# PhEDEx Data Service

**Ricky Egeland, University of Minnesota**
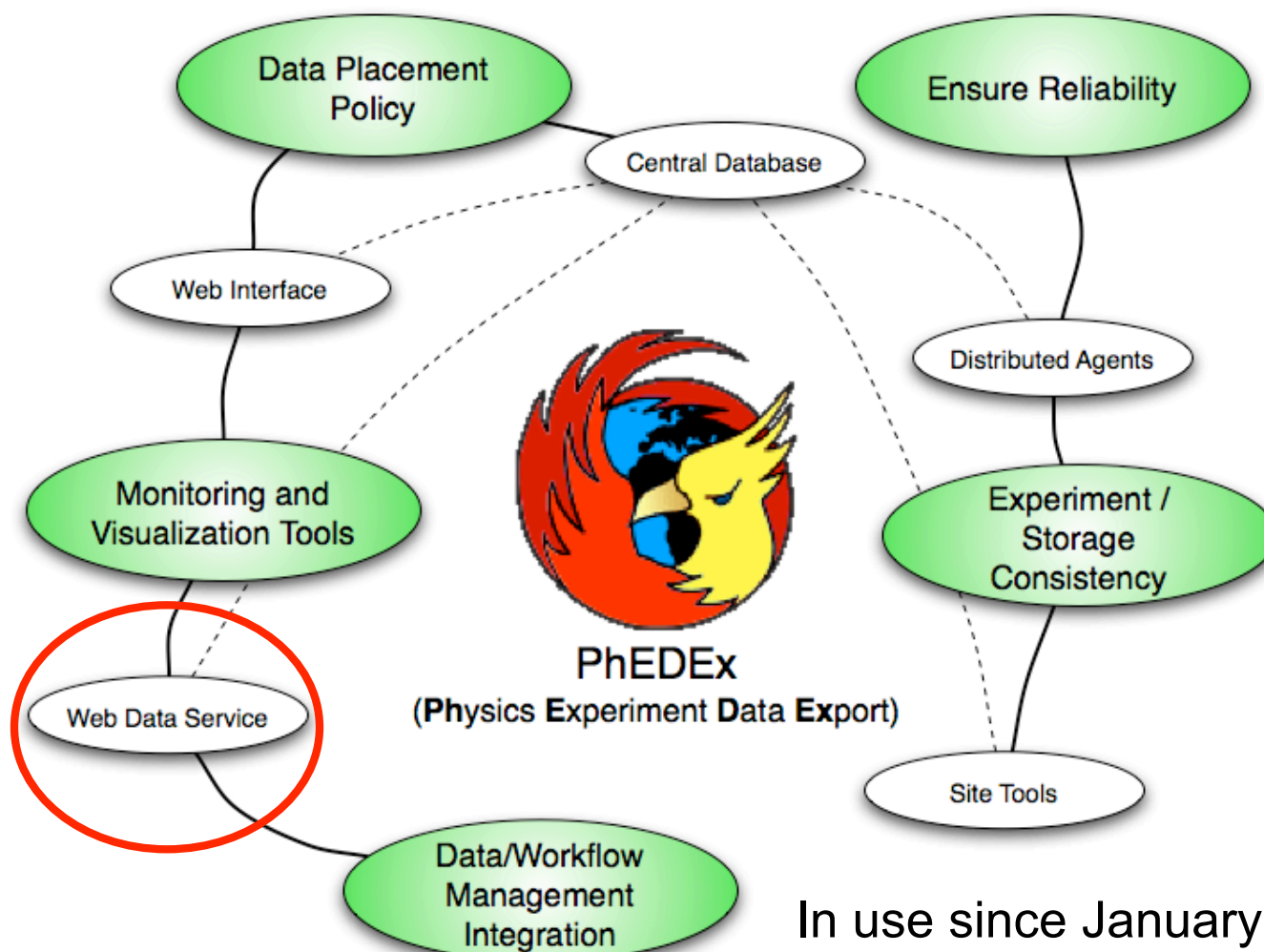
on behalf of numerous PhEDEx contributors
and the CMS collaboration

