# Any Data, Anytime, Anywhere

Dan Bradley <dan@hep.wisc.edu>
representing the AAA Team
At OSG All Hands Meeting
March 2013, Indianapolis

# AAA Project

## Goal

- Use resources more effectively through remote data access in CMS

## Sub-goals

- Low-ceremony/latency access to any single event

- Reduce data access error rate

- Overflow jobs from busy sites to less busy ones

- Use opportunistic resources

- Make life at T3s easier

# xrootd: Federating Storage Systems

- Step 1: deploy seamless global storage interface

- <span style="color:orange">But preserve site autonomy</span>:

  - xrootd plugin maps from global logical filename to physical filename at site

    - Mapping is typically trivial in CMS:
      /store/* → /store/*

  - xrootd plugin reads from site storage system

    - Example: HDFS

  - User authentication also pluggable

    - But we use standard GSI + lcmaps + GUMS
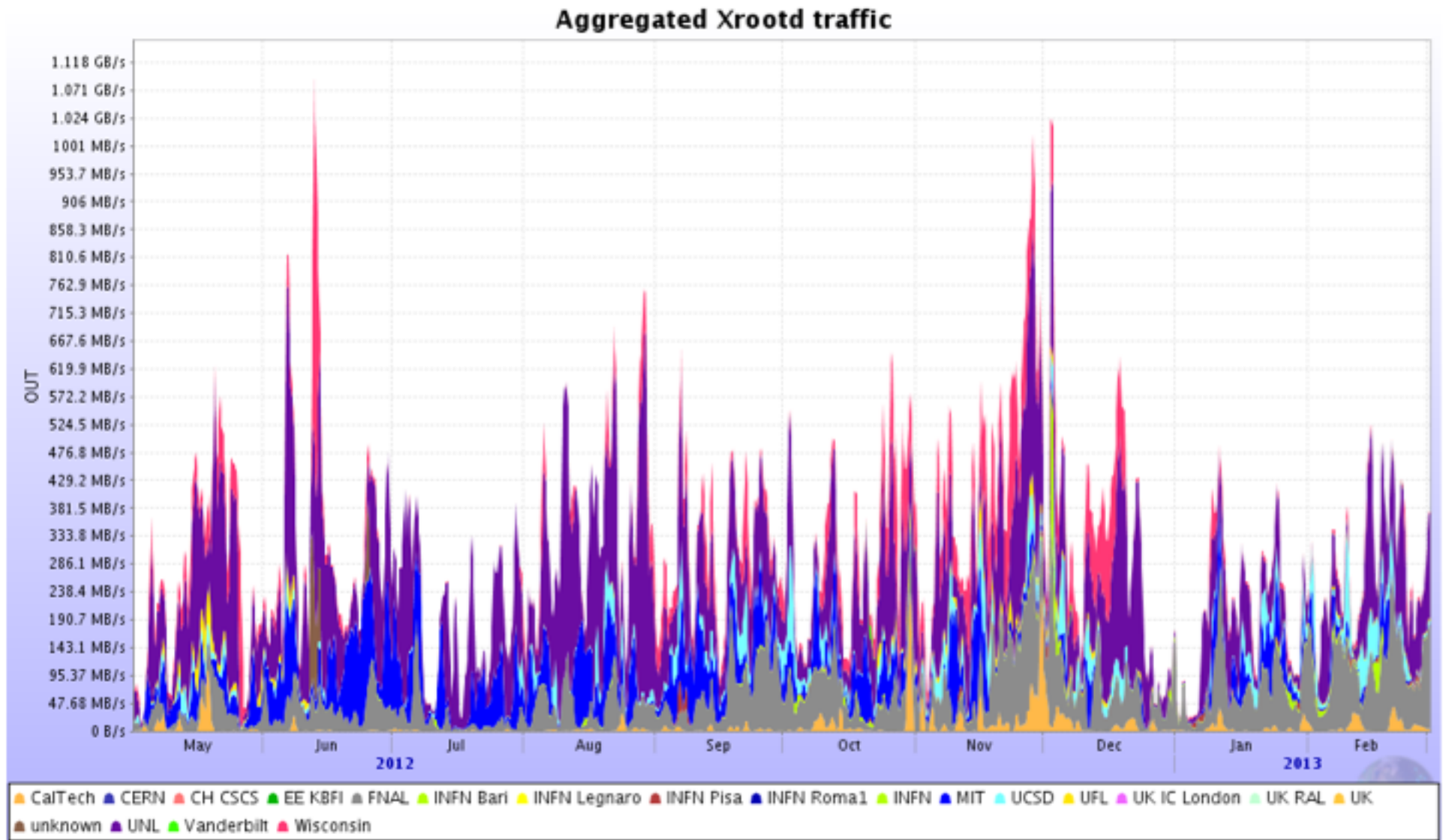
# Status of CMS Federation

