MB/s (no HI data)
Outbound to T0            :   2.5MB/s

Tier-0

disk buffer

| ESD1 | AODm1 |
|---|---|
| 0.5 GB/file | 500 MB/file |
| 0.02 Hz | 0.04 Hz |
| 1.7k F/day | 3.4k F/day |
| 20 MB/s | 20 MB/s |
| 0.8 TB/day | 1.6 TB/day |

| RAW2 | AOD2 |
|---|---|
| 1.6 GB/file | 10 MB/file |
| 0.02 Hz | 0.2 Hz |
| 1.7k F/day | 17k F/day |
| 32 MB/s | 2 MB/s |
| 2.7 TB/day | 0.16 TB/day |

| ESD2 | AOD2 | AODm2 |
|---|---|---|
| 0.5 GB/file | 10 MB/file | 500 MB/file |
| 0.02 Hz | 0.2 Hz | 0.004 Hz |
| 1.7k F/day | 17k F/day | 0.34k F/day |
| 10 MB/s | 2 MB/s | 2 MB/s |
| 0.8 TB/day | 0.16 TB/day | 0.16 TB/day |

| AODm1 | AODm2 |
|---|---|
| 500 MB/file | 500 MB/file |
| 0.04 Hz | 0.04 Hz |
| 3.4k F/day | 3.4k F/day |
| 20 MB/s | 20 MB/s |
| 1.6 TB/day | 1.6 TB/day |

| ESD2 | AODm2 |
|---|---|
| 0.5 GB/file | 500 MB/file |
| 0.02 Hz | 0.036 Hz |
| 1.7k F/day | 3.1k F/day |
| 10 MB/s | 18 MB/s |
| 0.8 TB/day | 1.44 TB/day |

CPU farm

| ESD2 | AODm2 |
|---|---|
| 0.5 GB/file | 500 MB/file |
| 0.02 Hz | 0.004 Hz |
| 1.7k F/day | 0.34k F/day |
| 10 MB/s | 2 MB/s |
| 0.8 TB/day | 0.16 TB/day |

Each Tier-2

**Data access for analysis: ESD, AODm**

| ESD2 | AODm2 |
|---|---|
| 0.5 GB/file | 500 MB/file |
| 0.02 Hz | 0.036 Hz |
| 1.7k F/day | 3.1k F/day |
| 10 MB/s | 18 MB/s |
| 0.8 TB/day | 1.44 TB/day |

