

AsyncStageOut: Distributed user data management for CMS Analysis

Hassen Riahi

Tony Wildish

Diego Ciangottini

CERN IT-SDC

Princeton University (US)

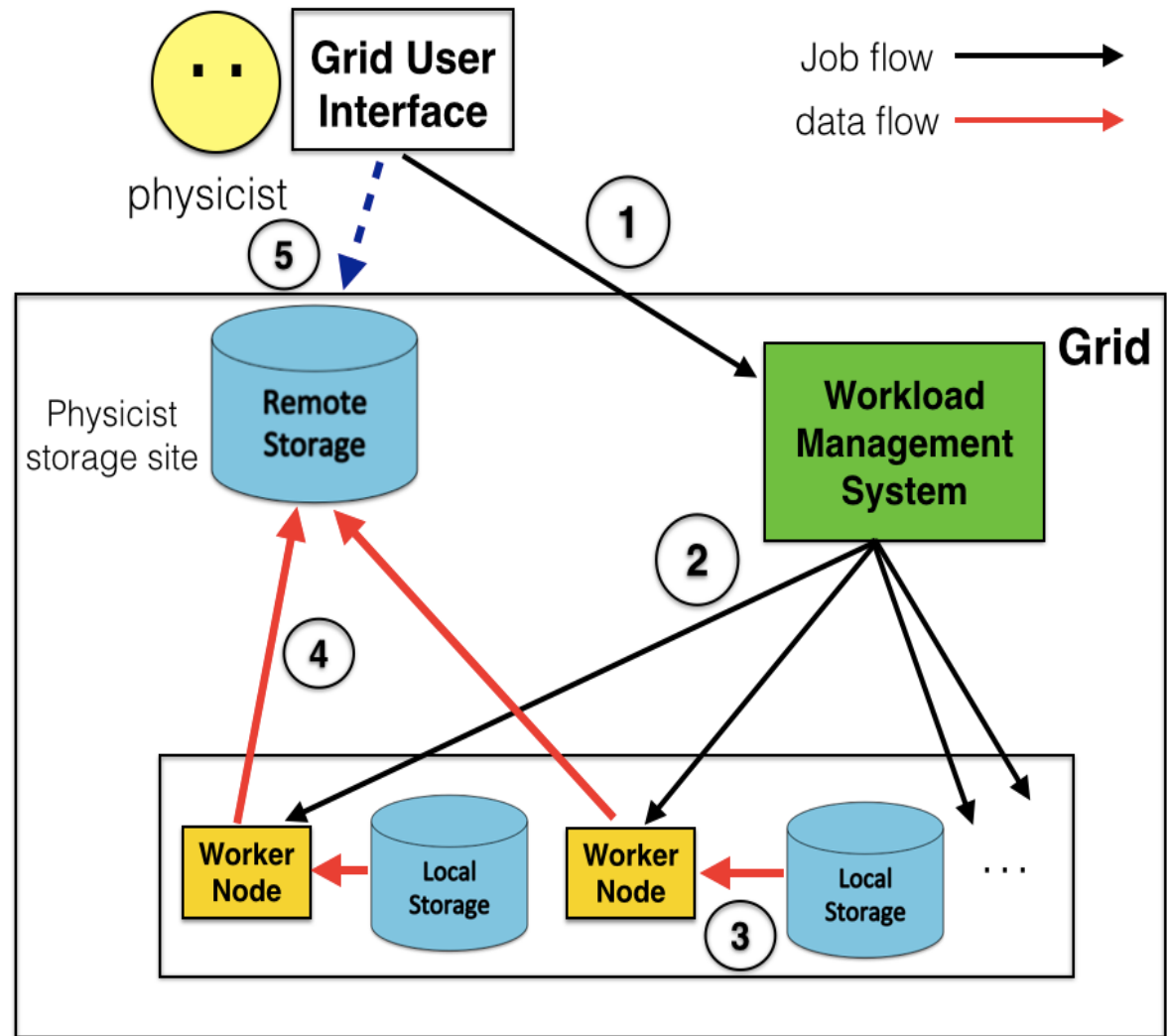
University & INFN Perugia (IT)

Outline

- Overview
- Problem and strategy
- Architecture
- Integration
- Results
- Conclusions

Distributed data analysis in CMS

- 1000 individual users per month
- More than 60 sites
- 20k jobs/hour
- Typically 1 file/job
 - Files vary in size
- 200k completed jobs per day
- Minimal latencies
- Chaotic environment

