
Panda, a Pilot-based workflow manager



Open Science Grid

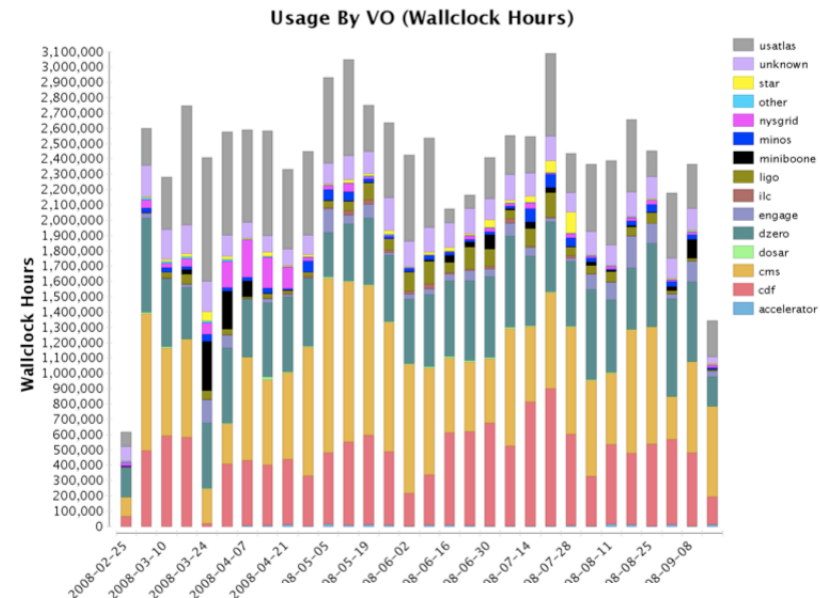
New Mexico Grid School – April 8, 2009

Marco Mambelli – University of Chicago

marco@hep.uchicago.edu

The ATLAS VO

- ▶ Virtual Organization in OSG (and other Grids)
 - ▶ In OSG since the beginning
 - ▶ <https://twiki.grid.iu.edu/bin/view/VO/ATLAS>
 - ▶ <https://lcg-voms.cern.ch:8443/vo/atlas/vomrs>
- ▶ Collaboration for the ATLAS experiment in the LHC at CERN
 - ▶ <http://atlas.ch/>
 - ▶ http://atlas.web.cern.ch/Atlas/ATLASreg_form.pdf



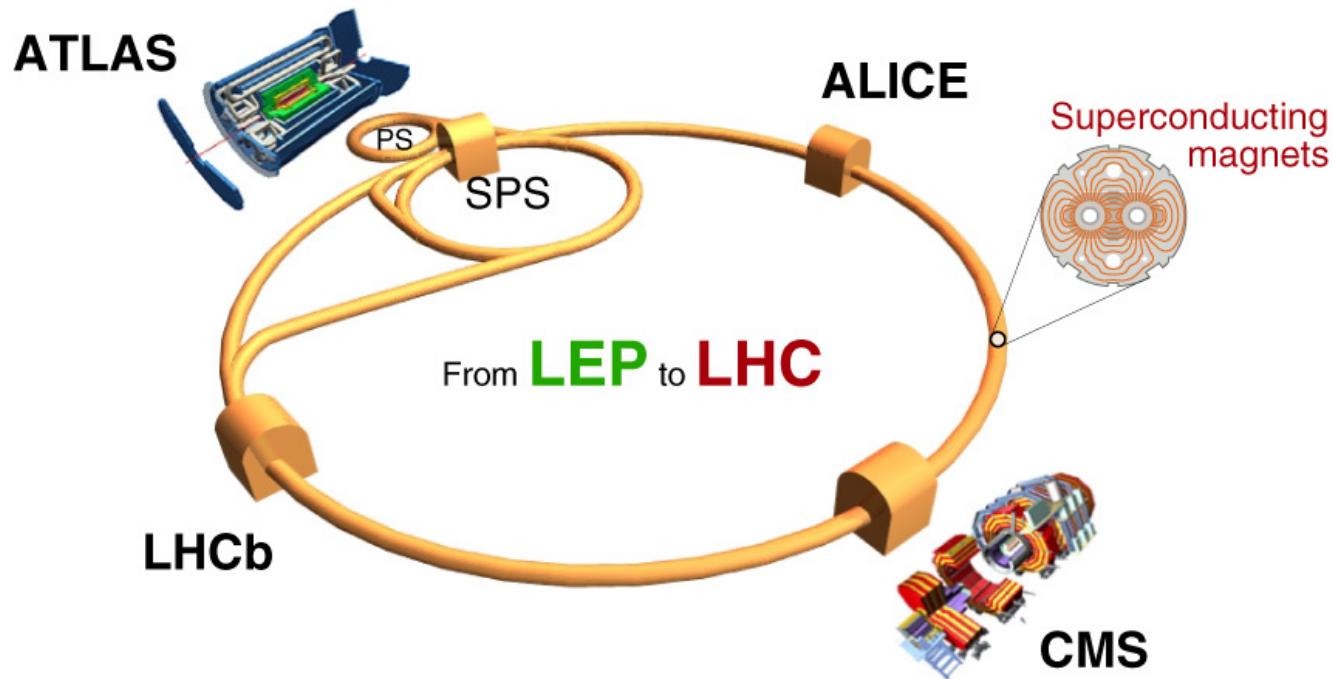
ATLAS REGISTRATION FORM



Please complete all sections and return this form to the
ATLAS SECRETARIAT, Bldg 40-4-D01, fax+41227678350