presented at
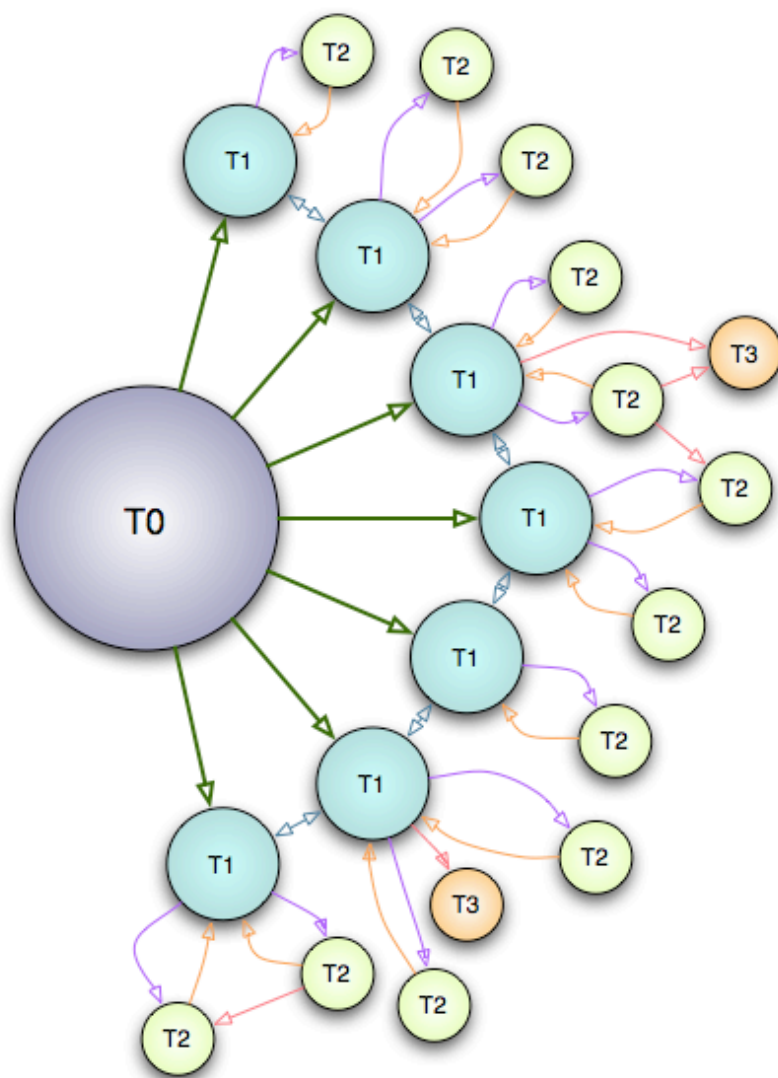**C**omputing in **H**igh **E**nergy **P**hysics (**CHEP**)
March 26th, 2009