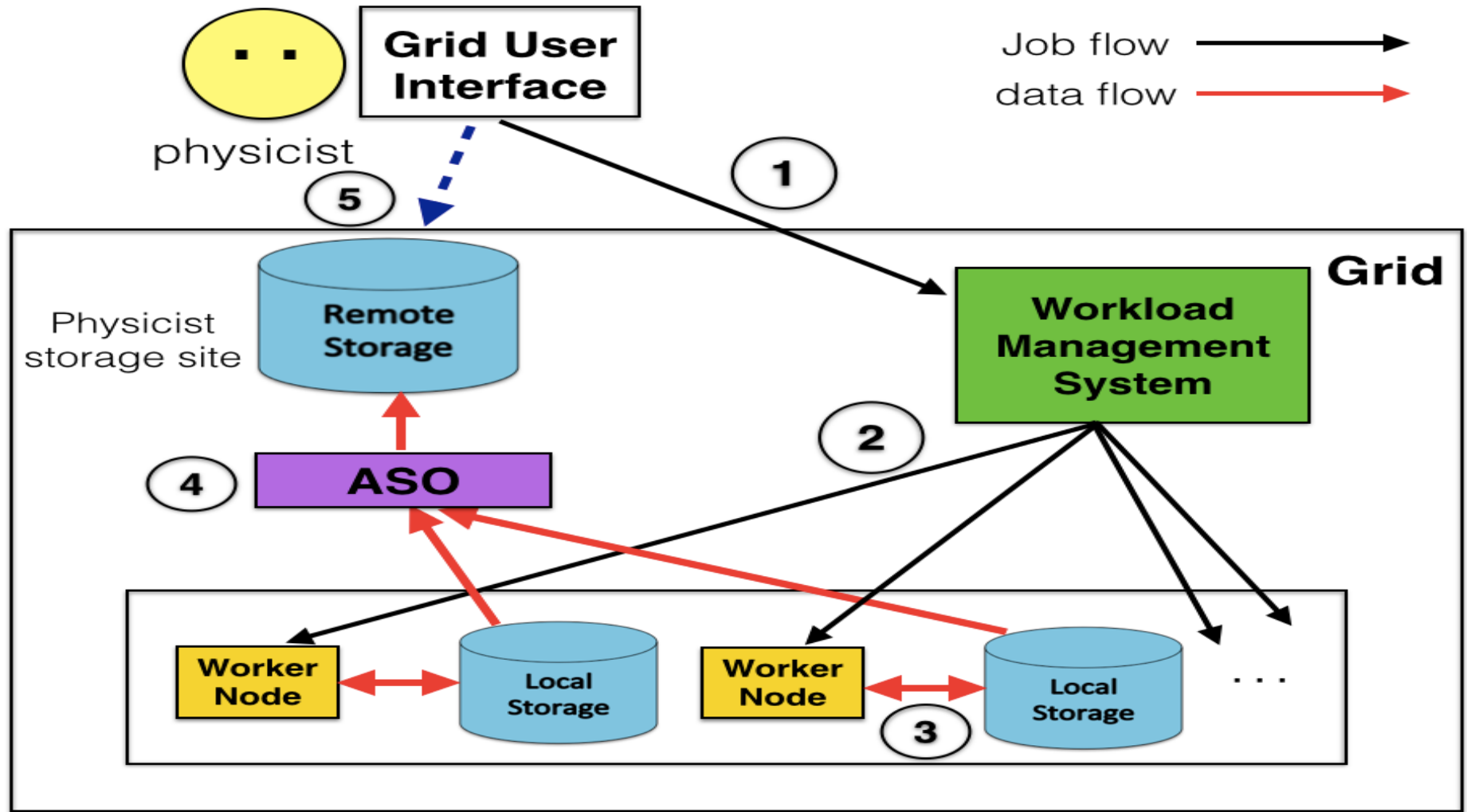


Problem

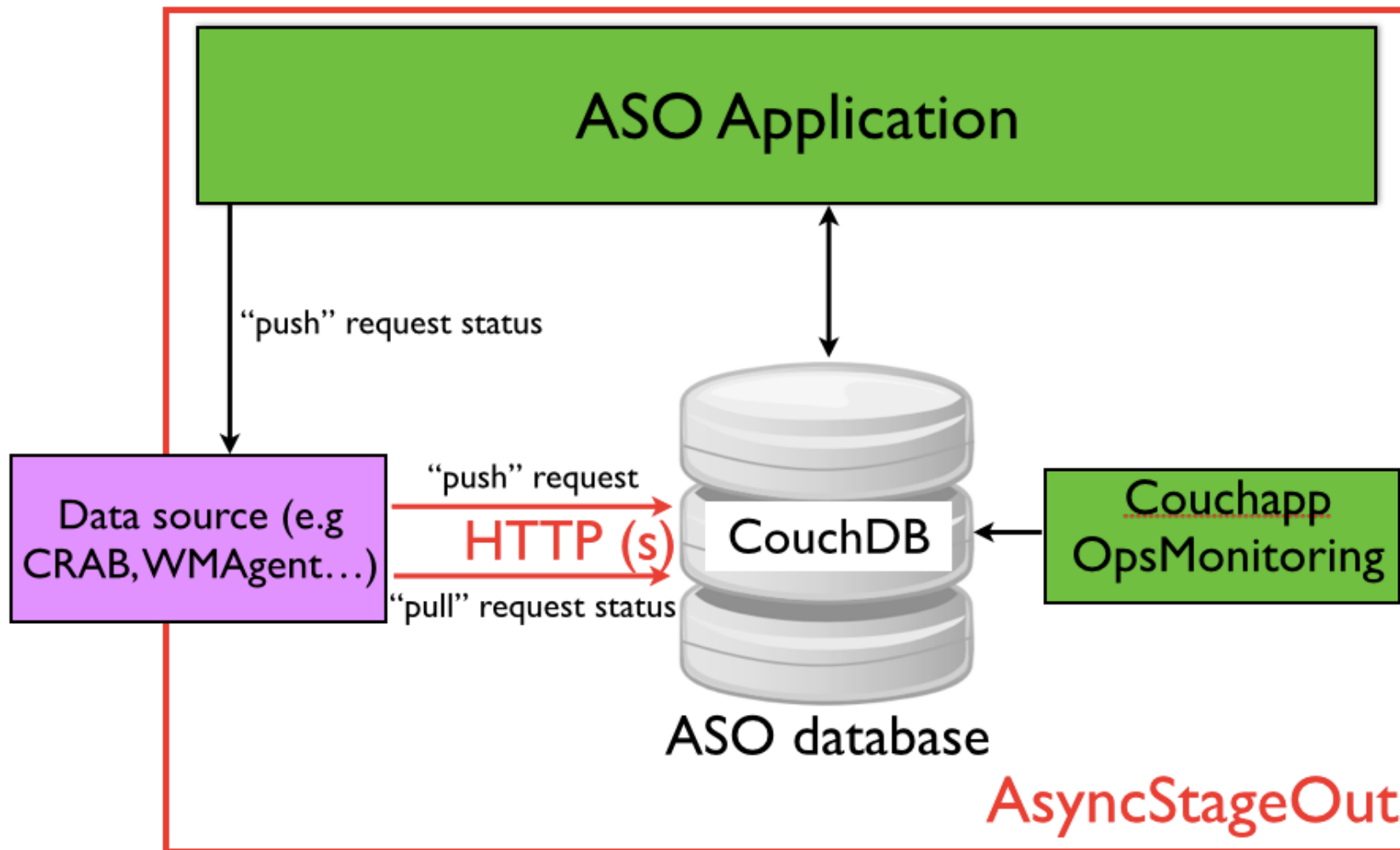
- 15% to 20% the jobs fail and about 30% to 50% of the failures are due to the jobs not being able to upload their output data to a remote disk storage
 - between 5% and 10% of jobs fail in the remote copy of outputs
 - the overall CPU loss is even higher than 5-10% since those jobs fail at the end of the processing after multiple retries
 - often it results in DDoS to CMS Tier-2 storage systems