*SURNAME:	*FIRST NAME:
*DATE OF BIRTH (DAY/MONTH/YEAR):	*SEX: M/F
*MARITAL STATUS: Single / Married / Divorced	*NATIONALITY:
*CERN IDENTIFICATION NUMBER:	*CERN PHONE NUMBER:
*CERN OFFICE:	
*PRIVATE HOME ADDRESS AND PHONE NUMBER:	
*HOME ADDRESS WHEN AT CERN (e.g. CERN HOSTEL MEYRIN or ST GENIS-POUILLY):	
*HOME INSTITUTE NAME AND ADDRESS:	
*HOME INSTITUTE OFFICE TELEPHONE AND FAX NO. (DIRECT):	
*EMAIL ADDRESS:	
*SUBDETECTOR OR DOMAIN OF ACTIVITY: (you may cross several)	

LHC experiment at CERN

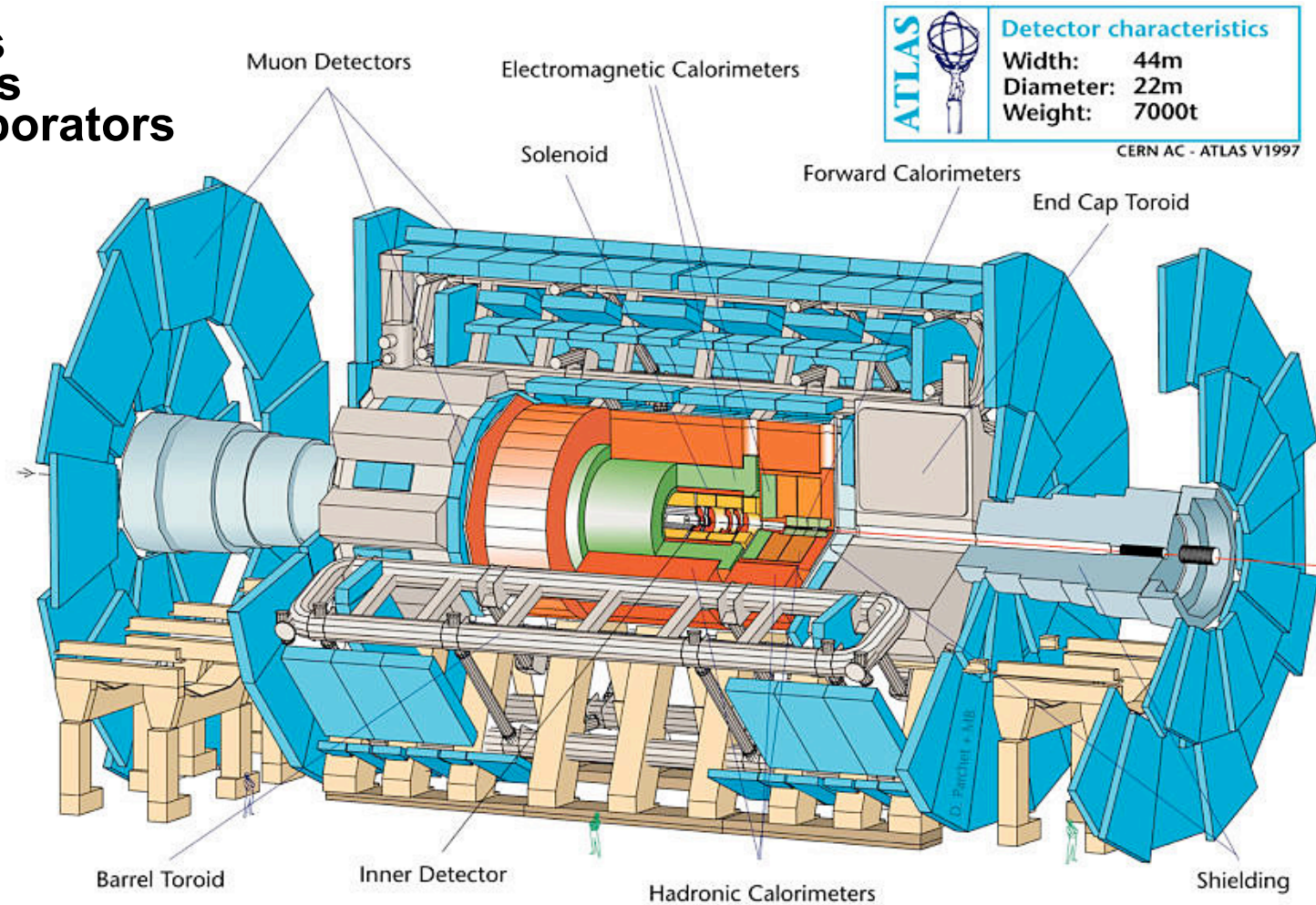


	Beams		Energy		Luminosity
LEP	e^+	e^-	200	GeV	$10^{32} \text{ cm}^{-2} \text{ s}^{-1}$
LHC	p	p	14	TeV	10^{34}

