



# Dynamo

## Intelligent Data Management

**Benedikt Maier for the Dynamo Team**

Dan Abercrombie, Max Goncharov, Benedikt Maier, Sid Narayanan, Christoph Paus

*March 27, 2018*

# First words

Dynamo is ...

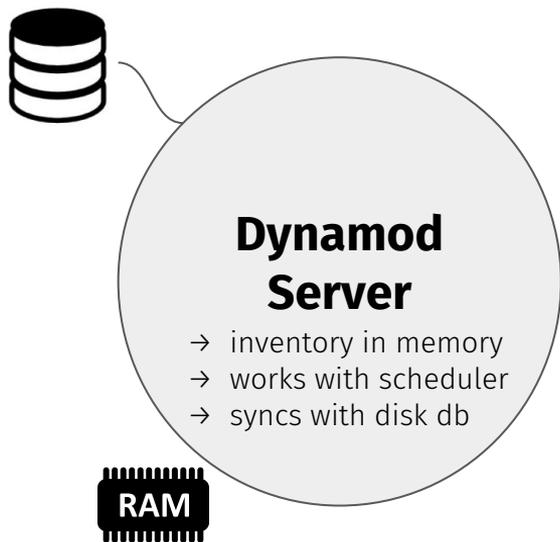
... an intelligent framework that manages data to facilitate easy access to the popular data by researchers,

... a collection of agents that can be scheduled to run in any configurable order. New agents can be easily defined and included into the framework,

... written in python,

... available at <https://github.com/SmartDataProjects>.

## Core element: the inventory



### **Inventory**

Information about placement of all datasets and blocks (and files) on disks and tape → many millions of objects, all sitting in memory

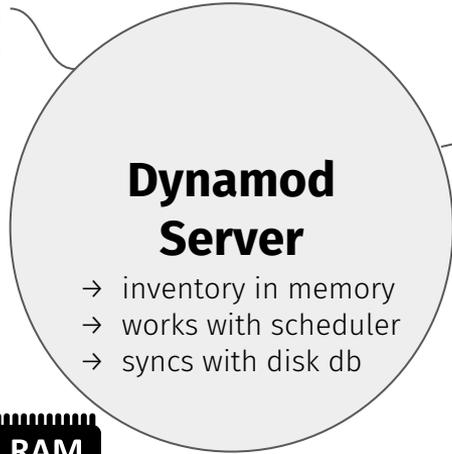
### **Scheduler**

Handles any task that queries or wants to change inventory status (grants read/write access)

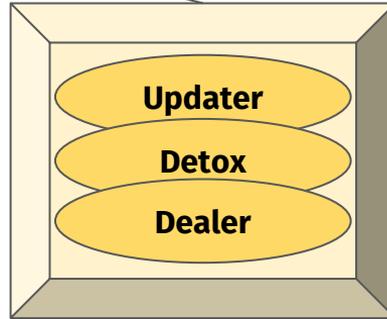
### **MySQL backend**

Whenever a write-request changed inventory status, the change is mirrored to persistent storage db

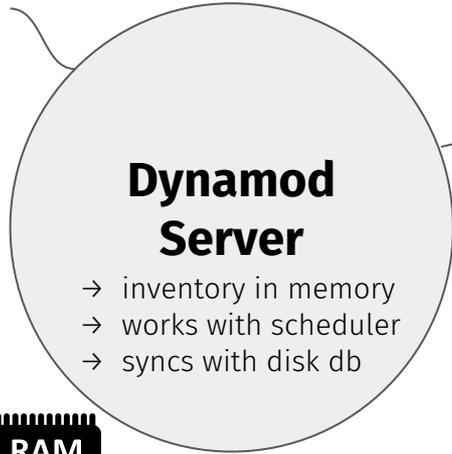
# Design



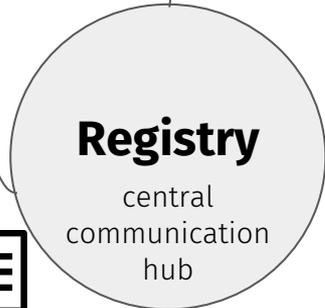
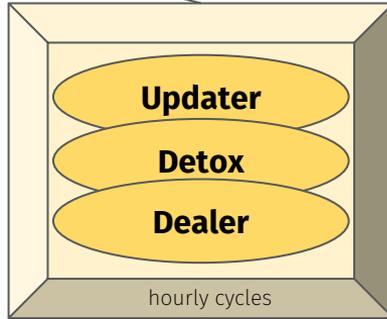
**scheduler**  
read-only  
read/write