AsyncStageOut (ASO) is implemented to reduce the **most common failure mode** of analysis jobs

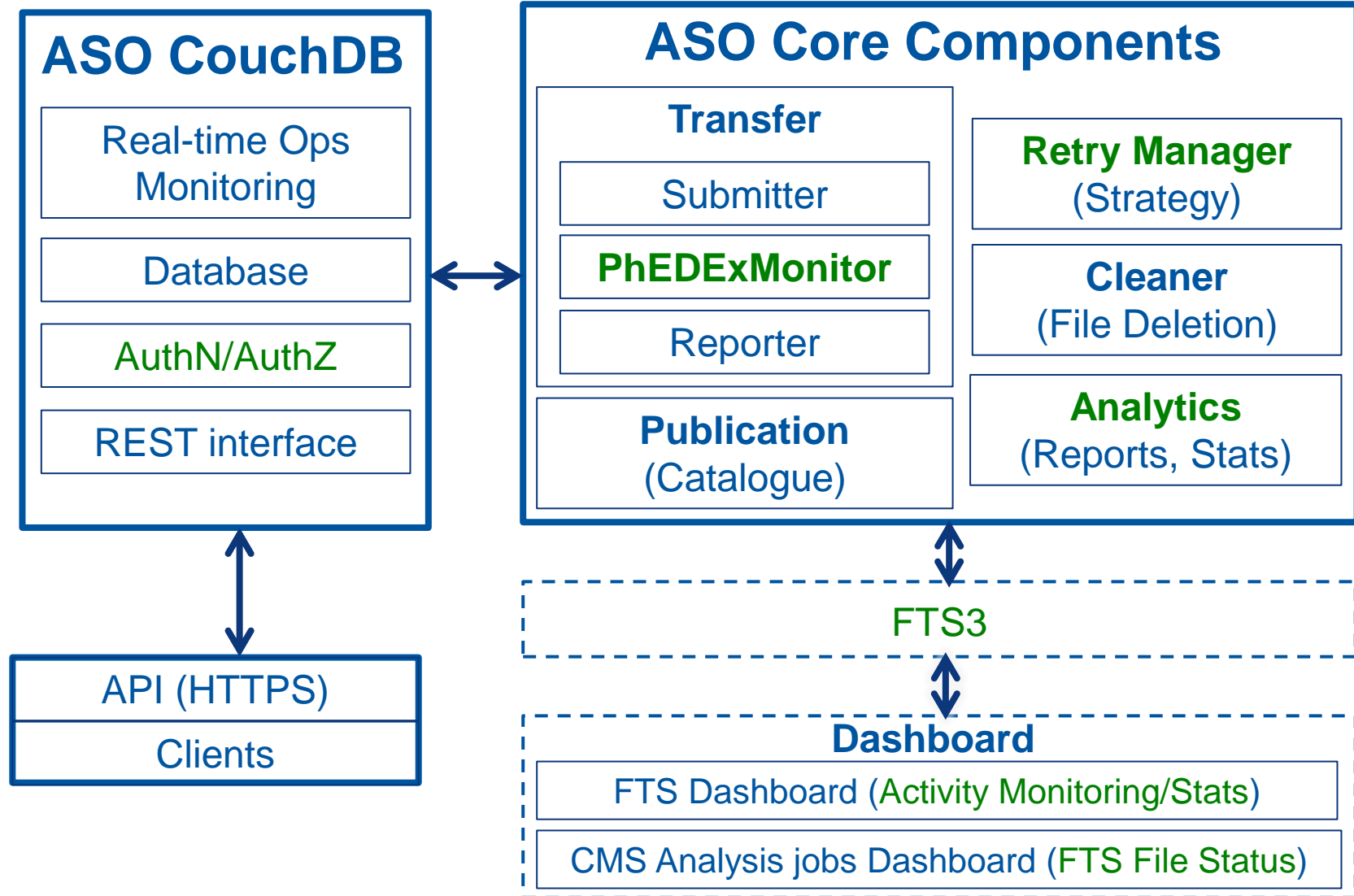
Asynchronous stage-out strategy



Architecture overview



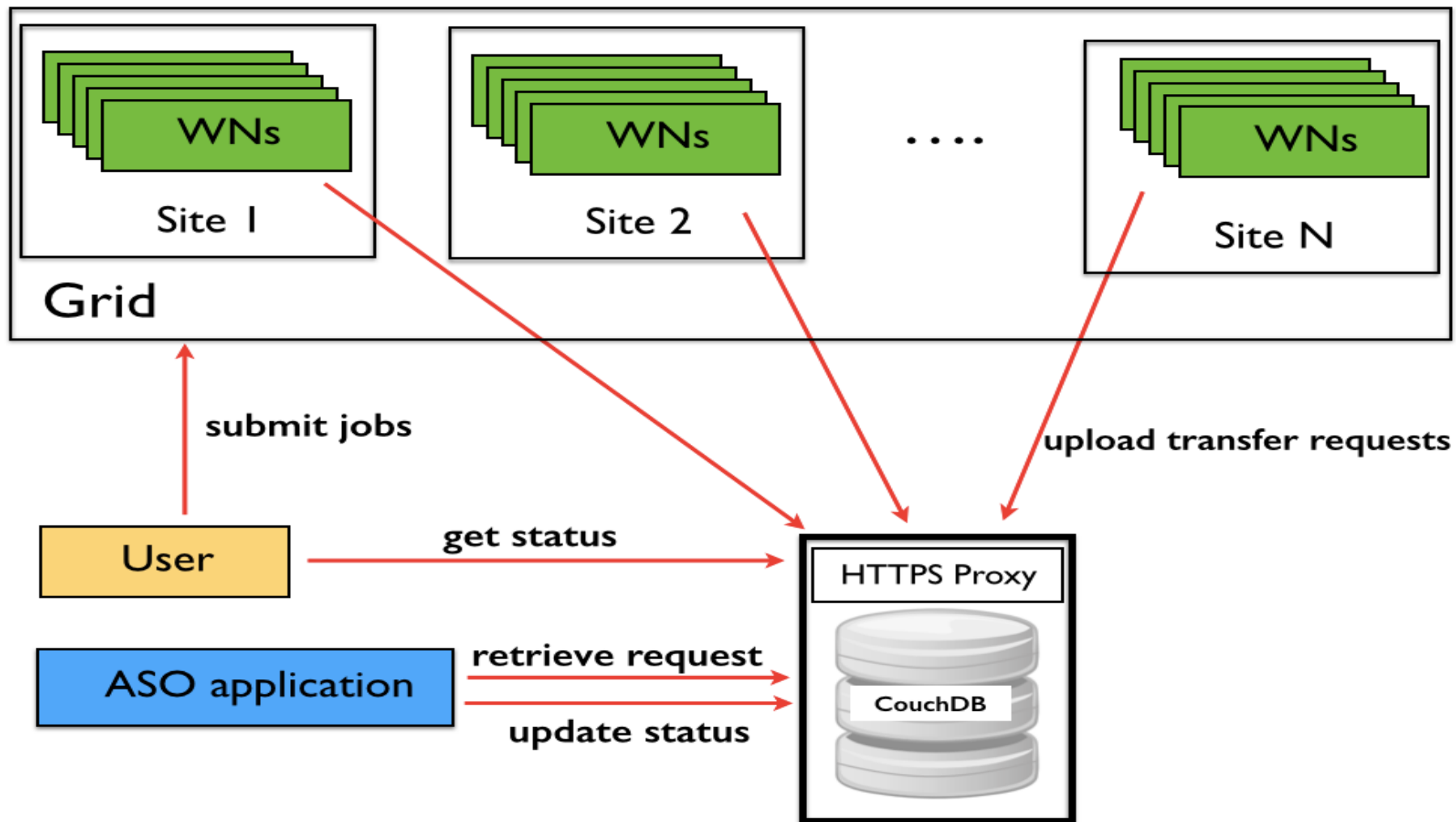
Architecture and evolution



CMS analysis jobs Dashboard

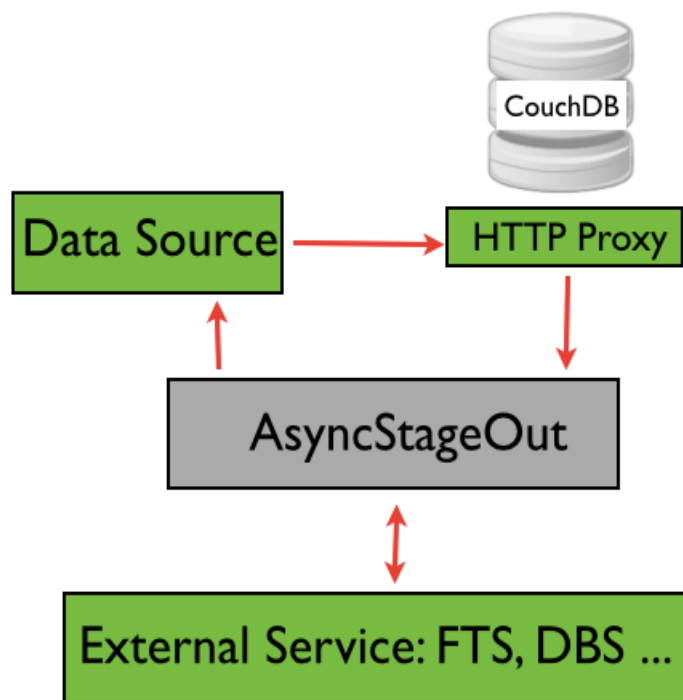
Data Charts Show 25 entries Task: 150306_135227:fromeo_crab_TTJetsbb_19_6 NJobTotal: 300 Pending: 0 Running: 0 Unknown: 0 Cancelled: 0 Success: 299 Failed: 1 WNPostProc: 0 ToRetry: 0 Search:											
Id	Status	AppExitCode	Site	Retries	Submitted	Started	Finished	Wall Time	Job Log	File Access	FTS File Status
1	finished	0	T1_US_FNAL	1	2015-03-06T13:53:47	2015-03-06T13:57:47	2015-03-06T14:11:07	00:13:20	Not available	File Info	FINISHED AT 2015-03-06 14:16:01
2	finished	0	T1_US_FNAL	1	2015-03-06T13:53:47	2015-03-06T13:57:45	2015-03-06T14:11:18	00:13:33	Not available	File Info	FINISHED AT 2015-03-06 14:16:09
3	finished	0	T1_US_FNAL	1	2015-03-06T13:53:47	2015-03-06T13:57:45	2015-03-06T14:11:36	00:13:51	Not available	File Info	FINISHED AT 2015-03-06 14:16:33
4	finished	0	T1_US_FNAL	1	2015-03-06T13:53:47	2015-03-06T13:57:50	2015-03-06T14:13:14	00:15:24	Not available	File Info	FINISHED AT 2015-03-06 14:19:16
5	finished	0	T2_DE_DESY	1	2015-03-06T13:53:47	2015-03-06T13:58:02	2015-03-06T14:06:21	00:08:19	Not available	File Info	FINISHED AT 2015-03-06 14:10:35
6	finished	0	T2_BE_IJHE	1	2015-03-06T13:53:47	2015-03-06T13:57:48	2015-03-06T14:15:52	00:18:04	Not available	File Info	FINISHED AT 2015-03-06 14:19:06