In use since January 2004
Over 67 PB of data transfers

2008 Average Per-Link Requirements

T0 ——→ T1
100 MB/s

T1 ←——→ T1
140 MB/s

T1 ——→ T2
25 MB/s

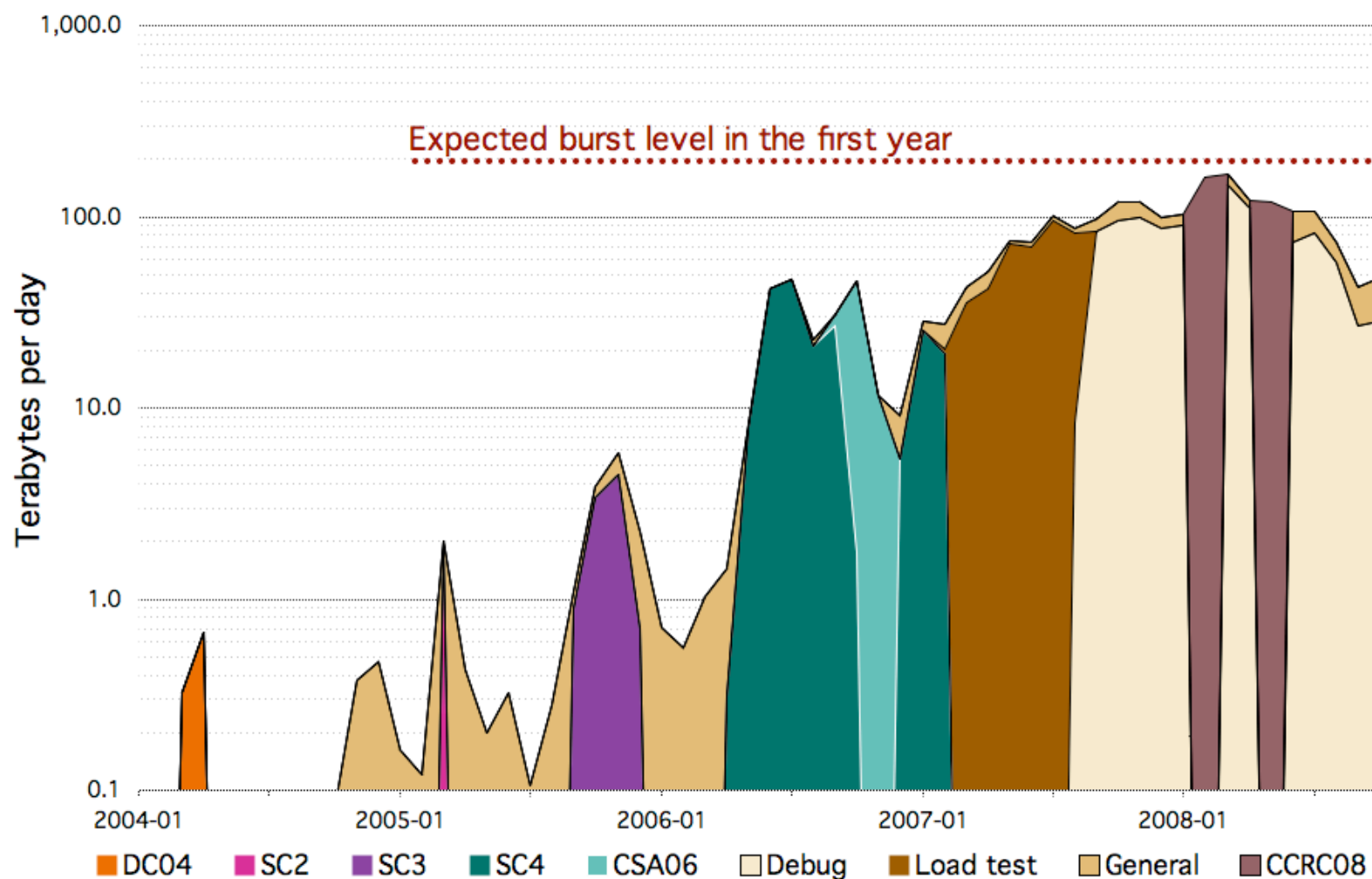T1 ←—— T2
6 MB/s

T2 ←——→ T2

T1 ——→ T3

T2 ——→ T3
undefined, small

production database: 110 nodes, 4.3 M files, 11.7 M replicas

Average data transfer volume

# Data available in PhEDEx

- Datasets known to the transfer system
  - datasets / blocks / files
  - filesize, checksum
- Locations of datasets
  - (node, block) => (node, file)
- Requested data transfers
  - who, what, where, when, why
- Current / Historical transfer statistics
  - actively transferring data
  - rate, quality history
- Current / Historical site usage statistics
  - resident data size, requested data size
  - by group or "custodial"
- Monitoring
  - Recent transfer errors
  - Consistency check tests
  - System monitoring (Agent uptime, status)
- System structure
  - Transfer topology
  - Node list

- ## Integration with other data management components
  - As painless as possible: only requirement is an HTTP client
  - PhEDEx became the "authorative source" for data location
  - Automatic dataset injection, subscription from production components. **The data service is not read-only.**

- ## Provide monitoring data to custom user scripts
  - As painless as possible: only requirement is an HTTP client
  - No database passwords to distribute