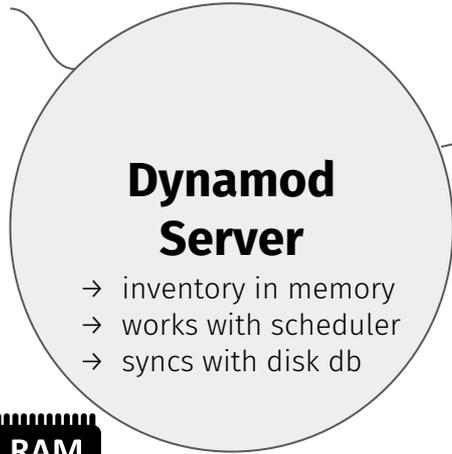
# Design



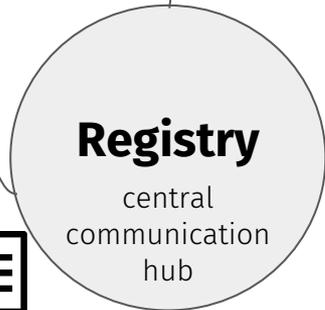
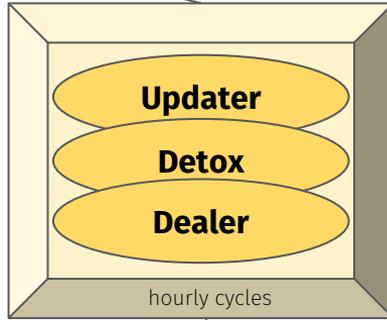
**scheduler**  
read-only  
read/write



# Design



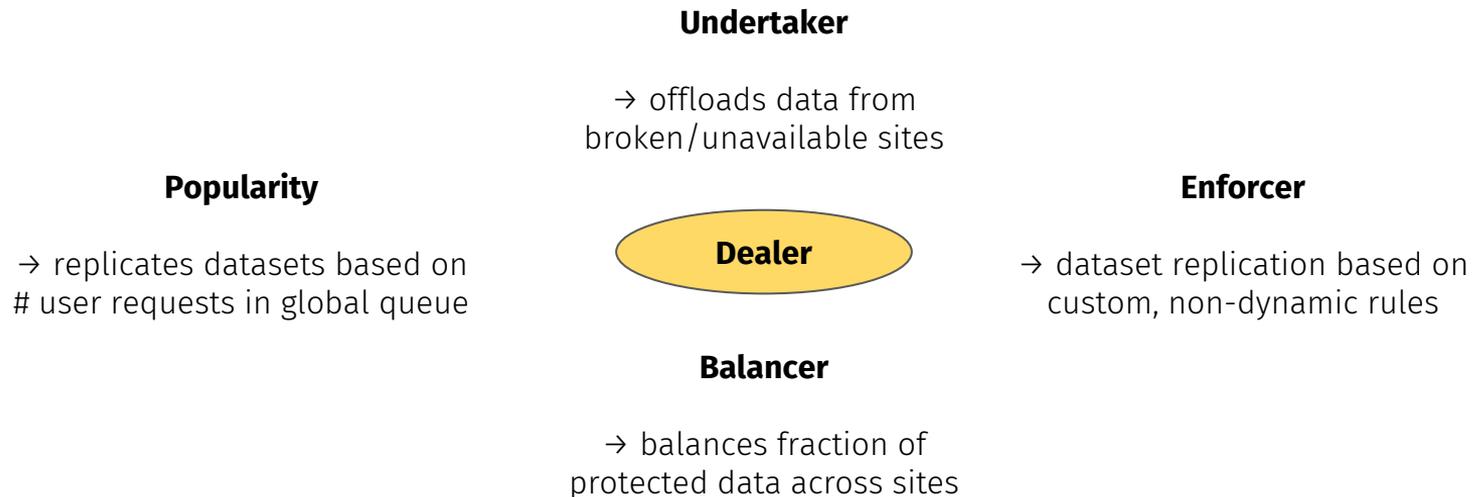
**scheduler**  
read-only  
read/write



**FTS**



# Dealer plugins



These all make decisions of the type: “Make a replica of dataset X at site Y”

# Dealer monitoring

Aggregated copied/missing volume



Individual datasets to sites

## T2\_FR\_GRIF\_IRFU

Click Request ID for corresponding PhEDEx API.  
Click Replica Name for corresponding progress graph.  
Red cells signify stuck transfers. (<math>\lt; 1\%</math> copied in past 5 days)  
Number of Replicas: 329