7	finished	0	T2_FL_HIP	1	2015-03-06T13:53:47	2015-03-06T13:57:39	2015-03-06T14:07:26	00:09:47	Not available	File Info	FINISHED AT 2015-03-06 14:11:19
8	finished	0	T2_CH_CSCS	1	2015-03-06T13:53:47	2015-03-06T13:58:13	2015-03-06T14:11:50	00:13:37	Not available	File Info	FINISHED AT 2015-03-06 14:16:30
9	finished	0	T2_BE_IJHE	1	2015-03-06T13:53:47	2015-03-06T13:57:42	2015-03-06T14:15:56	00:18:14	Not available	File Info	FINISHED AT 2015-03-06 14:19:06
10	finished	0	T2_FL_HIP	1	2015-03-06T13:53:47	2015-03-06T13:57:37	2015-03-06T14:09:25	00:11:48	Not available	File Info	FINISHED AT 2015-03-06 14:16:01
11	finished	0	T2_BE_IJHE	1	2015-03-06T13:53:47	2015-03-06T13:57:44	2015-03-06T14:15:53	00:18:09	Not available	File Info	FINISHED AT 2015-03-06 14:19:06
12	finished	0	T2_BE_IJHE	1	2015-03-06T13:53:47	2015-03-06T13:57:52	2015-03-06T14:16:12	00:18:20	Not available	File Info	FINISHED AT 2015-03-06 14:24:55
13	finished	0	T2_BE_IJHE	1	2015-03-06T13:53:47	2015-03-06T13:57:43	2015-03-06T14:15:40	00:17:57	Not available	File Info	FINISHED AT 2015-03-06 14:22:03
14	finished	0	T2_BE_IJHE	1	2015-03-06T13:53:47	2015-03-06T13:57:52	2015-03-06T14:16:06	00:18:14	Not available	File Info	FINISHED AT 2015-03-06 14:24:46
15	finished	0	T2_BE_IJHE	1	2015-03-06T13:53:47	2015-03-06T13:57:55	2015-03-06T14:16:05	00:18:10	Not available	File Info	FINISHED AT 2015-03-06 14:24:46
16	finished	0	T2_BE_IJHE	1	2015-03-06T13:53:47	2015-03-06T13:57:52	2015-03-06T14:16:05	00:18:13	Not available	File Info	FINISHED AT 2015-03-06 14:24:46
