# Towards end-to-end debugging for data transfers

Gavin McCance
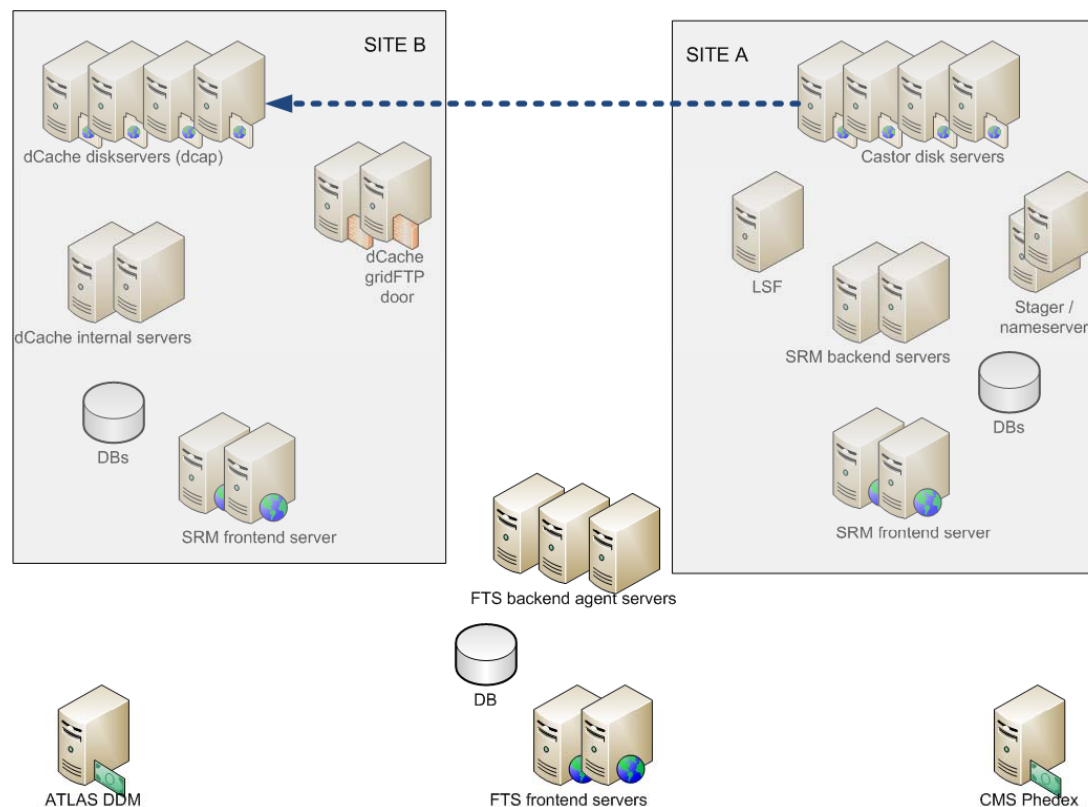
Javier Conejero Banon
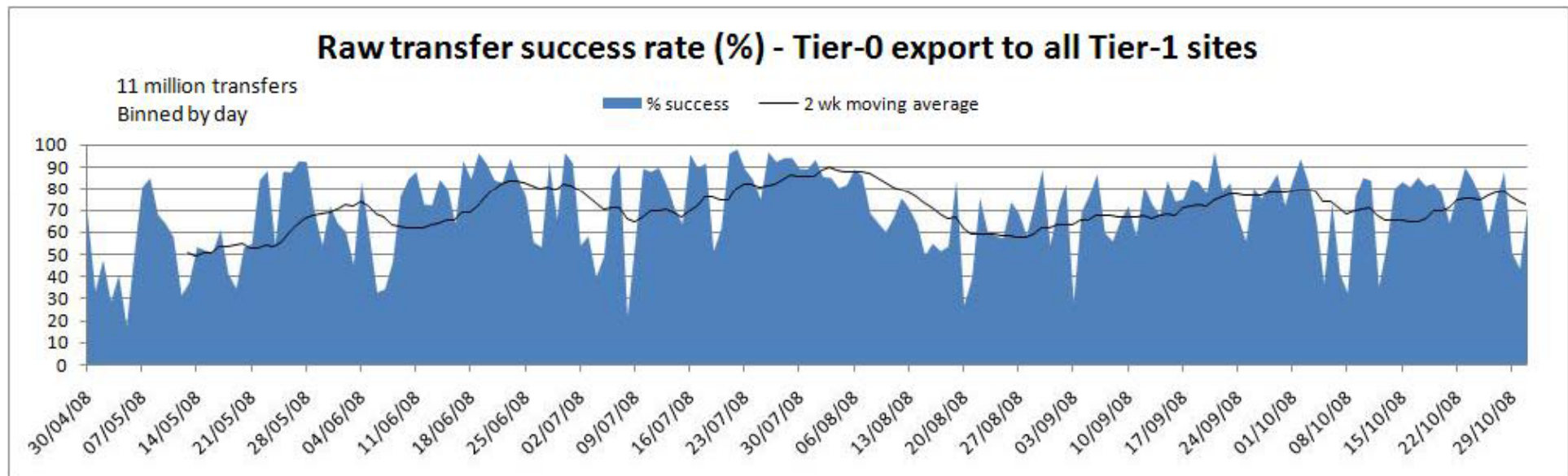
Sophie Lemaitre