Request ID:	Replica Name:	Copied (TB):	Total (TB):
1139612	/SingleMuonRun2017D-v1/RAW	46.7162	46.7517
1146121	/MuOnia1/n2012-22Jan2013-v1/AOD	0	19.8384
1146121	/DoubleEGRun2016B-16Apr2017_ver2-v1/MNIAOD	0	3.4458
1144555	/MinBias_TuneCUETPBM1_13TeV-pythia8/Run1Spring15DR74-NoPURealisticRecodebug_741_p1_mcRun2_Realistic_50ns_v0-v2/GEN-SIM-RECODEBUG	27.6836	44.6587
1144686	/WJetsToLNu_TuneCUETPBM1_13TeV-amcatnloFXFX-pythia8/Run1Fall15DR76-PU25nsData2015v1_76X_mcRun2_asymptotic_v12_ext2-v1/AODSIM	32.8776	39.491
1145655	/HMinimumBias7H/Run2015-02May2016-v1/AOD	16.6338	31.4865
1145755	/WJetsToLNu_Pt80ToInf_TuneCUETPBM1_13TeV-amcatnloFXFX-pythia8/Run1Spring16MiniAODv2-PUSpring16_80X_mcRun2_asymptotic_2016_miniAODv2_v0-v1/MNIAODSIM	0	0.0325
1145755	/DMV_NNPDF30_Axial_Mpfi-2500_Mchi-1_gSM-0p25_gDM-1p0_v2_13TeV-powheg/Run1Summer16DR80Premix-PUMorond17_80X_mcRun2_asymptotic_2016_TracheIV_v6-v1/AODSIM	0.0831	0.0942
1145755	/VBFHToTauTau_M125_13TeV_powheg_pythia8/Up/Run1Spring16MiniAODv2-PUSpring16RAWAODSIM_reHLT_80X_mcRun2_asymptotic_v14-v1/MNIAODSIM	0.0333	0.0432
1145755	/QstarToJJ_M_9000_TuneCUETPBM1_13TeV_pythia8/Run1Summer16DR80Premix-PUMorond17_80X_mcRun2_asymptotic_2016_TracheIV_v6-v1/AODSIM	0.0262	0.0474
1145755	/VBFHHTo2B2G_CV_1_C2V_1_C3_0_13TeV-madgraph/Run1Summer16DR80Premix-PUMorond17_80X_mcRun2_asymptotic_2016_TracheIV_v6-v1/AODSIM	0.0062	0.0342
1145755	/SeesawTypeII_SIGMAplusSIGMAminusWW_M-760_13TeV-madgraph-pythia8/Run1Summer16DR80Premix-PUMorond17_80X_mcRun2_asymptotic_2016_TracheIV_v6-v1/AODSIM	0.0022	0.0311
1145755	/SingleMuonRun2016C-MuAIOverlaps-07Aug17-v1/ALCARECO	0.0534	0.249
1145755	/WJetsToLNu_TuneCUETPBM1_13TeV-madgraphMLM-pythia8/Run1Summer16MiniAODv2-FlatPU28to2HcalNZSRW_80X_mcRun2_asymptotic_2016_TracheIV_v6-v1/MNIAODSIM	0.0139	0.045
1145755	/HeavyNeutrino_trilepton_M-10_V-0p01_2_NLO/Run1Summer16DR80Premix-PUMorond17_80X_mcRun2_asymptotic_2016_TracheIV_v6-v1/AODSIM	0.0005	0.0287
1145755	/WprimeToTB_plusSM_T0Lco_M-3100_LH_TuneCUETPBM1_13TeV-comphep-pythia8/Run1Summer16DR80Premix-PUMorond17_80X_mcRun2_asymptotic_2016_TracheIV_v6-v1/AODSIM	0	0.0628
1145755	/JTGamma_SingleLeptFromT_TuneCUETPBM2T4_13TeV-amcatnlo-1srdown-pythia8-pythia8/Run1Summer16MiniAODv2-PUMorond17_80X_mcRun2_asymptotic_2016_TracheIV_v6-v2/MNIAODSIM	0.041	0.2239
1145755	/SinglePhoton_Pt-200/PhaseII/TDRSpring17DR-PU140CaloAging4500Ultimate_91X_upgrade2023_realistic_v3-v2/GEN-SIM-RECO	0.0351	0.1197
1145755	/ZH_HTO20_270A8_M120_13TeV_powheg_pythia8/Run1Summer16DR80Premix-PUMorond17_80X_mcRun2_asymptotic_2016_TracheIV_v6-v1/AODSIM	0	0.0317

Click on a bar

# Enforcer

→ dataset replication based on custom, non-dynamic rules

*e.g. US-CMS decision to have one replica of widely used datasets at US sites*

Configurable via JSON

```
"usa_miniaod": {  
  "datasets": ["/**/MINIAOD*"],  
  "num_copies": 1,  
  "sites": ["T1_US_FNAL_Disk", "T2_US_*"]  
}
```