17	finished	0	T2_BE_IJHE	1	2015-03-06T13:53:47	2015-03-06T13:57:45	2015-03-06T14:16:00	00:18:15	Not available	File Info	FINISHED AT 2015-03-06 14:24:46
18	finished	0	T3_US_Colorado	1	2015-03-06T13:53:47	2015-03-06T13:58:01	2015-03-06T14:08:35	00:10:34	Not available	File Info	FINISHED AT 2015-03-06 14:16:18
19	finished	0	T1_US_FNAL	1	2015-03-06T13:53:47	2015-03-06T13:58:00	2015-03-06T14:04:12	00:06:12	Not available	File Info	FINISHED AT 2015-03-06 14:10:43
20	finished	0	T2_BE_IJHE	1	2015-03-06T13:53:47	2015-03-06T13:57:47	2015-03-06T14:04:54	00:07:07	Not available	File Info	FINISHED AT 2015-03-06 14:11:25
21	finished	0	T2_US_Florida	1	2015-03-06T13:53:47	2015-03-06T13:58:00	2015-03-06T13:59:37	00:01:37	Not available	File Info	FINISHED AT 2015-03-06 14:05:36
22	finished	0	T2_UK_London_IC	1	2015-03-06T13:53:47	2015-03-06T13:57:50	2015-03-06T14:08:15	00:10:25	Not available	File Info	FINISHED AT 2015-03-06 14:10:47
23	finished	0	T2_US_Florida	1	2015-03-06T13:53:47	2015-03-06T13:59:59	2015-03-06T14:14:15	00:14:16	Not available	File Info	FINISHED AT 2015-03-06 14:19:55
24	finished	0	T1_US_FNAL	1	2015-03-06T13:53:47	2015-03-06T14:01:43	2015-03-06T14:14:52	00:13:09	Not available	File Info	FINISHED AT 2015-03-06 14:15:11
25	finished	0	T1_US_FNAL	1	2015-03-06T13:53:47	2015-03-06T14:01:42	2015-03-06T14:15:13	00:13:31	Not available	File Info	FINISHED AT 2015-03-06 14:19:18
Id	Status	AppExitCode	Site	Retries	Submitted	Started	Finished	Wall Time	Job Log	File Access	FTS File Status
Showing 1 to 25 of 300 entries											
										First	Previous Page 1 of 12 Next Last