- ## Integrate transfer data with other monitoring services

- ## Provide data for next-generation website

- ## Maximize code re-use within our project

https://cmsweb.cern.ch/phedex/datasvc/format/instance/api?options

https://cmsweb.cern.ch/phedex/datasvc/doc/api
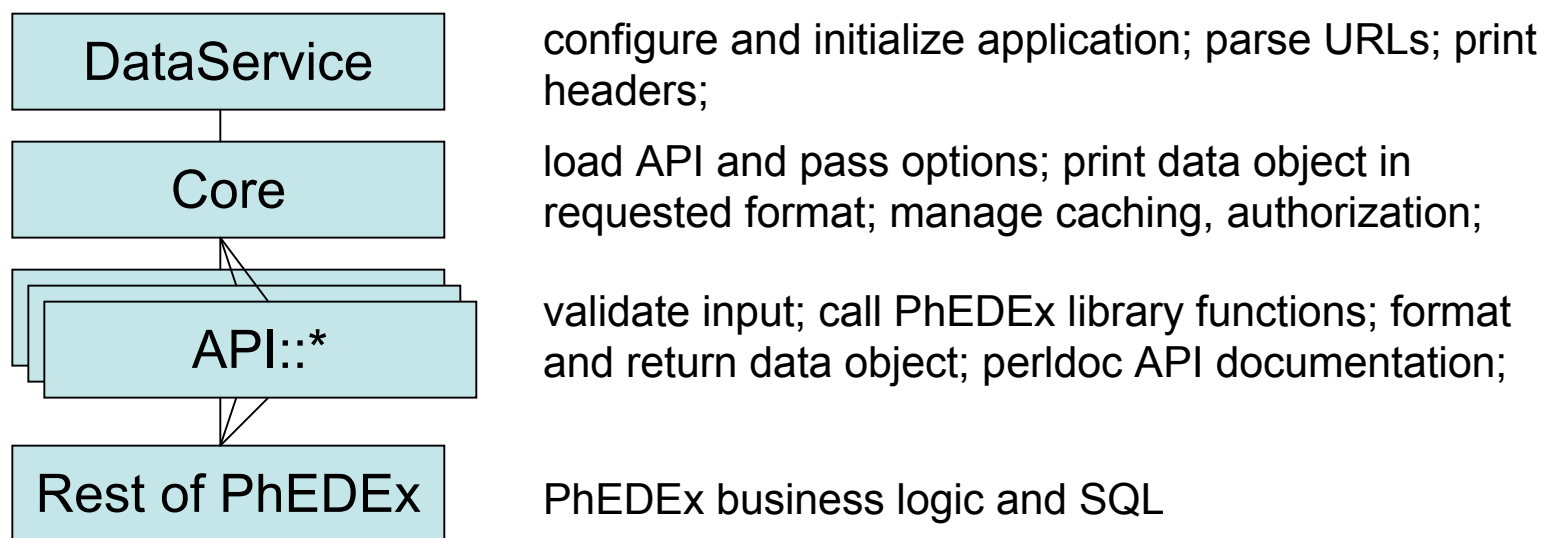
format : output format – xml, json, perl

instance : database instance – prod, debug, test

api : API to use – blockreplicas, subscribe

options : API options (mostly filters) – block=/X/Y/Z#123

doc : output documentation

| | |
|---|---|
| DataService | configure and initialize application; parse URLs; print headers; |
| Core | load API and pass options; print data object in requested format; manage caching, authorization; |
| API::* | validate input; call PhEDEx library functions; format and return data object; perldoc API documentation; |
| Rest of PhEDEx | PhEDEx business logic and SQL |

***An API call returns only one data structure.*** An API call does one thing and one thing only.  No option shall change the format of the returned data.  This is to ensure that clients cannot be surprised by results and know what to expect.

***Common entities have required attributes across all API calls.***  Entities with unique IDs shall always have that ID as an attribute.  Basic attributes (e.g. number of files in a block) will appear with that entity no matter what the context is.  This is to allow for client-side correlation of results from separate calls.