US

- T1 (disk) + 7/7 T2s federated

- Covers 100% of the data for analysis

- Does not cover files only on tape

World

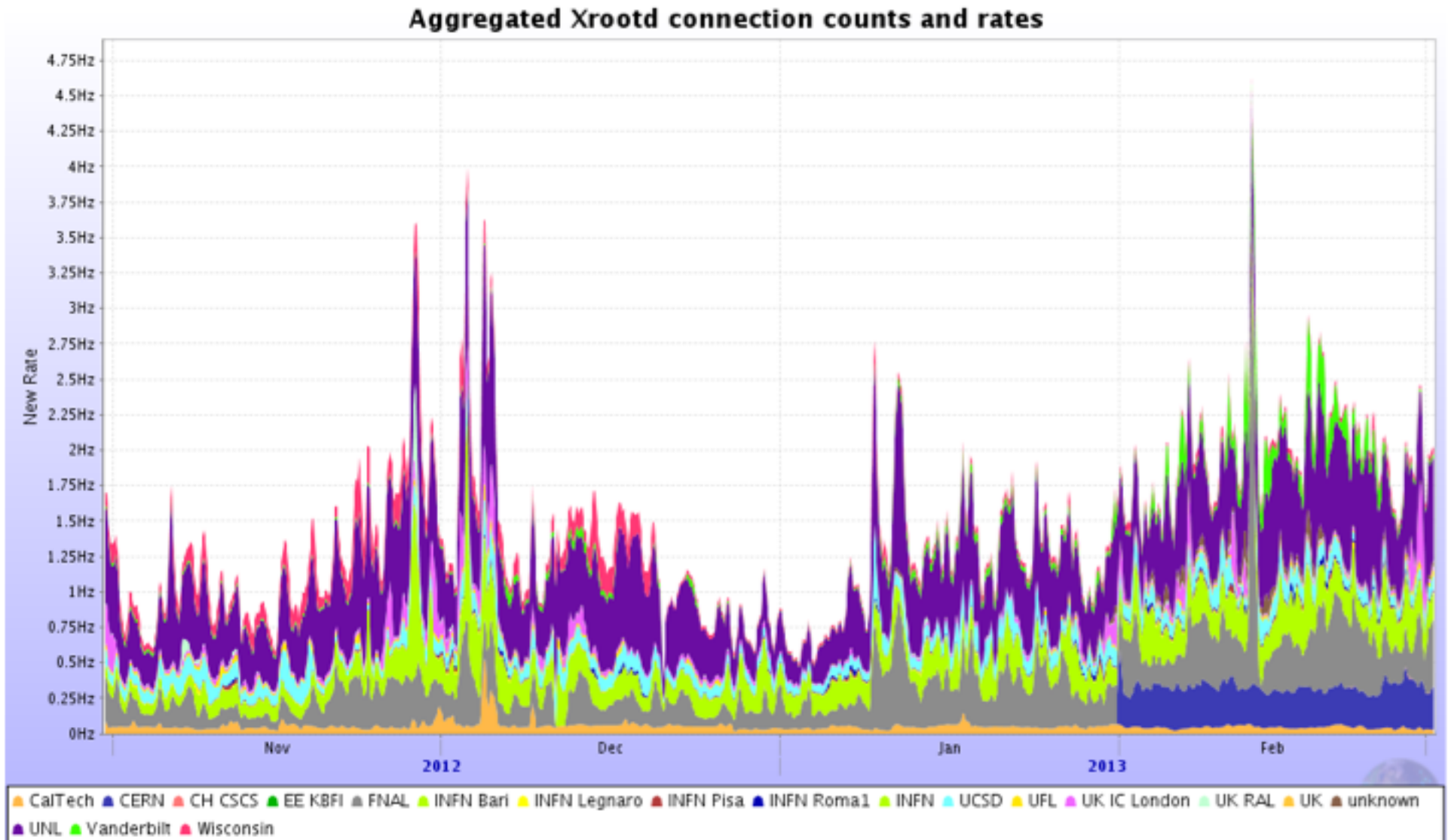- 2 T1s + 1/3 T2s accessible

- Monitored but not a "turns your site red" service (yet)

# WAN xrootd traffic



**Aggregated Xrootd traffic**

CalTech ■ CERN ■ CH CSCS ■ EE KBFI ■ FNAL ■ INFN Bari ■ INFN Legnaro ■ INFN Pisa ■ INFN Roma1 ■ INFN ■ MIT ■ UCSD ■ UFL ■ UK IC London ■ UK RAL ■ UK
unknown ■ UNL ■ Vanderbilt ■ Wisconsin

# Opening Files



Aggregated Xrootd connection counts and rates

# Microscopic View

# Problem

- Access via xrootd overloads site storage system

  - Florida to Federation, "We are seceding!"