# Enforcer

→ dataset replication based on custom, non-dynamic rules

*e.g. US-CMS decision to have one replica of widely used datasets at US sites*

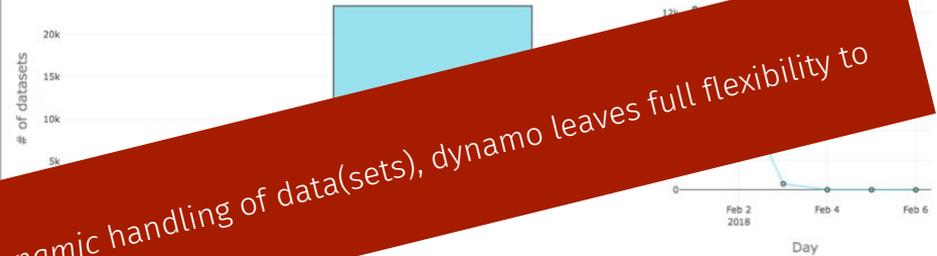
Configurable via JSON

```
"us": {  
  "datasets": ["/*/*/MINIAOD*"],  
  "num_copies": 1,  
  "sites": ["T1_US_FNAL_Disk", "T2_US_*"]  
}
```

## Enforcer

Choose rule

Rule: usa\_miniaod



Considered sites



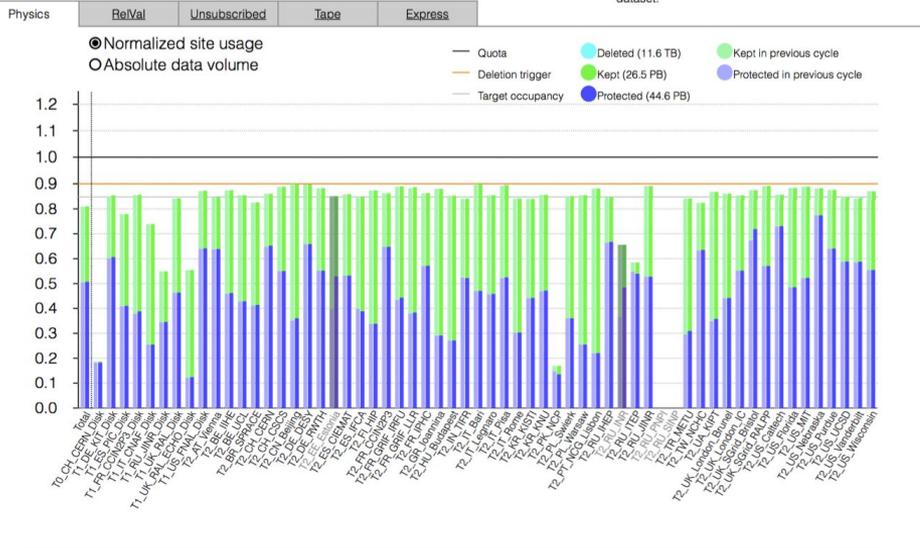
Take-away message:  
Despite the underlying paradigm being the dynamic handling of data(sets), dynamo leaves full flexibility to impose custom rules.

# Detox deletion results

Cycle 15118 (Policy v62, 2018-03-15 04:05:12)

Cycle: Previous Next  Go  
PhEDEx request:  Go

Find dataset:  Search



## Get details:

- [Download deletion list](#)
- [PhEDEx requests page](#)

```
### Deletion trigger
When site.occupancy > 0.9
Until site.occupancy < 0.85
```

```
### Replica protection / deletion policies
Ignore site.name == T*_MSS
```

```
ProtectBlock blockreplica.is_locked
...
ProtectBlock blockreplica.is_last_transfer_source
...
Dismiss dataset.usage_rank > 200
...
# Default decision
Dismiss
```

## Dynamo deletion planning

You are currently using this tool as user 'bmaier'.

Email address results should be sent to:

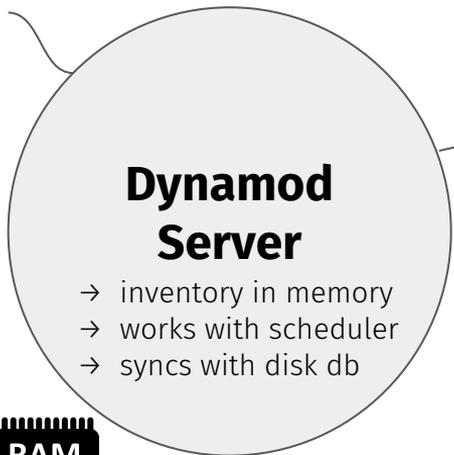
Policy file (must be .txt)

No file selected

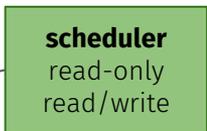
## Detox

- Everything is based on a human-readable policy file. In principle, anything that is not explicitly protected can be deleted.
- Virtually any DatasetReplica/BlockReplica attribute can be defined/added and used to manage deletions.
- On top of dynamic deletions, policy files for dedicated deletion campaigns can easily be spelled out (and tested).
- The entire history of deletions is automatically saved so we can always say why something got deleted where and when.