Other Tier-1s

disk storage

Other Tier-1s

imentov                          Ann        ser's meeting '08

# DQ2 Concepts

- 'Dataset':
  - an aggregation of data (spanning more than one physical file!), which are processed together and serve collectively as input or output of a computation or data acquisition process.
  - Flexible definition:
    - … can be used for grouping related data (e.g. RAW from a run with a given luminosity)
    - … can be used for data movement purposes
  - Dataset concept is extended to all ATLAS data (MC, DAQ, DB releases, etc)
- 'File':
  - constituent of a dataset
    - Identified by Logical File Name (LFN) and GUID
- 'Site'
  - A computing site providing storage facilities for ATLAS
    - … which may be a federated site
- 'Subscription'
  - Mechanism to request updates of a dataset to be delivered to a site
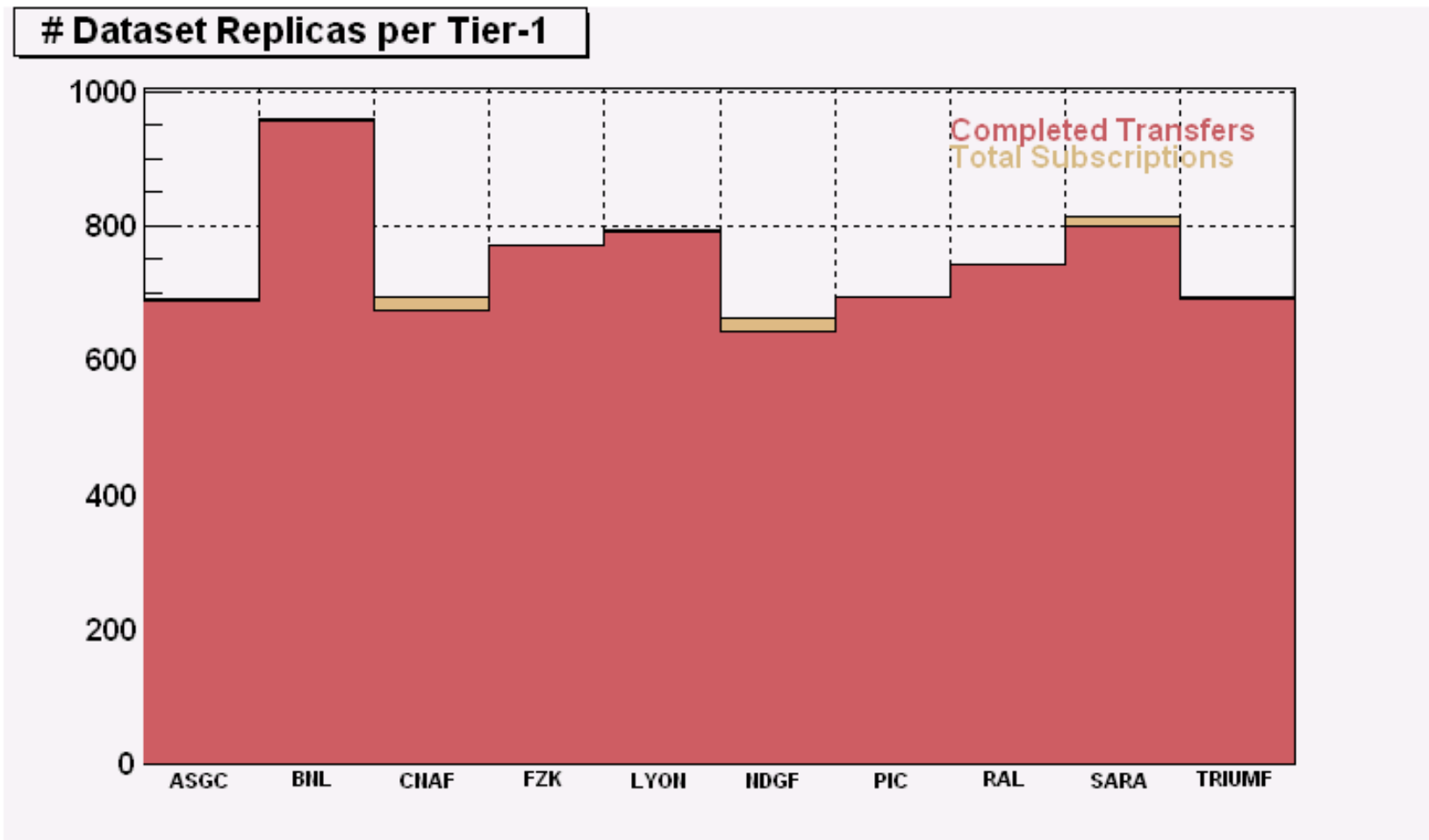
*See M.Lassnig talk CHEP Sep 3, GM2 session*