***Utilize hierarchy wherever possible.***  Do not flatten results, even where it seems convenient.  The full context of data entities should be a part of the result, and the client should not have to rely on the options to the call to successfully interpret the results.
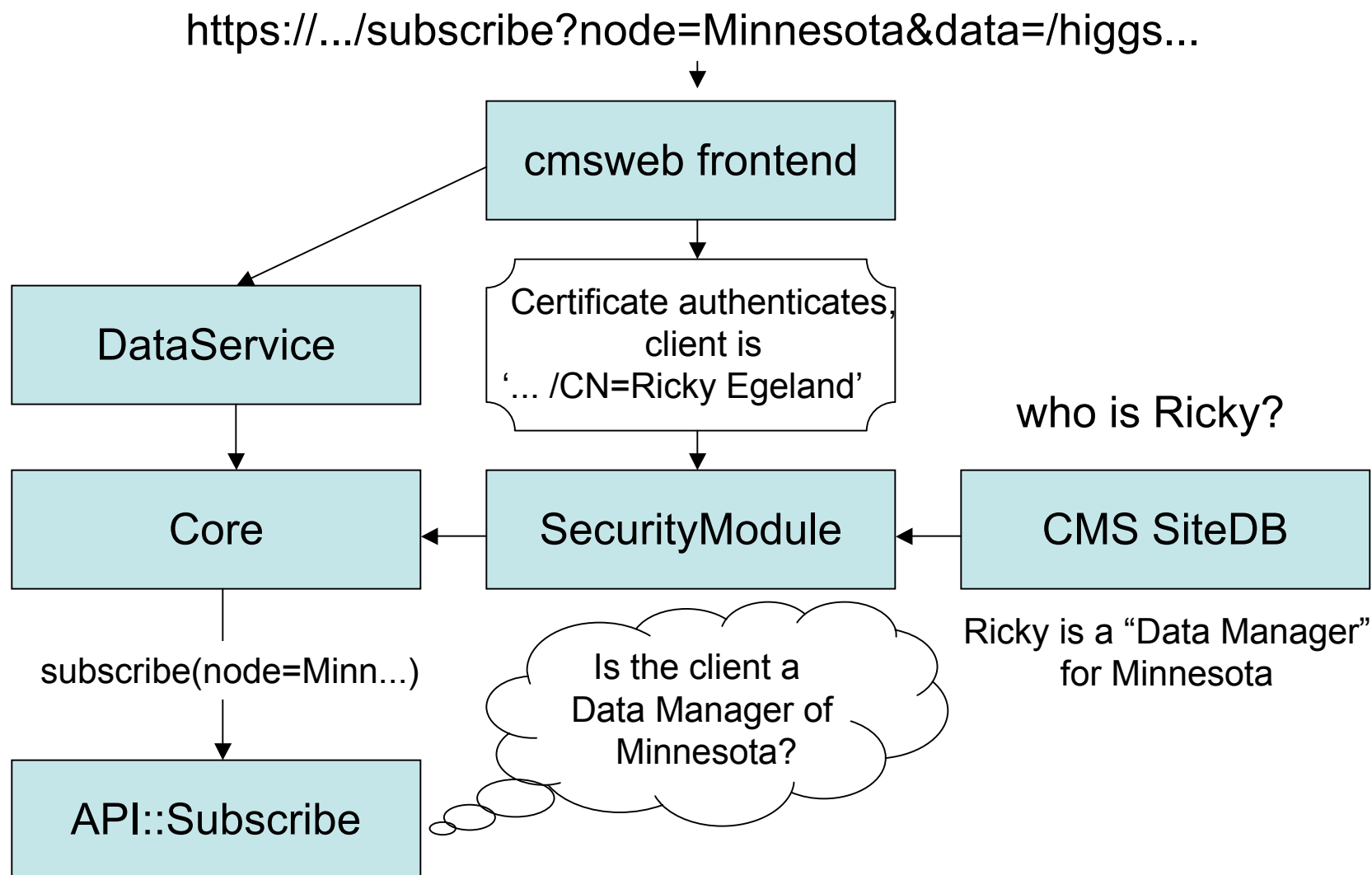
***Consistent call/response semantics.***  An attribute which is filterable should have the same name in the response as in the input.
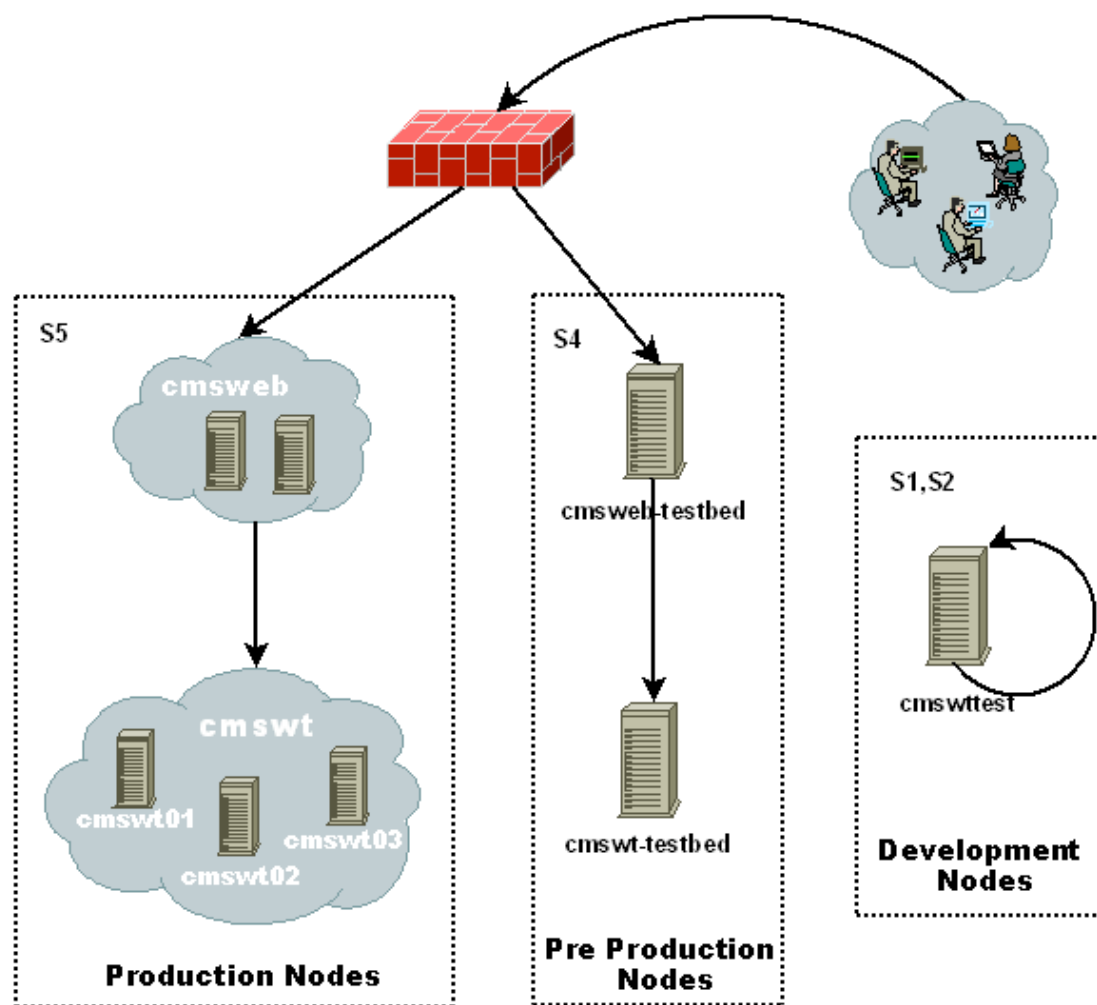
*We develop the data service as if we were developing a library for future unknown and inexperienced developers.*