# Design



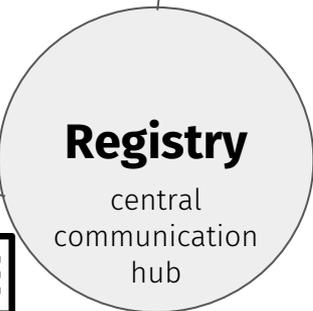
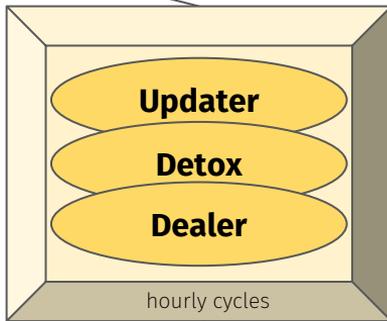
monitoring



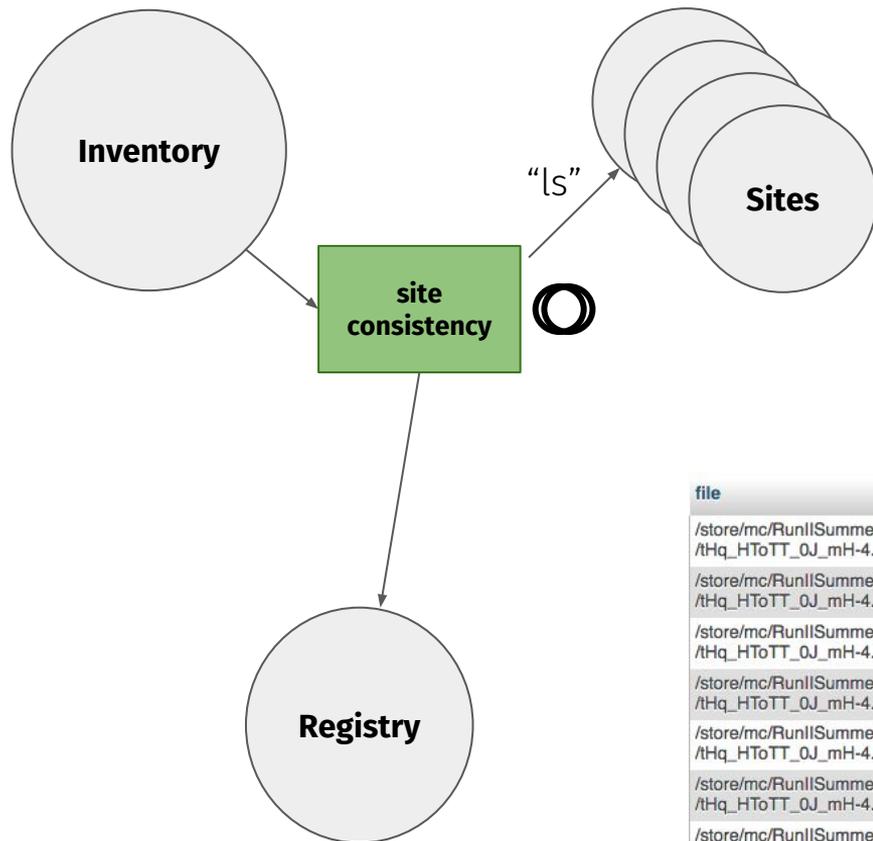
script upload  
web interface



command line  
analysis



# Site consistency



Aligns what should be on the site with what *is* on the sites.

(Before, this was done not very frequently for T0/T1, and only occasionally for T2s)

**Important to avoid stuck transfers, problems with user access, or wasted space.**

→ catalog obtained from inventory

→ comparing this against results from remotely probing site content via xrd fs (10min - 2hrs per site)

→ Ingests “missing” files to be transferred and “orphan” files to be deleted into registry, also publishes ASCII files in [web](#)

file	site_from	site_to	status	created
/store/mc/RunIIISummer16MiniAODv2/Hq_HToTT_0J_mH-4...	T2_FI_HIP	T1_UK_RAL_Disk	new	2018-02-14 22:06:37
/store/mc/RunIIISummer16MiniAODv2/Hq_HToTT_0J_mH-4...	T1_US_FNAL_Disk	T1_UK_RAL_Disk	new	2018-02-14 22:06:37
/store/mc/RunIIISummer16MiniAODv2/Hq_HToTT_0J_mH-4...	T1_FR_CCIN2P3_Disk	T1_UK_RAL_Disk	new	2018-02-14 22:06:37
/store/mc/RunIIISummer16MiniAODv2/Hq_HToTT_0J_mH-4...	T2_UA_KIPT	T1_UK_RAL_Disk	new	2018-02-14 22:06:37
/store/mc/RunIIISummer16MiniAODv2/Hq_HToTT_0J_mH-4...	T2_FI_HIP	T1_UK_RAL_Disk	new	2018-02-14 22:06:37
/store/mc/RunIIISummer16MiniAODv2/Hq_HToTT_0J_mH-4...	T1_US_FNAL_Disk	T1_UK_RAL_Disk	new	2018-02-14 22:06:37
/store/mc/RunIIISummer16MiniAODv2/Hq_HToTT_0J_mH-4...	T1_FR_CCIN2P3_Disk	T1_UK_RAL_Disk	new	2018-02-14 22:06:37