# DDM/DQ2 more than just s/w development

- DDM forced the introduction of many concepts, defined in the Computing Model, onto the middleware:
  - ATLAS Association between Tier-1/Tier-2s
  - Distinction between temporary (e.g. disk) and archival (e.g. tape) areas
  - Datasets as the unit of data handling
- Not all ATLAS concepts were originally supported by the GRID middleware.
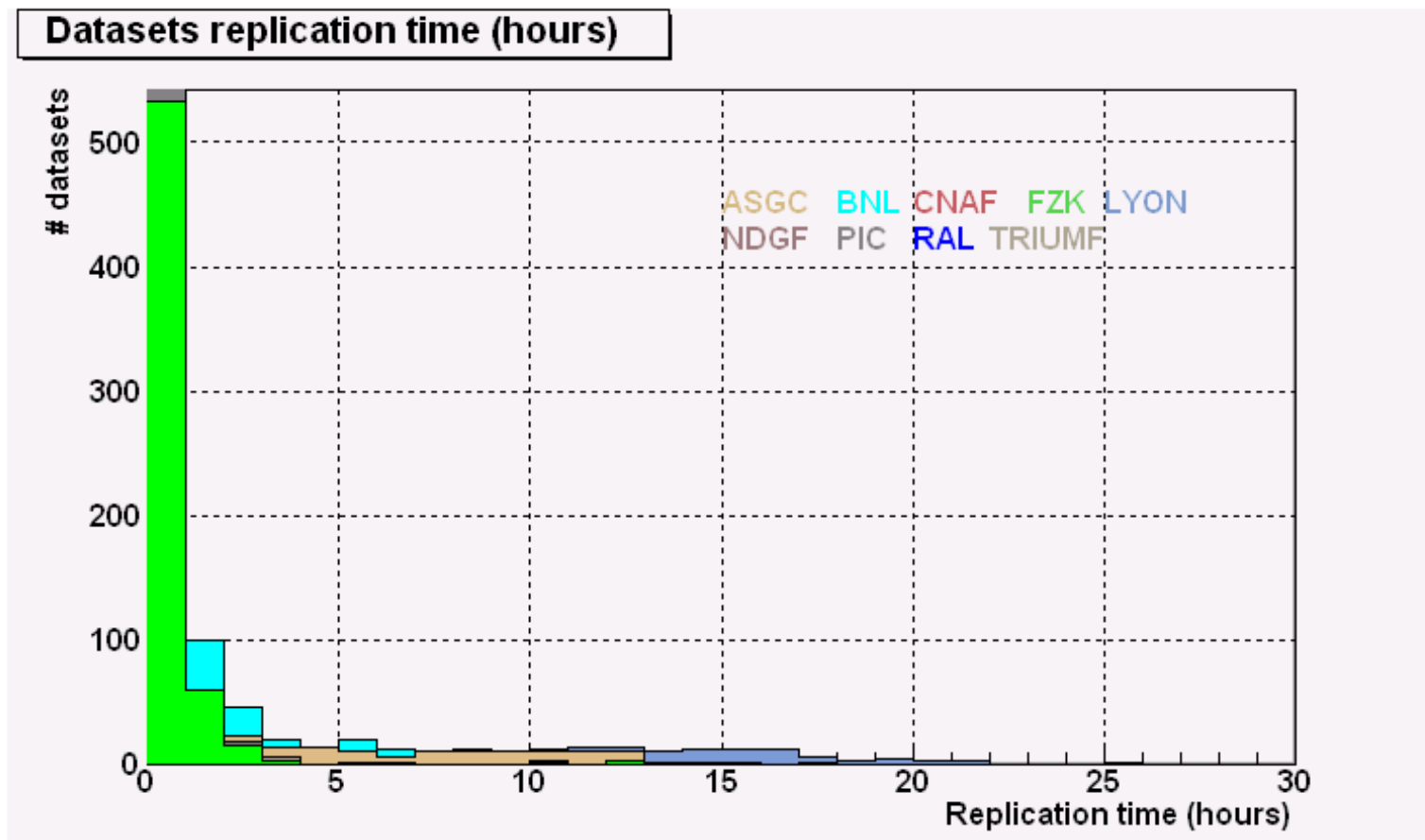
*See M.Lassnig talk CHEP Sep 3, GM2 session*

# Dataset Replicas per Tier-1



Completed Transfers
Total Subscriptions

A.Klimentov, K.Leffhalm Jun 10, 2008

CCRC08 May 2008 (week 19). Replication Processing

Datasets replication time (hours)

A.Klimentov, K.Leffhalm Jun 10, 2008

# DDM Glossary

- 'Dataset':
  - an aggregation of data (spanning more than one physical file!), which are processed together and serve collectively as input or output of a computation or data acquisition process.
  - Flexible definition:
    - … can be used for grouping related data (e.g. RAW from a run with a given luminosity)
    - … can be used for data movement purposes
  - Dataset concept is extended to all ATLAS data (MC, DAQ, DB releases, etc)

  *Duality and many potential and existing problems are coming form it*
    - *Dataset is a unit of replication*
    - *Dataset is a unit of physics data organization*

- 'File':
  - constituent of a dataset
    - Identified by Logical File Name (LFN) and GUID

- 'Site'
  - A computing site providing storage facilities for ATLAS
    - … which may be a federated site

- 'Subscription'
  - Mechanism to request updates of a dataset to be delivered to a site