# Design: SecurityModule

https://.../subscribe?node=Minnesota&data=/higgs...

cmsweb frontend

DataService

Certificate authenticates, client is
'... /CN=Ricky Egeland'

who is Ricky?

Core ← SecurityModule ← CMS SiteDB

subscribe(node=Minn...)

Is the client a Data Manager of Minnesota?

Ricky is a "Data Manager" for Minnesota

API::Subscribe

- Load-balanced frontends providing certificate authentication

- Load-balanced backends hosting the data service

- Backends only accept requests from the frontends

(diagram by Patricia Bittencourt Sampaio, UERJ, Brazil)

# Security & Deployment Summary

- cmsweb frontend provides certificate authentication
- SecurityModule+SiteDB provides authorization data
- Data Service implements fine-grained, secure access to APIs
- Only the frontends are accessible to the outside world
- The data service only accepts requests from the frontend
  - enforced both in the host firewall and in the server configuration

- First officially supported client, 'phedex'
- Uses LWP::UserAgent HTTP library
- Prints data "as-is" from the data service
- "report format" plugins for human-readable output

```
> phedex --format xml nodes > nodes.xml

> phedex --format xml blockreplicas --block '...' > myblock.xml

> phedex --format report nodes
NAME, SE, KIND, TECHNOLOGY
T0_CH_CERN_Export, (undef), Buffer, Castor
T0_CH_CERN_MSS, srm-cms.cern.ch, MSS, Castor
...

> phedex --format perl blockreplicas --block '...' | perl -e '...'
```