# Site consistency

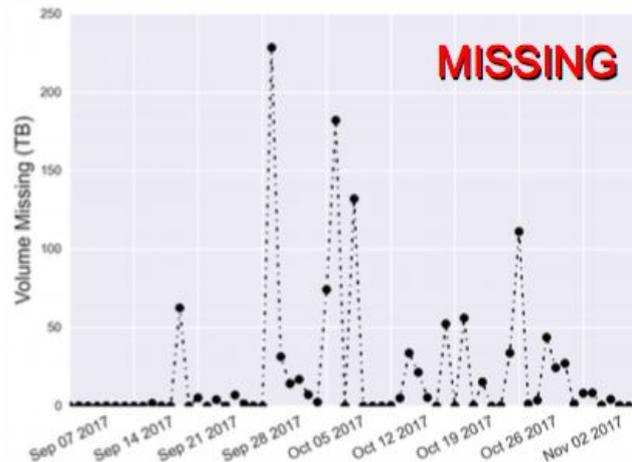
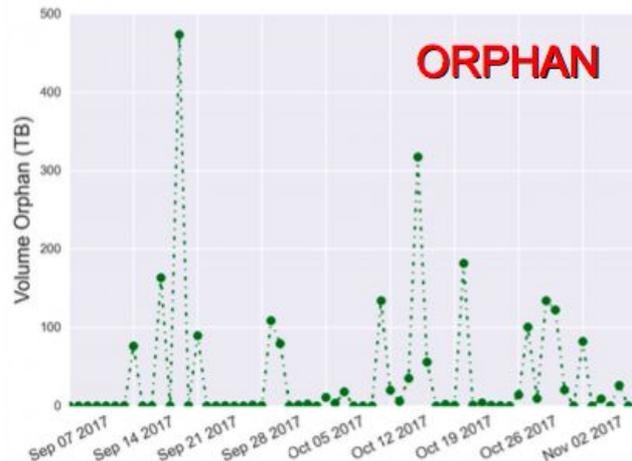
In the roll-out phase last fall, we

- deleted O(1 PB) orphan files
- recovered O(1 PB) missing files
- removed several million of empty directories that previous deletion processes left behind

Spikes in the right plots appear when we probed a new site for the first time.

Now, orphan/missing files appear occasionally.

Constitutes the very first use case for file transfers fully operated by dynamo.



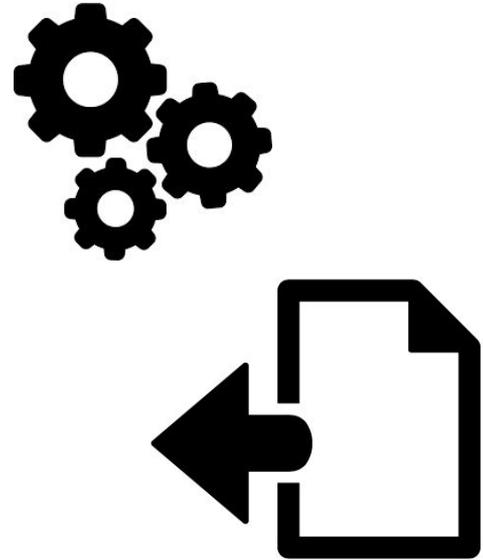
# Consistency Check Summary

Debugged sites	Sites that need debugged or hand checking	All sites
----------------	---	-----------