<http://www.youtube.com/watch?v=j50ZssEojtM>

The ATLAS experiment

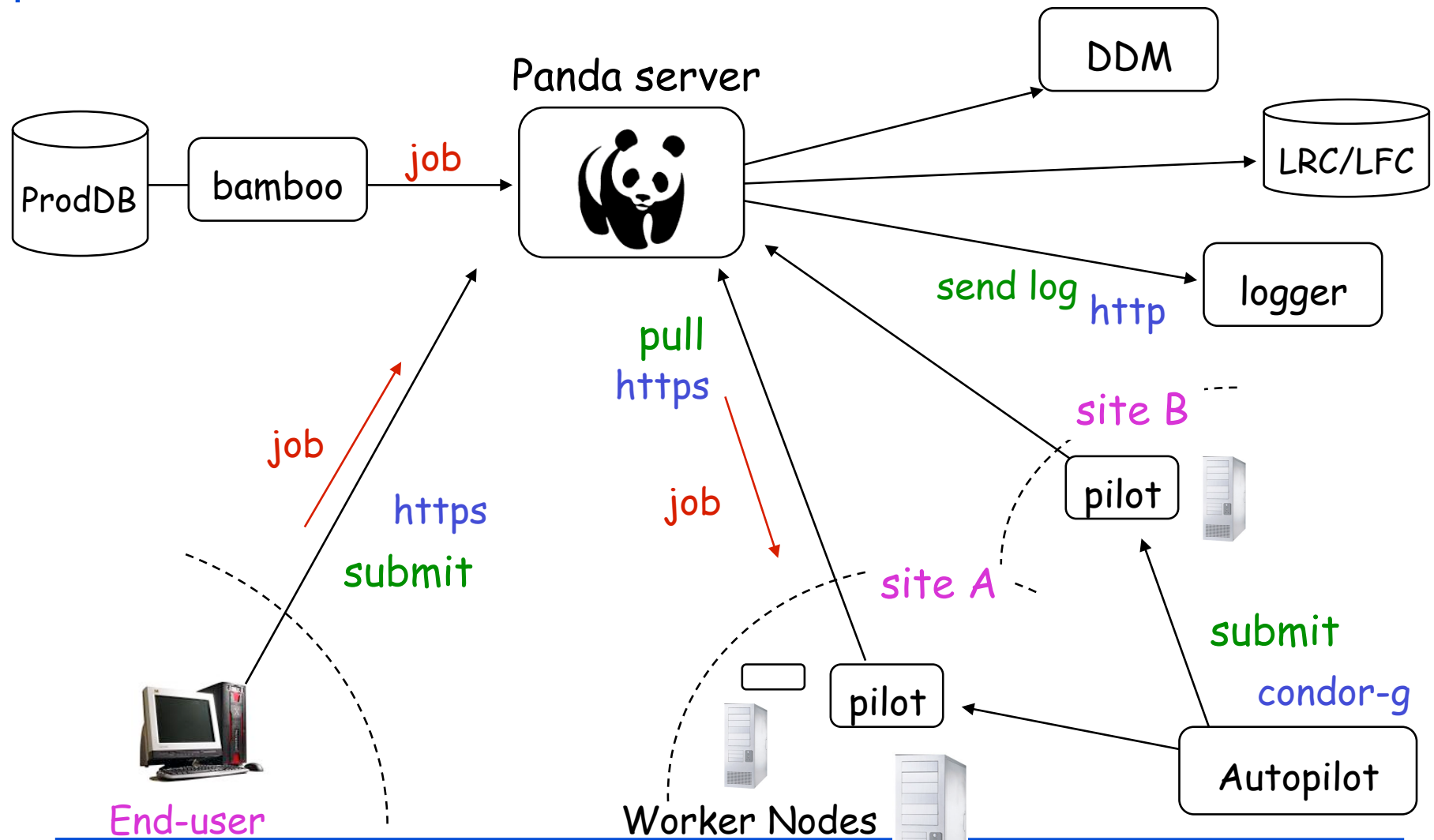
37 Countries
167 Institutes
~2000 Collaborators