CERN IT/FIO

CHEP 2009

- The challenge
- Our problem
- Our solution

- **The challenge**
- Our problem
- Our solution

- The data service model we have in the WLCG is all a bit complex
- There are many layers of software involved
  - Expt framework<->Transfers<->SRMs<->Stagers<->gridFTP<->Tape
- Individual components are quite internally complex

# Efficiency?

- WLCG planning meeting November 2008:
- Examined efficiency of the whole data management stack
  - All the files get there in the end! (multiple retries)
  - RAW transfer rate (#successes / total # attempts, per day)
    - Failure can and do happen in any layer, at both ends of a transfer

- About ¼ of all transfer *attempts* fail due to storage errors ! ☹



Raw transfer success rate (%) - Tier-0 export to all Tier-1 sites

11 million transfers
Binned by day

% success      2 wk moving average

# How can we improve this?

1. **Summary "dashboards"** collect 'events' and provides summary views, sliced in different ways

   – e.g. Current quality on transfers per site

2. **Distributed debug tracing** allows you to follow a specific operation through all the middleware components that process it

   – Show me the transfer for *this* file

- Service operations staff typically use

   – the first one to look for problems

   – second one to drill down and understand them

- The challenge
- **Our problem**
- Our solution

- **Operational cost of distributed debugging is still too high**

- Currently it's grep intensive over all across multiple services
  - wassh –l root –c gridfts grep reqID /var/tmp/*failed/glite*.log
  - wassh -l root -c gridsrm/atlas zgrep "49b4fa78-0000-1000-ad97-fa78af063a57" /var/spool/srm/log.4.gz

- This impacts sites and experiment shifters
  - Training curve is rather steep for new operations staff
  - Inaccessible files (~hours)
  - gridFTP mysteriously timing out (bouncing emails/phone calls back a couple of times to the other site) (~hours)
  - "We see reduced transfer rates, please could you check" (~hours)
  - Performance variance is typically very large and not well understood
    - Some files transfer at 10MB/s, some go a 200KB/s, same site, same time

- **Better debug tools can reduce operations cost!**

1. A support ticket comes in
   - "We see lots of transfers timing out"
   - Example file:
     - /castor/cern.ch/grid/atlas/atlasdatadisk/data08_cosmag/ESD/data08_cosmag.00090272.physics_RPCwBeam.reco. ESD.o4_r560_tid027478/ESD.027478._00769.pool.root.1"

2. Submit request to debug transfer for this file

3. Picture will be built up asynchronously as data is returned from the various sources, like a web-page loading

Fabric Infrastructure and Operations

CERN IT Department

Network utilization

- eth0 in   aver:280.9M   max:1.1G   min:47.4k   curr:374.6M
- eth0 out  aver:666.0M   max:1.8G   min:41.4k   curr:1.0G

transfer summary:
ROD, RAL-T1, atlas, "gridFTP: the server timed out"

Trace detail:
pareToGet -> CERN: detail
srmGetStatusOfGet -> CERN: detail
Srm gsiftp returned CERN: gsiftp://lxfsrc0203.c

**SRM service**: received call     **Stager:**  request schedule o
Scheduled on stager
TURL determined                   **LSF scheduler:**  diskserve

CPU utilization

- User     aver:225.8m   max:409.1m   min:173.6m   curr:196.0m
- System   aver:1.1      max:4.2      min:50.0m    curr:51.8m
- Nice     aver:111.8m   max:467.5m   min:0.0      curr:0.0
- Idle     aver:92.4     max:99.7     min:63.0     curr:99.6
- IO Wait  aver:5.2      max:34.9     min:87.5m    curr:113.5m
- IRQ      aver:38.0m    max:116.3m   min:0.0      curr:0.0
- Soft IRQ aver:901.0m   max:3.9      min:0.0      curr:0.0

srmPrepareToPut -> RAL: detail
srmGetStatusOfGet -> RAL: detail
Srm gsiftp returned RAL: gsiftp://dispool0023.r

**SRM service**: received call   **Stager:**  request schedule of job of diskserver
Scheduled on stager
TURL determined                 **Scheduler:**  diskserver access scheduled

gridFTP 3rd party call:
CERN -> RAL: detail
:
**GridFTP RAL: FTS client connect**
**Opening data connection to other side on port X**
**Timeout!**
:

# Integration problem

- It's an integration problem



- Multiple logfile / database / feed formats to be parsed
- Logs located on multiple machines (O(1000) nodes @CERN)

- The challenge
- Our problem
- **Our solution**

- Our previous attempts focused on recording *all events*
  - You collect all events from all sources, all the time, parse them, and put them in an index database
    - Specific debug searches can be run over the database
  - Approach used by Splunk
  - Common logging instrumentation: approach taken by netlogger

- While this does work, routine parsing, collecting and joining can be expensive
  - Parsing 10's GB's of logs from O(1000) machines
  - It's overkill for this application
    - A typical service manager will probably run no more than 0(100) debug trace queries a day, and we know what queries will be run

- We prefer to parse on demand
  - Can make use of debug trace databases if they are available

# Our approach

- **On-demand extraction from data sources (request / response)**
  - Send out requests to all data sources that might know something, get them to parse and return what they know
  - If sufficiently detailed summary or trace logging databases are available, use them
  - Integrate other feeds (fabric monitoring, network monitoring data)

- **Integrate** (join) the data from the various sources for that specific debug request
  - The flow is asynchronous, i.e. the picture of what happened is built up as information is returned
  - Even with missing information, the picture obtained is still useful for debugging

- Based on message-oriented middleware
- This handles the request / response reliably and easily

Publish-Subscribe / Broadcast

Point-to-point



- Send a query to all nodes that might know something
- e.g. all SRM nodes in a load-balanced alias

- Send a query to the one diskserver that we know handled the transfer
- e.g. gridFTP logs

Broadcast to all SRM servers that might know something

FTS service DB transfer overview

SRM servers x20

Determine queries to ask
Send out queries
Gather responses
Integrate
Present result

Admin requests debug of a file transfer

Point to point message to the one disk-server that does know something

Disk-servers (gridFTP) x1000

Response messages returned for integration and presentation

- The message system handles the plumbing and the reliable delivery of the messages – local agents do the parsing

# Messaging technology

- **Using the MSG messaging framework**
  - Technology already used in WLCG production in EGEE/OSG for grid site-monitoring data
  - See EGEE User Forum for details of MSG:
  - http://indico.cern.ch/contributionDisplay.py?contribId=136&sessionId=9&confId=40435
  - Uses Apache ActiveMQ: open source, easy to use

- ✓ Throughput requirements
  - Isn't really an issue for administrator initiated requests: O(100) per day
- ✓ Latency requirements
  - Needs to deliver fast – we don't want to be waiting too long
- ✓ Reliability requirements
  - We do care that the messages get there in order to build up a full picture
- ✓ Scaling requirements
  - We need it to scale up to O(1000) nodes so that we can run this over all our diskservers

- Planning to integrate all data at CERN from:
  - File Transfer Service (Tier-0 physics data export service)
  - Castor SRM and Castor core components
  - Lemon fabric monitoring service

- Aim: tool usable by service managers in summer to help with the transfer debugging problem

- Future:
  - Add data feeds from other sites (other SRMs): collaboration with external sites. Add network monitoring data
  - Tool itself useful for other sites?
  - Re-use components for distributed workload-management services?

Fabric
Infrastructure
and
Operations

CERN
IT
Department

# Flexible architecture

- Future re-plumbing is easy: the architecture allows us to easily change the data sources as software develops
  - *Decide we want to collect and archive gridFTP logs on 10 central machines*
    - Move the gridFTP agents off all your diskservers to just these 10 machines instead, to answer the same request
    - The rest of the system remains unchanged
  - *Next version of one component comes with a detailed-enough trace database?*
    - Unplug all the log-mining agents and plug on an agent to answer the same request from the trace database instead
    - The rest of the system remains unchanged
  - *Want to add in network flow data ?*
    - Write another feed to make this data available and add it in

# Summary

✓ **Aim: to reduce operations cost of running complex distributed services**

● Developing a flexible architecture based on messaging for trace-debugging of distributed services

  – Parse logs as data sources

  – Use trace database sources if available

● Integrate data **on-demand** from various sources instead of routine parsing

✓ **Will have a usable tool to help with the transfer debugging problem by summer**

# Backup

- Using common formats and even better a common logging trace schema for all components involved is a great idea!

- Easier to do if you control all the components
  - e.g. most components of Castor drop trace info into a distributed tracing component (DLF database)
  - *Netlogger* calls can be added to the code to send data streams out
- Hard for other components
  - Some bits of the code we don't 'own' (Castor: *LSF, gridFTP*), so it can be hard to add trace info at the level needed
  - Why should FTS, dCache, Lemon, Nagios log into the same format?

- While this is a good goal we prefer to deal with the integration problem we have directly