Integration

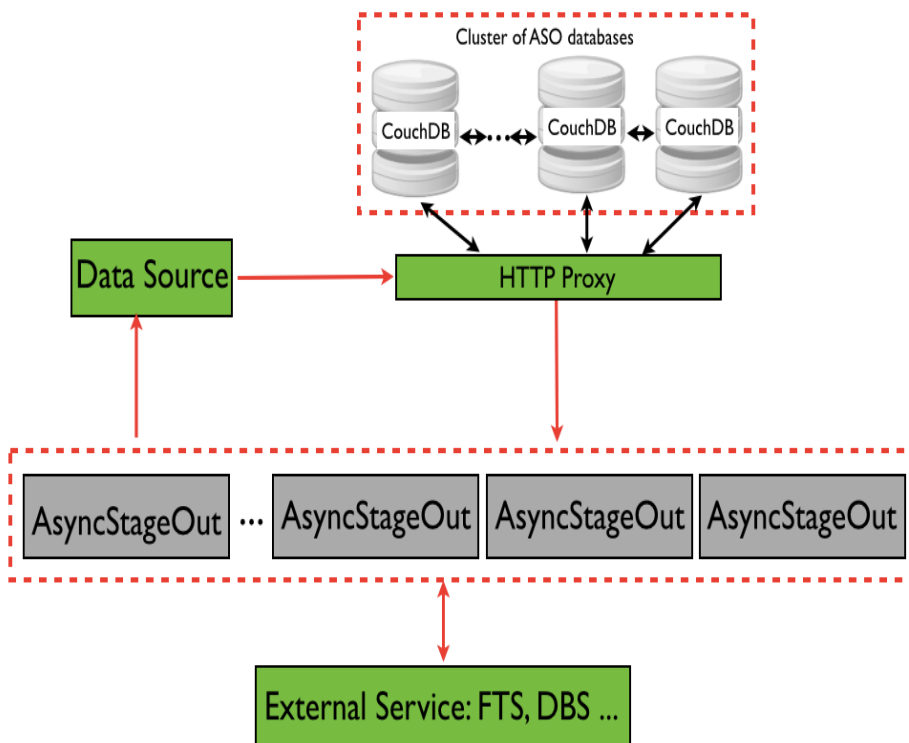


Deployment models

Production model



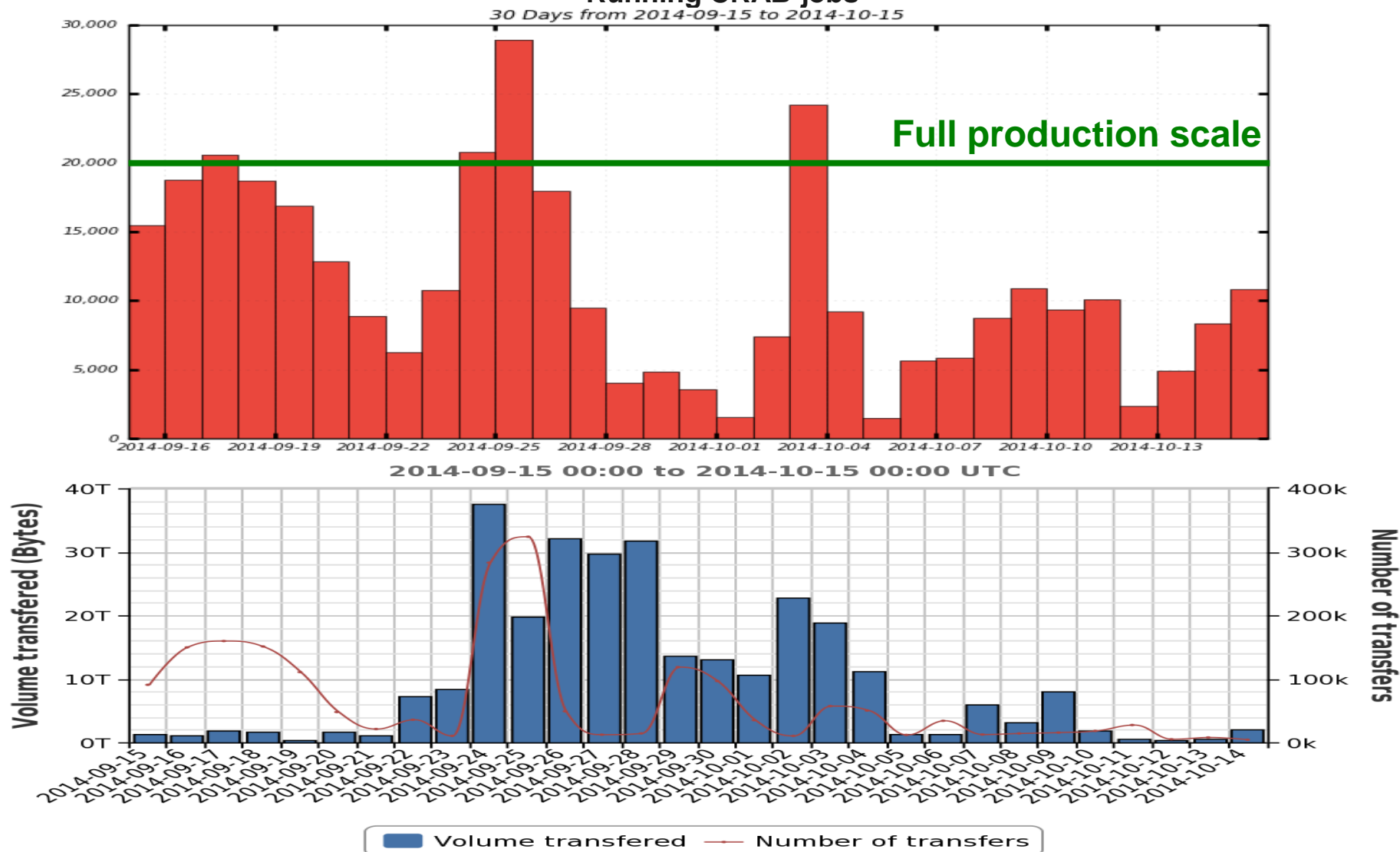
Highly scalable/available model based on CouchDB v2.0