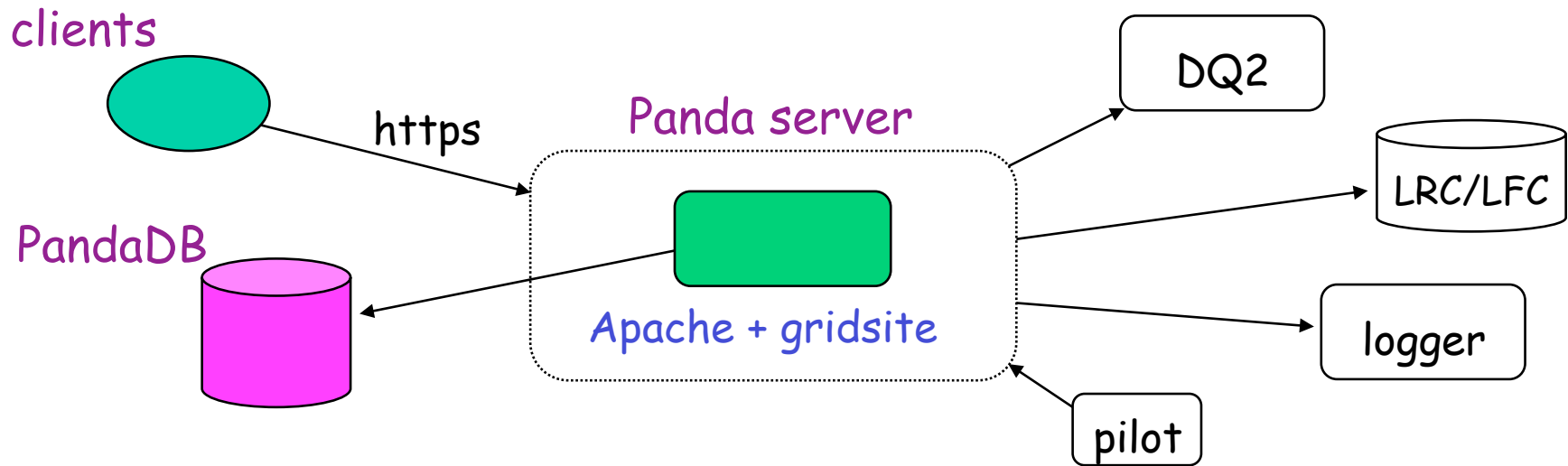
PANDA

- ▶ PANDA = Production ANd Distributed Analysis system
 - ▶ Designed for analysis as well as production for High Energy Physics
 - ▶ Works both with OSG and EGEE middleware
- ▶ A single task queue and pilots
 - ▶ Apache-based Central Server
 - ▶ Pilots retrieve jobs from the server as soon as CPU is available
→ late scheduling
- ▶ Highly automated, has an integrated monitoring system
- ▶ Integrated with ATLAS Distributed Data Management (DDM) system
- ▶ Not exclusively ATLAS: has its first OSG user in CHARMM (Chemistry at HARvard Molecular Mechanics)

Panda System

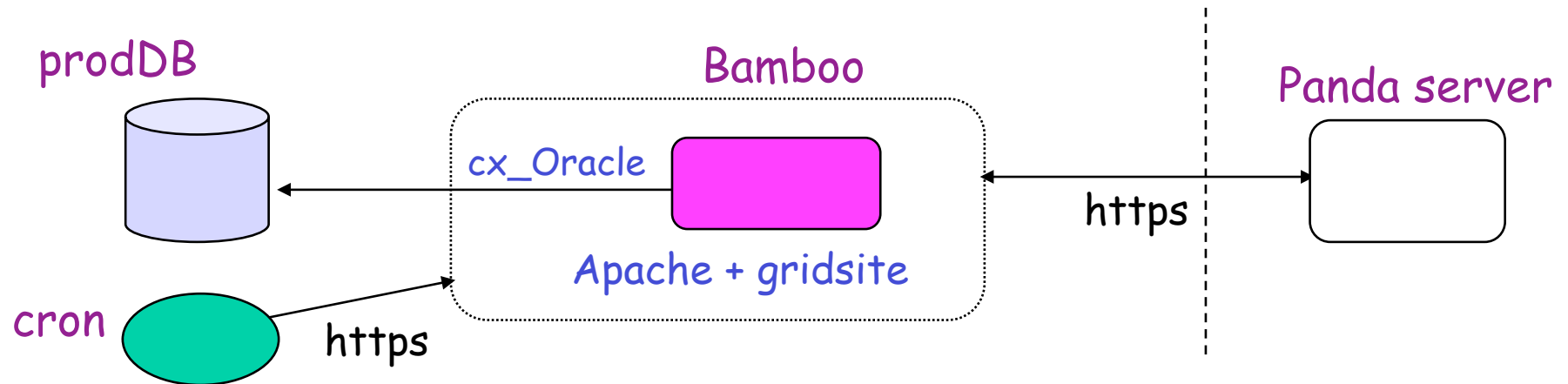


Panda Server



- ▶ Central queue for all kinds of jobs
- ▶ Assign jobs to sites (brokerage)
- ▶ Setup input/output datasets
 - ▶ Create them when jobs are submitted
 - ▶ Add files to output datasets when jobs are finished
- ▶ Dispatch jobs

