

# CMS Input for Data Management TEG

Brian Bockelman

# Forewarning

- Given the number of questions, amount of time to collect input, and the breadth of the collaboration.
  - ... this input isn't written in gold,
  - ... this input hasn't been signed off by every soul in CMS computing,
  - ... but I think reflects the accurate status of CMS data management and contains some crystal ball gazing.
- I hope that, by volunteering to go first, I get the chance to do a follow-up for all the things I royally mess-up.

# Background

- CMS's Data Management is underpinned by the dataset placement service (PhEDEx) and the dataset bookkeeping service (DBS).
  - PhEDEx handles the interaction with the transfer layer, block location information, and dataset/block subscription.
    - Notice that file location is only tracked for in-flight blocks.
  - DBS handles physics metadata information.
  - PhEDEx has no knowledge of physics; DBS has no knowledge of

# Background (Philosophy)

- CMS attempts to not micro-manage sites.  
Examples:
  - Site may use any space management technique they'd like.
  - No quota systems required. Sites "data managers" must approve all data transfers before they happen.
  - This has downsides too – we require active, engaged sites. Can't really "plug it in and walk away", although the operational cost drops considerably after initial setup.
- We feel the existing system is functional, but is overly costly in terms of complexity and cost.

# On file catalogs

- LFC is not used.
  - Our file catalog has about 300k HTTP requests / day with 600ms average response time.
  - Total request count is far reduced by the fact our transfer system uses direct Oracle access.
- Observation: there has been lots of effort to collaborate on location catalogs, but none on physics catalogs. Why?
  - Have we been burned too many times here?

# On Namespaces

- CMS has a global namespace. Each file has one, unique LFN.
  - Sites provide a set of mapping rules; these are used when jobs start, or right before the transfer is sent to FTS, to translate to a PFN.
    - It is a touch tricky communicating between sites – the remote site needs to a priori tell you the PFN, or you need to find its mapping rules.
  - CMS has no concept of a SURL, TURL, GUID, etc. Just LFN and PFN.
    - We do not track SURLs in our catalogs. The blocks we track are associated with an opaque string (usually a hostname) associated with the site.
- Observation: In our xrootd work, we've found it very convenient to have the WAN protocol speak "LFNs" instead of "PFNs". This way, your knowledge of the remote site only needs to be the endpoint.
  - It would be nice if such a thing crossed over to other protocols. Why does SRM or GridFTP need to expose site internals?

# CMS Site Storage Requirements

- Stripped to the minimum:
  - Site must support wide area transfers via FTS.
  - Site must support one of the stock ROOT local protocols *or* provide an implementation in CMSSW.
  - Site must use one of our plugins for file stageout *or* provide one of their own.
- Note: a T2 site with a powerful gridftp server and a cluster file system would probably be just fine.

# Space Management

- CMS delegates space management to the sites. They make the decision of whether to take on more data, and police free space.
  - If the best way to do that is via SRM space tokens, sites can do that. In practice, zero-to-little use of SRM for space management.
  - CMS provides tools to discover discrepancies with the catalog and perform checksums.
    - Moderate-to-good uptake at the T1s.
    - Relatively little uptake at the T2s. Perhaps this isn't necessary at T2s if data federations are widely available?



# Authorization Management

- Authorization management is handled by the sites.
  - We provide guidance on what roles/users should be able to access what, and the sites implement it.
    - And bug the sites until they are mostly in compliance.
  - Our DM tools prevent users from causing problems – *if* they are using our tools.
  - No strong, verifiable means to know that a site is protected against a malicious user (a-la ALICE security tokens).
  - Number one concern is keeping CMS data private from other organizations/ the world.
- Observation: we would benefit if intra-VO authorizations could be taken from VOMS, even if only for grid-based write mechanisms (SRM/GridFTP).
  - It would probably simplify things for sites, as the storage would need to be configured at the VO-level, not the sub-VO-level.
  - This is not a priority, however.

# The SRM and GridFTP Question

- CMS's heaviest use of SRM is to load-balance GridFTP.
  - We can work with space tokens for the sites that desire SRM space management.
- Can CMS live without SRM? Yes!
  - However, experience with Bestman2 shows SRM can be made “small enough” to not be a problem
  - Replacing SRM would mean an appropriate replacement needs to be found...

# HTTP vs GridFTP

- HTTP and GridFTP can both load-balance themselves, but may require more network expertise at sites.
- Many have observed that GridFTP can be replaced by HTTP.
  - Just need to figure out third-party-copy.
  - Oh, and delegation.
  - Oh, and standardize a mechanism for integrity checks without encryption
  - At this point, is it a “standardized protocol” or “building a new SRM”?
- This would be an interesting investigation (HTTP is *much* more flexible and widespread than GridFTP), but aren't ready to jump ship yet.

# The HSM Question

- AKA, should we hard partition archive and disk so HSM is no longer needed.
- Our data system can handle this, with small adaption.
  - Definitely need the ability to interact with archives – will need to do organized recalls, and need to know when a file is on tape.
- We feel this should be driven by the needs of the sites.
  - However, if this is a way to increase reliability and simplify site layout, we support it.