CSA14 exercise



Running CRAB jobs

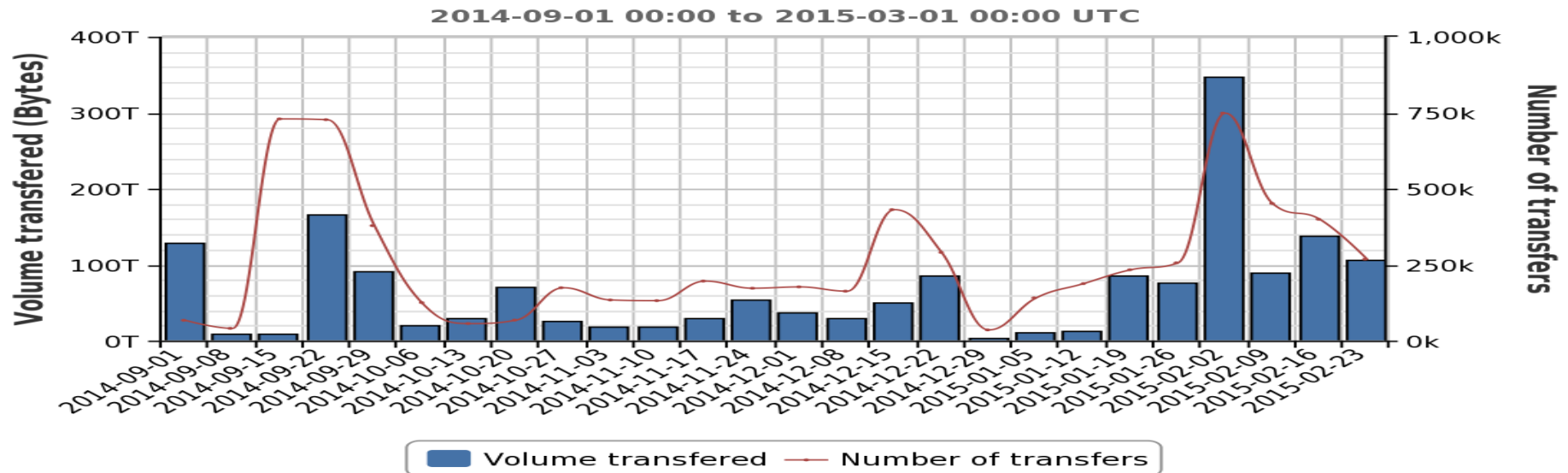


Production results

- ASO has started production in June 2014
 - More than 2 PB/7 M files transferred during the last 6 months
 - Peaks of more than 700k transfers per week



Volume transfered / Number of transfers (cms)



Scale tests

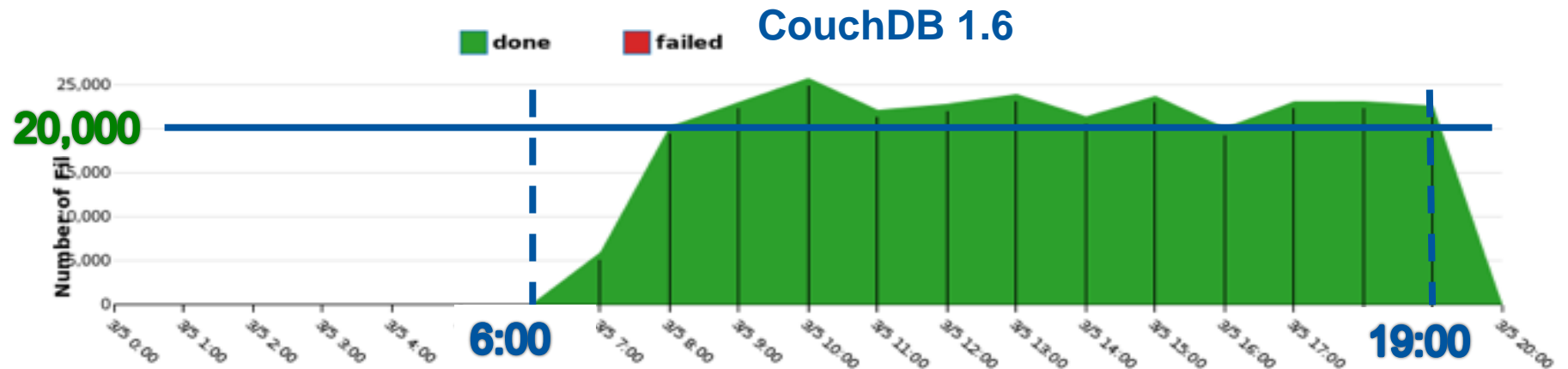
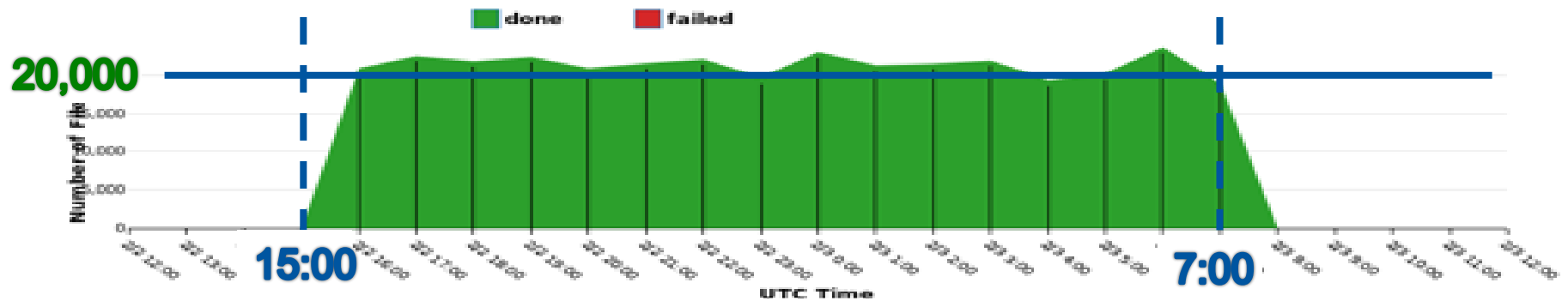
- Goal:
 - Simulate data transfers using PhEDEx LifeCycle Agent
 - Explore the scalability limits of ASO by scaling-up to 2-4 times the design load within the production and the new version of CouchDB
 - The design load is ~ 200 k completed analysis jobs/day
 - ~ 300 k completed files/day (~ 200k outputs + ~ 100k logs)
- Configuration:
 - 200 parallel users
 - 100 files per FTS job
 - 60 sites
 - File size: 2 GB
 - Transfer throughput: 100 MB/s
 - Run the system for more than 12 hours