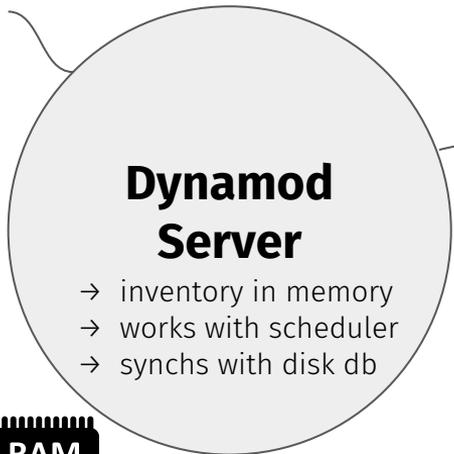
Last Update	Site Name	Time [min]	Number Files	Number Nodes	Unlisted	Empty Nodes	Num Missing	No Source	Size [TB]	Num Orphan	Size [TB]
<a href="#">2018-03-14 22:24:29</a>	<a href="#">T2_TW_NCHC</a>	53	<a href="#">128536</a>	18866	0	10	<a href="#">0 [blocks]</a>	<a href="#">0</a>	0.00	<a href="#">0</a>	0.00
<a href="#">2018-03-14 21:39:20</a>	<a href="#">T2_FR_GRIF_IRFU</a>	68	<a href="#">354694</a>	29003	0	2804	<a href="#">0 [blocks]</a>	<a href="#">0</a>	0.00	<a href="#">0</a>	0.00
<a href="#">2018-03-14 20:49:20</a>	<a href="#">T1_DE_KIT_Disk</a>	47	<a href="#">705434</a>	39063	0	26	<a href="#">0 [blocks]</a>	<a href="#">0</a>	0.00	<a href="#">0</a>	0.00
<a href="#">2018-03-14 19:39:17</a>	<a href="#">T2_AT_Vienna</a>	8	<a href="#">88548</a>	1031	0	13	<a href="#">0 [blocks]</a>	<a href="#">0</a>	0.00	<a href="#">0</a>	0.00
<a href="#">2018-03-14 17:53:17</a>	<a href="#">T2_IT_Rome</a>	22	<a href="#">171489</a>	9468	0	478	<a href="#">4559 [blocks]</a>	<a href="#">1746</a>	12.88	<a href="#">0</a>	0.00
<a href="#">2018-03-14 17:51:59</a>	<a href="#">T2_US_Purdue</a>	81	<a href="#">618567</a>	84345	0	7734	<a href="#">0 [blocks]</a>	<a href="#">0</a>	0.00	<a href="#">0</a>	0.00
<a href="#">2018-03-14 17:07:41</a>	<a href="#">T2_DE_DESY</a>	96	<a href="#">1349167</a>	73689	0	4743	<a href="#">0 [blocks]</a>	<a href="#">0</a>	0.00	<a href="#">0</a>	0.00
<a href="#">2018-03-14 14:04:31</a>	<a href="#">T2_BR_SPRACE</a>	33	<a href="#">273535</a>	22082	0	3590	<a href="#">0 [blocks]</a>	<a href="#">0</a>	0.00	<a href="#">0</a>	0.00
<a href="#">2018-03-14 12:50:56</a>	<a href="#">T2_UK_SGrid_Bristol</a>	20	<a href="#">109198</a>	8117	0	541	<a href="#">0 [blocks]</a>	<a href="#">0</a>	0.00	<a href="#">0</a>	0.00
<a href="#">2018-03-14 11:46:31</a>	<a href="#">T2_PK_NCP</a>	15	<a href="#">17592</a>	1964	0	23	<a href="#">0 [blocks]</a>	<a href="#">0</a>	0.00	<a href="#">1966</a>	2.40
<a href="#">2018-03-14 10:43:04</a>	<a href="#">T2_FR_CCIN2P3</a>	12	<a href="#">97088</a>	7140	0	219	<a href="#">0 [blocks]</a>	<a href="#">0</a>	0.00	<a href="#">0</a>	0.00
<a href="#">2018-03-14 10:39:47</a>	<a href="#">T2_PL_Warsaw</a>	68	<a href="#">85320</a>	19548	0	1798	<a href="#">16464 [blocks]</a>	<a href="#">2039</a>	51.57	<a href="#">0</a>	0.00
<a href="#">2018-03-14 07:55:46</a>	<a href="#">T2_UK_SGrid_RALPP</a>	24	<a href="#">201675</a>	21943	0	1487	<a href="#">0 [blocks]</a>	<a href="#">0</a>	0.00	<a href="#">0</a>	0.00
<a href="#">2018-03-14 07:02:19</a>	<a href="#">T2_US_Florida</a>	31	<a href="#">636861</a>	62861	0	4454	<a href="#">2 [blocks]</a>	<a href="#">2</a>	0.01	<a href="#">0</a>	0.00
<a href="#">2018-03-14 06:11:53</a>	<a href="#">T2_FR_GRIF_LLRL</a>	41	<a href="#">228493</a>	22369	0	4695	<a href="#">10 [blocks]</a>	<a href="#">10</a>	0.02	<a href="#">1</a>	0.00

# File operations

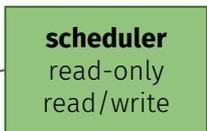
- Module developed in tandem with Site Consistency
  - Deletion: immediate action with `gfal-rm`
  - Transfer: form FTS job and submit to CERN FTS
- Disk-to-disk transfers available since last year
- Functionality to stage from tape established last month