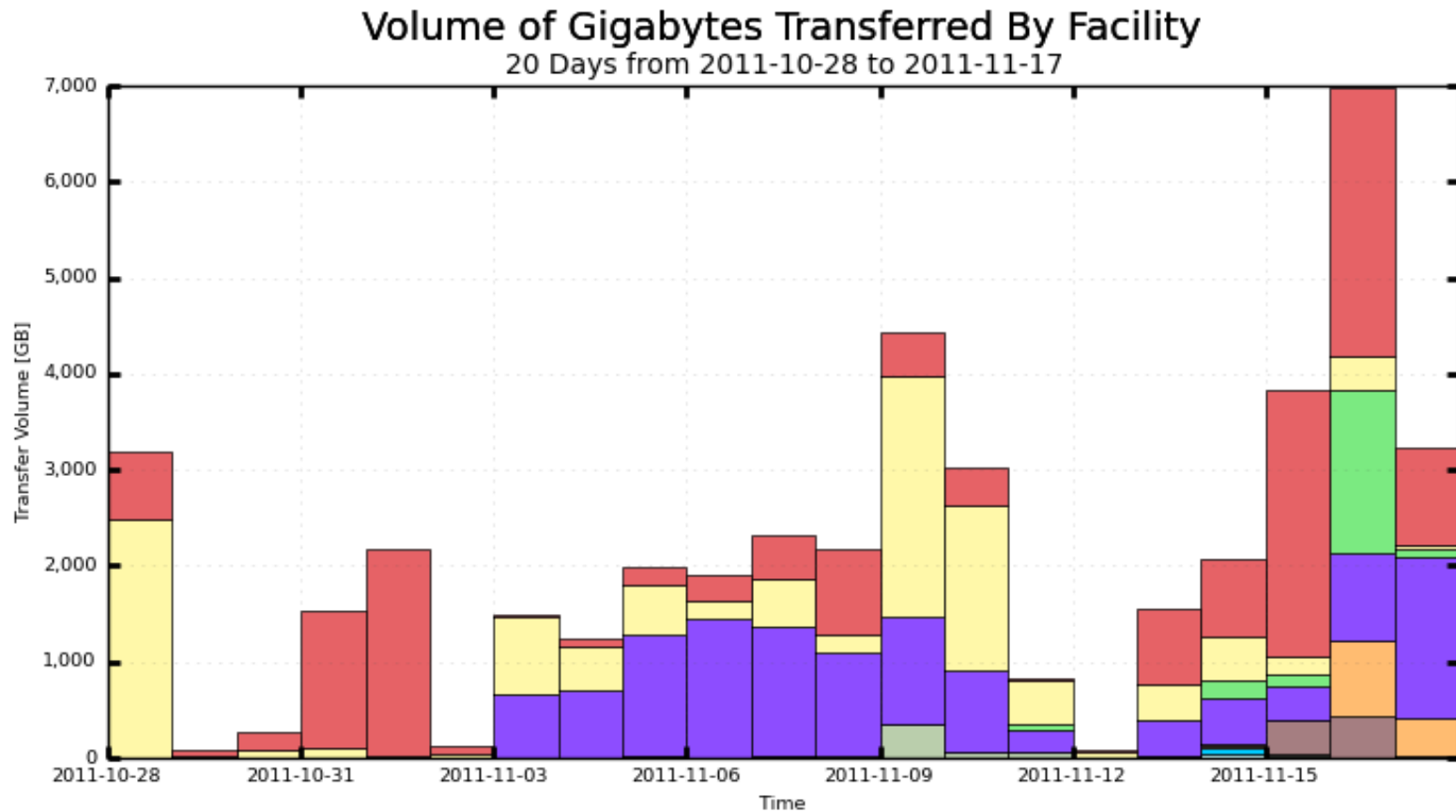
# Data Federations

- What is a data federation?
  - Hopefully to be better defined at the related workshop next week!
  - Loosely, it's the ability for users to access data across multiple source sites uniformly from a single endpoint.
    - Note that access implies reads, not writes.
    - Some differences between sites are hidden by the federation; user never has to do a location lookup, know local namespaces, or know site endpoints a priori.

# Data Federations in CMS

- Each “network region” in CMS is building a data federation.
  - Best coverage is in US, where all T2s and T1 are participating.
- Provides:
  - Fallback access: if files are missing or damaged, applications transparently failover.
  - Interactive access: a mechanism for users to directly interact with the storage over WAN for small tasks (purposely has scaling limitations).
  - Overflow access: if extra CPU is available, send jobs there regardless of whether the necessary data is present (a small optimization on our data distribution).
  - Data for T3s. A slightly larger-scale “interactive access”.

# Teaser Slide



Daily monitoring from our data federation – each bar represents the activities of one source site.

# Data Federations and Caching

- You can have a federation without a cache, and a cache without a federation.
- We do not currently see a future where the whole experiment is cache based.
  - Therefore, any use of caching will be alongside the existing system. Does this cut complexity?
  - Possibly useful in limited form – ie, T3s and data sharing.
- It's possible to recast the current system as more “cache-like”.
  - For example, our systems assume that when they place data somewhere, it stays there indefinitely.
  - Not a safe assumption for T3s (unfortunately not safe for some T2s also!).
  - Can we take the current system, and start asking ourselves what needs to be changed if “cache eviction” was a normal, not exceptional, event.



# Complexity Kills

- From the MB: Where should the complexity and intelligence lie - the experiments or the infrastructure? How do you view the balance now, and how would you like to see this change (if at all)?
  - Our gut is that we want to see sites be “simple and reliable”, and have the complexity managed by the experiments.
  - However – it strikes us that this opinion does not do much to cut “complexity and cost”.
- Not sure how to resolve this, hopefully the TEG has some ideas!