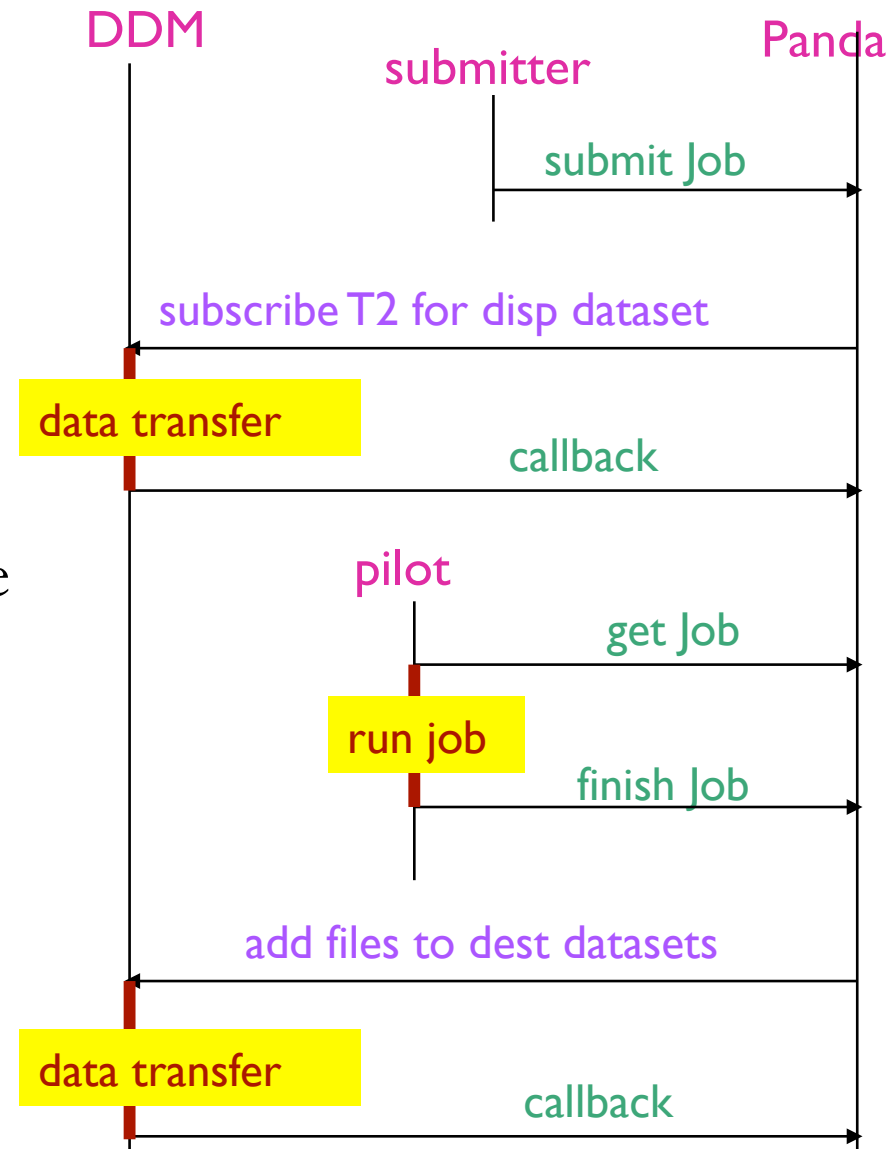
Bamboo



- ▶ Get jobs from prodDB to submit them to Panda
 - ▶ Update job status in prodDB
 - ▶ Assign tasks to clouds dynamically
 - ▶ Kill TOBEABORTED jobs
-
- ▶ A cron triggers the above procedures every 10 min

Panda Job Timeline

- ▶ Rely on ATLAS DDM
 - ▶ Panda sends requests to DDM
 - ▶ DDM moves files and sends notifications back to Panda
 - ▶ Panda and DDM work asynchronously
- ▶ Dispatch input files to execution sites and aggregate output files to destination
- ▶ Jobs get 'activated' when all input files are copied, and pilots pick them up
 - ▶ Pilots don't have to transfer data (asynchronous)
 - ▶ Data-transfers and Job-executions can run in parallel



How the pilot works

- ▶ Sends the several parameters to Panda server for job matching (HTTP request)
 - ▶ CPU speed
 - ▶ Available memory size on the WN
 - ▶ List of available ATLAS releases at the site
- ▶ Retrieves an `activated` job (HTTP response of the above request)
 - ▶ activated → running
- ▶ Runs the job immediately because all input files should be already available at the site
- ▶ Sends heartbeat every 30min
- ▶ Copy output files to local Storage Element and register them to Local Replica Catalog

Pilot vs ATLAS Job

Pilot

- Submitted by factories
 - remote submit hosts
 - local cluster factories
- Managed by factories
- Python code to support ATLAS Job execution
- Submitted continuously
- Partially accounted
 - no big deal if some fail