2 times scale

600 k completed files/day
CouchDB 1.1 (production)

Injection starts

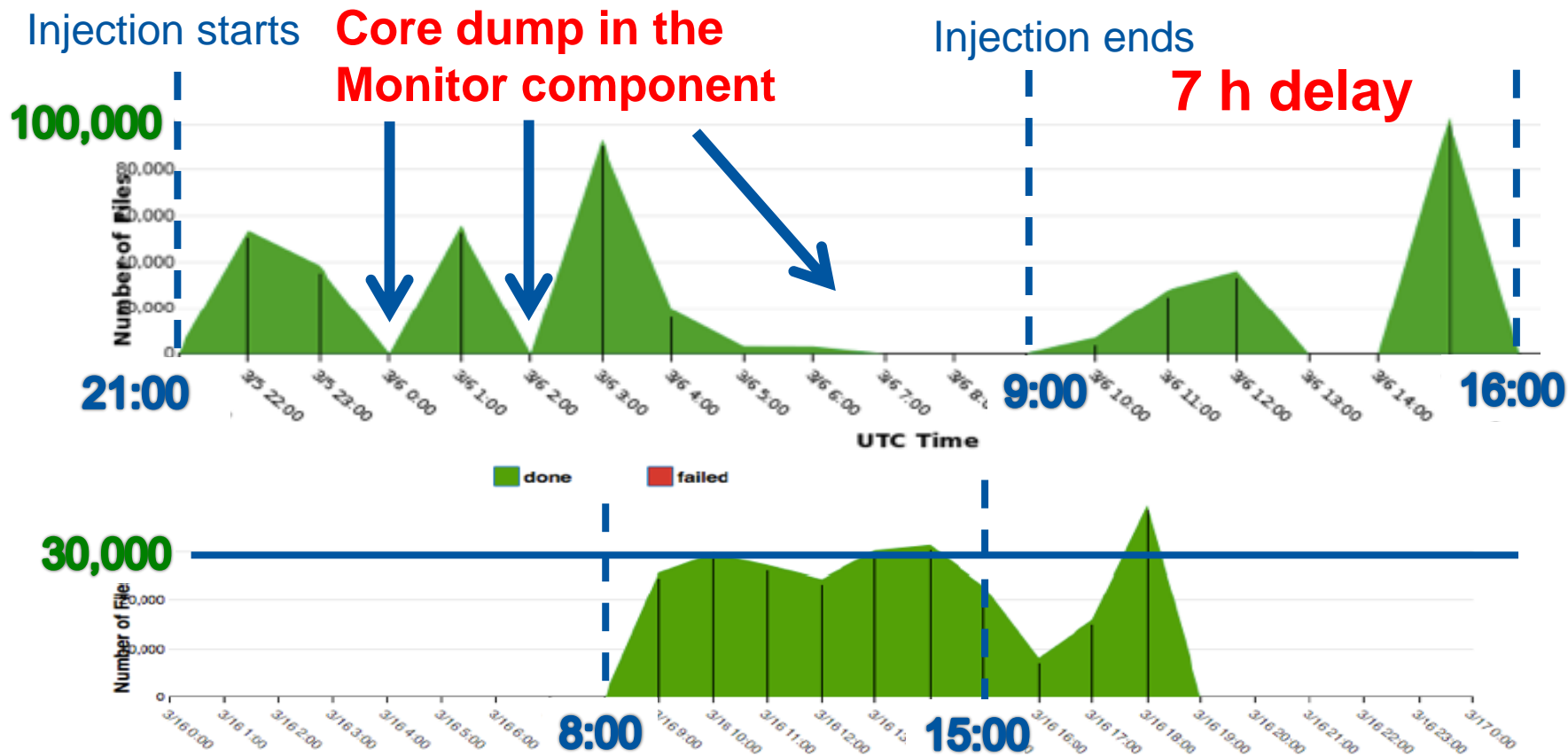
Injection ends



- ✓ Independently of the version of CouchDB, ASO can manage 2 times scale the design load

4 times scale

1,2 M completed files/day



- ✓ ASO can manage ~ 3 times scale the design load with tuned parameters

Summary and conclusions

- ASO has evolved from a limited prototype to a highly adaptable and scalable service
- ASO has shown good performance during commissioning and production
- ASO can manage 2 times scale of the design load
 - The management of 4 times scale is possible with an accurate tune of ASO and system parameters
- Re-use of design and components from PhEDEx point the way to a more modular architecture for data-management tools in CMS
 - Long-term maintainability, performance and adaptability