- For development of APIs without the overhead of running the full apache / mod_perl stack

- A subclass of the 'phedex' command-line tool which fakes the **_server response_**
    - faked authentication
    - http server access replaced by direct DataService library invocation
    - Exercises everything in the client/server stack

- Uncaught exceptions are detected immediately and reported at the command prompt
    - no log trawling

- Possible to run in a debugger

## DLS = Data Location Service

dls.getLocations per-block delay vs clients



configuration:  10 server children

(plot and study courtesy Antonio Delgado Peris, CIEMAT, Spain)

- **CMS Computing projects**
  - DLS Client / CRAB (Grid Analysis)
  - Data Discovery Page
  - Tier-0
- **Other projects**
  - Custom site consistency checking tools
  - Netvibes "PhEDEx Download Status" widget
- **"Other projects" makeup the majority of hits to the service**
  - Great!

**Monthly Accesses to PhEDEx Data Service**

155,816

146,507

80,947

Dec-08    Jan-09    Feb-09

Hits

Month

**Requested Data Format in PhEDEX Data Service**
(from Dec-08 though Feb-09)

json
0.03%

xml
36%

perl
64%

# Usage : DBS Discovery

**DBS** = **D**ataset **B**ookkeeping **S**ervice



https://cmsweb.cern.ch/dbs_discovery

# Usage : Netvibes Widget



http://www.netvibes.com/gonis#Higgs_V11                    widget by Isidro Gonzàlez – CIEMAT, Spain

## Accessing PhEDEx Data Service Programmatically (in Perl)

We've developed WebTools::PhedexSvc, a Perl module, in order to access the PhEDEx data service programmatically. The class and a few tools that we discuss below can be found in tools.tar. Download and uncompress the tar-ball suitably, edit tools/setup.sh to adjust BASEDIR and source setup.sh in order to be able to access the service.

## Find blocks for a fraction of a dataset

There are occasions when on user request we transfer only a small fraction of a large number of datasets to a T2 site (temporarily in most cases). While subscribing only a fraction of a dataset, we must have a list of blocks that corresponds to the fraction and doing that manually for many a dataset is pretty tedius. Here is a tool that we have developed in Pisa.

First of all prepare an input file (e.g dataset.txt) that contains the (dataset, fraction, SE) tuple as shown below

```
/W0jet-alpgen/CMSSW_1_5_2-CSA07-2203/GEN-SIM-DIGI-RECO              0.102 cmssrm.fnal.gov
/Z0jet-alpgen/CMSSW_1_5_2-CSA07-1193756147/GEN-SIM-DIGI-RECO        0.102 ccsrm.in2p3.fr
/W5jet_100ptw300-alpgen/CMSSW_1_5_2-CSA07-2225/GEN-SIM-DIGI-RECO    0.102 srm-cms.gridpp.rl.ac.uk
```

We can now query the PhEDEx data service as follows

```
[phedex@phedex ~]$ cd bin
[phedex@phedex bin]$ perl -w findBlocks.pl dataset.txt
```

https://twiki.cern.ch/twiki/bin/view/CMS/ItalianT2ToolsDataTransfers

tools by Subir Sarkar – INFN-Pisa, Italy

- 2ⁿᵈ official client:  PhEDEx web site

| More datasvc APIs for official website | → | Interesting UIs from outside developers | → | Integration back into official website |
|---|---|---|---|---|

- More RESTful?
  - Current APIs represent a search, not a resource
- Implement caching
  - Caching results of a wildcard search on a dynamic source leads to a low hit ratio
  - Caching results of distinct entity accesses would be more feasible
  - Find a performance balance between bulk accesses and cacheable resources

- PhEDEx data service satisfies integration requirements of CMS computing
- Provides a platform for increased developer involvement
- Allows for increased code re-use within and outside the project
- The data service is planned to become an even more integral component to the PhEDEx project

More info:
PhEDEx Data Service:     http://cmsweb.cern.ch/phedex/datasvc/doc
PhEDEx Web Site:         http://cmsweb.cern.ch/phedex
contact PhEDEx:          cms-phedex-admins@cern.ch