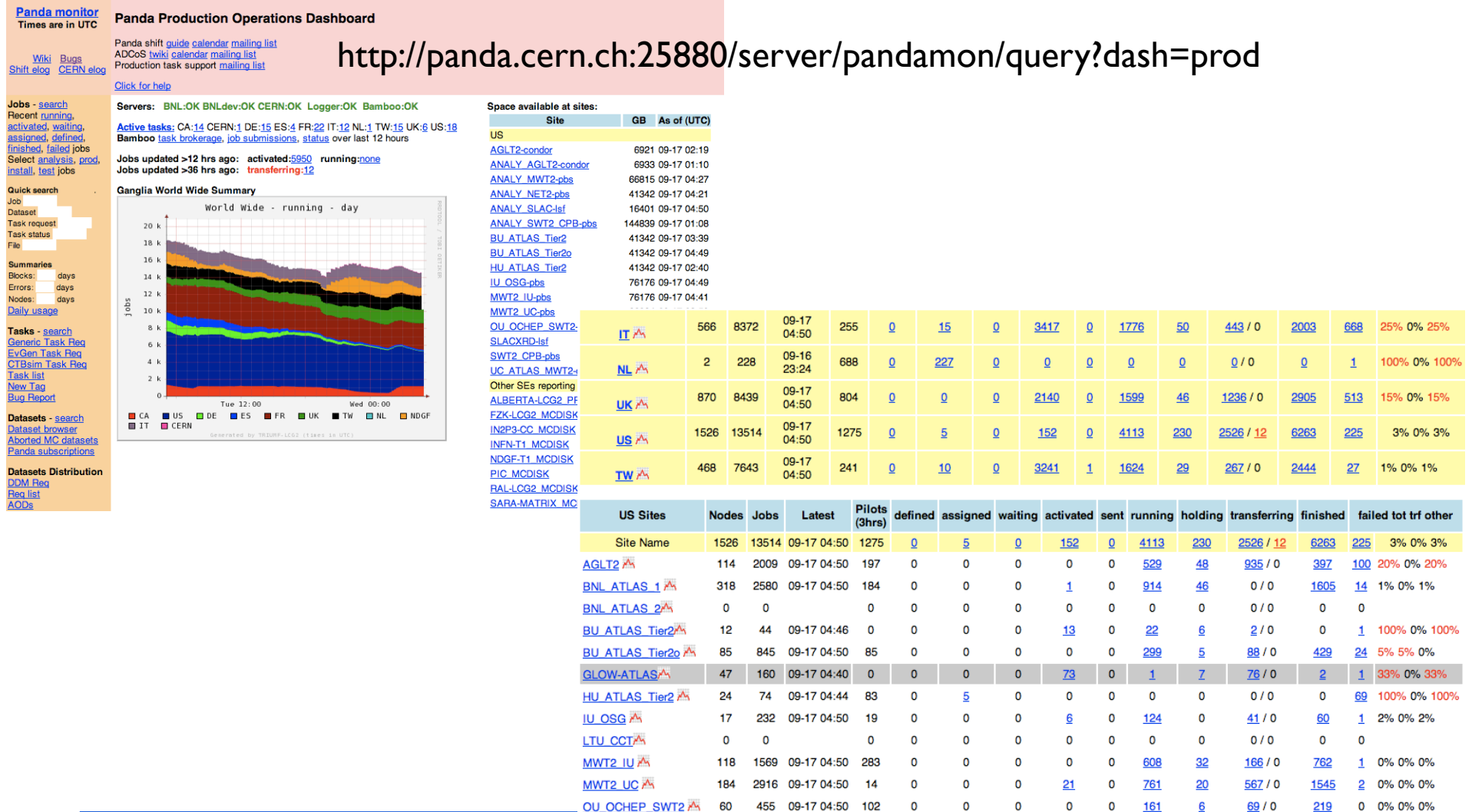
ATLAS Job

- Submitted by users or production managers (Bamboo)
- Managed by Panda Server
- Runs Athena software (ATLAS libraries)
- Submitted when needed
- Fully accounted
 - error statistics are important

Some monitoring resources

- ▶ The following pages present some monitoring example
- ▶ Screenshots are just example pages, actual content varies
- ▶ URLs are one of the possible URLs providing a similar page
 - ▶ e.g. queries may vary the actual Site or Time interval
- ▶ Main URLs:
 - ▶ DDM Dashboard: <http://dashb-atlas-data-test.cern.ch/dashboard/request.py/site>
 - ▶ Panda Monitor: <http://panda.cern.ch:25880/> or <http://panda.atlascomp.org/?redirect=pandamon>
(hostname may change since there are multiple servers)
- ▶ Take time to navigate Panda Monitor and the Dashboard

Panda Monitor: production dashboard



Panda Monitor: Dataset browser

[Configuration](#)
[Update](#)
[Panda monitor](#)
[Quick guide](#), [twiki](#)

[Production](#) [Clouds](#) [DDM](#) [PandaMover](#) [AutoPilot](#) [Sites & Grids](#) [Analysis](#) [Physics data](#) [Usage & Quotas](#) [Plots](#) [ProdDash](#) [DDMDash](#)