- Terms of the Feb 2013 treaty:

  - Addition of local xrootd I/O load monitoring

  - Site can configure automatic throttles

    – When load too high, rejects new transfer requests

    – End-user only sees error if file unavailable elsewhere in federation

- But these policies are intended for the exception, not the norm, because ...

# Regulation of Requests

- To 1$^{st}$ order, jobs still run at sites with the data

  - ~0.25 GB/s average remote read rate

  - O(10) GB/s average local read rate

  - ~1.5 GB/s PhEDEx transfer rate

- Cases where data is read remotely:

  - Interactive          - limited by # humans

  - Fallback             - limited by error rate opening files

  - Overflow             - limited by scheduling policy

  - Opportunistic        - limited by scheduling policy

  - T3                   - watching this

# At the Campus Scale

Some sites are using xrootd for access to data from across a campus grid

- Examples: Nebraska, Purdue, Wisconsin



Any data, Anytime, Anywhere

# More on Fallback

- On file open error, CMS software can retry via alternate location/protocol
  - Configured by site admin
  - We fall back to regional xrootd federation
    - US, EU
    - Could also have inter-region fallback
      - Have not configured this … yet
- Can recover from missing file error, but not missing block within file error (more on this later)
- Has more uses than just error recovery ...

# More about Overflow

- GlideinWMS scheduling policy
  - Candidates for overflow:
    - Idle jobs with wait time above threshold (6h)
    - Desired data available in a region supporting overflow
  - Regulation of overflow:
    - Limited number of overflow glideins submitted per source site
- Data access
  - No reconfiguration of job required
    - Uses fallback mechanism
    - Try local access, fall back to remote access on failure

# Overflow

- Small but steady overflow in US region

# Running Opportunistically

- To run CMS jobs at non-CMS sites, we need

  - Outbound network access

  - Access to CMS datafiles
    - Xrootd remote access

  - Access to conditions data
    - http proxy

  - Access to CMS software
    - CVMFS (also needs http proxy)

# CVMFS Anywhere

But non-CMS sites might not happen to mount the CMS CVMFS repository

→ Run the job under Parrot (from cctools)

- Can now access CVMFS without FUSE mount
- Also gives us identity boxing
  - Privilege separation between glidein and user job
- Has worked well for guinea pig analysis users
  - Working on extending it to more users

What about in the cloud?

- If you control the VM image, just mount CVMFS

# Fallback++

- Today we can recover when file is missing from local storage system

- But missing blocks within files cause jobs to fail

  - And job may come back and fail again ...

  - Admin may need to intervene to recover the data

  - User may need to resubmit the job

- Can we do better?

# Yes, We Hope

- Concept
  - Fall back on read error
  - Cache remotely read data
  - Insert downloaded data back into storage system

# File Healing



XRD REDIRECTOR

Some Other Site

3.2

3.1

3.3

Xrd Site Master
(not relevant)

Xrd Caching Proxy

Site Storage

Disk Buffer

4.2 Store to local disk buffer

2.

4.1 Serve data to
the original job

1. Local acces
fails

4.x And to any
other client

5. Move to Site Storage

Any data, Anytime, Anywhere

18

# File Healing Status

- Currently have it working via whole-file caching

  - Still only triggered by file open error

- Plans to support partial-file healing

  - Will need to fall back to local xrootd proxy on all read failures

  - Current implementation is HDFS-specific

    – Modifies HDFS client to do the fallback to xrootd

    – But it's not CMS-specific

# Cross-site Replication

- Once we have partial-file healing …

  - Could reduce HDFS replication level from 2 to 1 and use cross-site redundancy instead

    – Would need to enforce the replication policy at higher level

    – May not be good idea for hot data

    – Need to consider impact on performance

# Performance

Mostly CMS application-specific stuff

- Improved remote read performance by combining multiple reads into vector reads

  – Eliminates many round-trips

- Working on bit-torrent-like capabilities in CMS application

  – Read from multiple xrootd sources

  – Balance load away from slower source

  – React in O(1) minute time frame

# HTCondor Integration

- Improved vanilla universe file transfer scheduling and monitoring in 7.9

  - Used to have one file transfer queue

    - One misconfigured workflow could starve everything else
    - Difficult to diagnose

  - Now one per user

    - Or per arbitrary attribute (e.g. target site)

  - Equal sharing between transfer queues in case of contention

  - Reporting transfer status, bandwidth usage, disk load, and network load

  - And now you can condor_rm those malformed jobs that are transferring GBs of files :)

# Summary

- xrootd storage federation rapidly expanding and proving useful within CMS

- We hope to do more

  - Automatic error recovery

  - Opportunistic usage

  - Improving performance