[Jobs - search](#)
Recent [running](#),
[activated](#), [waiting](#),
[assigned](#), [defined](#),
[finished](#), [failed](#) jobs
Select [analysis](#),
[production](#), [test](#) jobs

Quick search
Job
Dataset
Task
File

Summaries
Blocks: days
Errors: days
Nodes: days
[Daily usage](#)

Tasks - search
[Generic Task Req](#)
[EvGen Task Req](#)
[CTBsim Task Req](#)
[Task list](#)
[Task browser](#)

Datasets - search
[Dataset browser](#)
[New datasets](#)
[Aborted MC datasets](#)

DQ2 dataset browser for csc datasets

[Click for help](#)

Dataset lists last updated 157 min ago

Select a project:

Or (the old way) select a dataset category *Counts are totals, exclusive of selections*

Category	Count	Description
All	147033	All datasets
DBrelease	12	DB release datasets
M3	386	M3 cosmics run
M4	4210	M3 cosmics run
conditions	35	Datasets for conditions data files
csc	4411	Computing system commissioning production
ctb	613	Combined testbeam production
dc2	6	Data Challenge 2 production
larg	53	LAr commissioning
mc	6135	MC validation production
other	71	Everything else
rome	210	Rome physics workshop production
testpanda	6439	Panda test datasets
tile	52	Tilecal commissioning
user	58456	User datasets
validation	642	Validation samples (testIdeal* etc)

Choose a site if you want to restrict dataset listings to site-resident datasets

CANADA	CERN	FRANCE	GERMANY	ITALY	NDGF	NL	SPAIN	TAIWAN	UK
ALBERTA	CERNCAF	FRTIER2S	CSCS	CNAF	IJST2	IHEP	IFAE	ASGC	Rutherford
MCGILL	CERNPROD	LPC	CYF	CNAFDISK	NDGFDISK	ITEP	IFIC	ASGCDISK	Belle
MONTREAL	TIER0DISK	LAL	DESY-HH	CNAFTAPE	NDGFT1	JINR	IFICDISK	ASGCDISK V2	Belle
SFU	TIER0TAPE	SACLAY	DESY-ZN	CNAFTTEST	NDGFT1DISK	NIKHEF	IFICTAPE	ASGCTAPE	SH
TORONTO	LBNL	FZK	INFN	NDGFT1TAPE	PNPI	LIP-COIMBRA	ASGCTAPE V2	JINR	SH

<http://panda.cern.ch:25880/server/pandamon/query?overview=dslist>

Panda Monitor: error reporting

Panda monitor
Times are in UTC

[Wiki](#) [Bugs](#)
[Shift elog](#) [CERN elog](#)

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](#)

Datasets - [search](#)
[Dataset browser](#)
[Aborted MC datasets](#)
[Panda subscriptions](#)

Datasets Distribution
[DDM Req](#)
[Req list](#)

Panda job error summary for last 24 hours (1.0 days)

All CEs and jobs. Show [production](#), [analysis](#), [test](#), all jobs/CEs

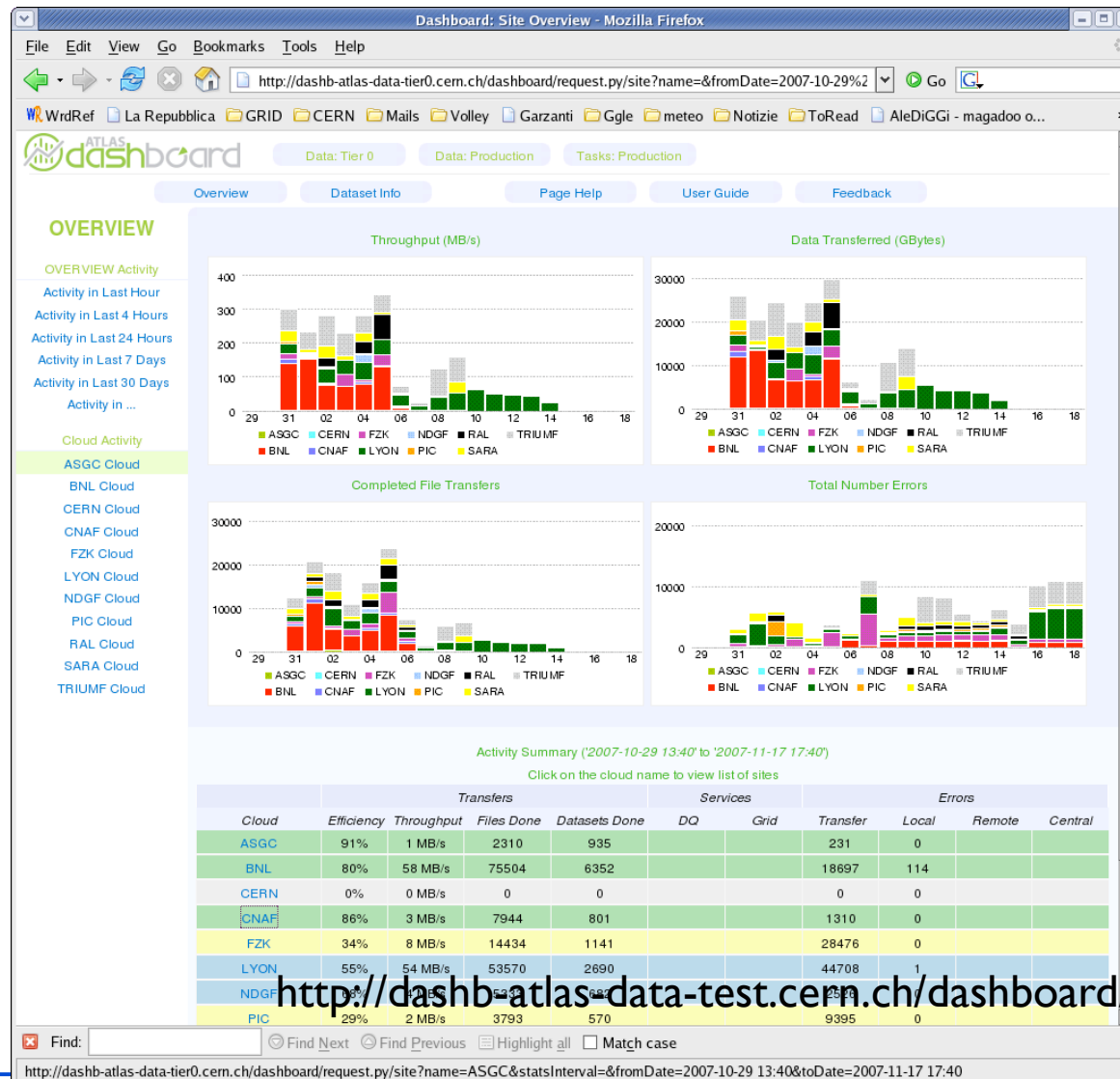
Job wall time: 317553 hrs Error losses: trans: 9971 (3.1%) panda: 8458 (2.7%) ddm: 3329 (1.0%) other: 1317 (0.4%)

Error type (type count)	Count	CPU-hrs	Latest	Code: Description
All	defined :708 failed :8358	assigned :264 (17.1%)	waiting :0	activated :19020 sent :0 running :10202 holding :1599 transferring :5359 finished :40421
brokerageErrorCode (120)	120	0.0	09-17 13:11	100 : Unknown error code
ddmErrorCode (6)	1	0.0	09-16 18:14	100 : DQ2 server error
ddmErrorCode (6)	5	14.0	09-17 13:54	200 : Could not add output files to dataset
exeErrorCode (1114)	2	2.6	09-16 13:45	1101 : LRC registration error: Connection refused
exeErrorCode (1114)	1	0.9	09-16 20:13	1114 : Put error: Failed to import LFC python module
exeErrorCode (1114)	4	30.3	09-16 18:57	1131 : Put function can not be called for staging out
exeErrorCode (1114)	31	13.6	09-17 04:50	1132 : LRC registration error (consult log file)
exeErrorCode (1114)	7	14.8	09-17 10:25	1133 : Put error: Fetching default storage URL failed
exeErrorCode (1114)	1	26.2	09-15 10:22	1135 : Could not get file size in job workdir
exeErrorCode (1114)	875	7494.9	09-16 22:32	1137 : Put error: Error in copying the file from job workdir to localSE
exeErrorCode (1114)	13	159.2	09-16 15:34	1154 : Failed to register log file
exeErrorCode (1114)	6	58.6	09-16 15:20	1155 : Failed to move output files for lost job
exeErrorCode (1114)	1	11.8	09-14 15:01	1176 : Pilot has no child process
exeErrorCode (1114)	1	22.1	09-15 07:10	1211 : Missing installation
exeErrorCode (1114)	3	51.7	09-17 13:52	60000 : segmentation violation
exeErrorCode (1114)	117	399.1	09-17 13:44	60010 : segmentation fault
exeErrorCode (1114)	5	92.2	09-17 10:47	61200 : ServiceManager Unavailable
exeErrorCode (1114)	6	107.3	09-17 13:36	62600 : AthenaCrash
exeErrorCode (1114)	30	94.8	09-17 10:27	64100 : Transform output file error
exeErrorCode (1114)	11	52.2	09-17 12:15	69999 : Unknown Transform error

OSG Errors for period 886-2007

<http://panda.cern.ch:25880/server/pandamon/query?days=1&overview=errorlist>

DDM Dashboard: overview



?

!

Client-Server Communication

- ▶ HTTP/S-based communication (curl+grid proxy+python)
- ▶ GSI authentication via mod_gridsite
- ▶ Most of communications are asynchronous
 - ▶ Panda server runs python threads as soon as it receives HTTP requests, and then sends responses back immediately. Threads do heavy procedures (e.g., DB access) in background → better throughput
 - ▶ Some are synchronous