# Design



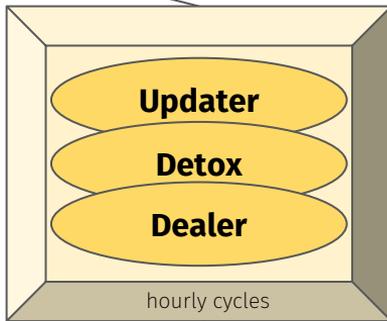
monitoring



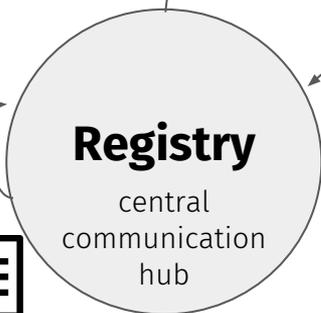
script upload  
web interface



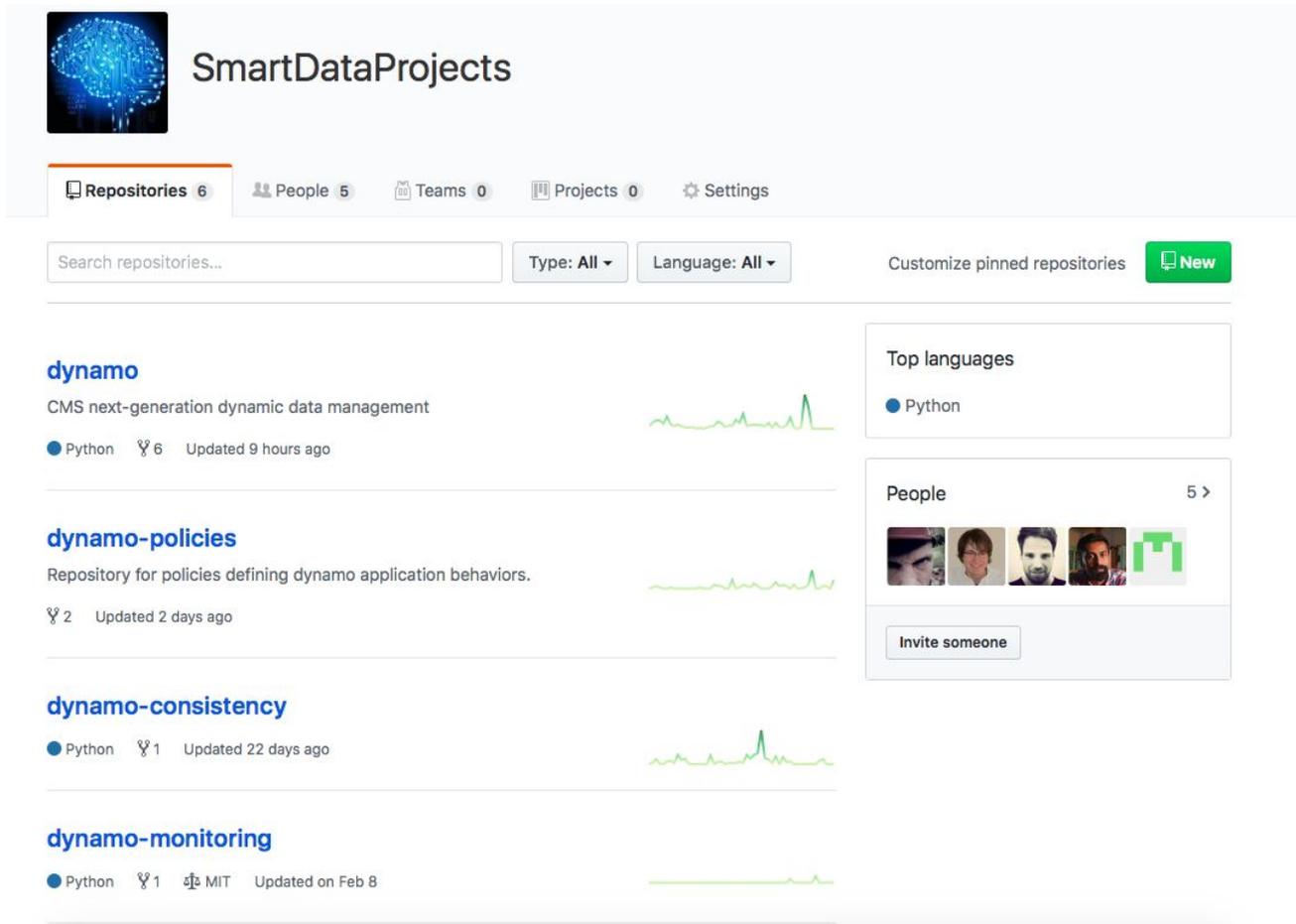
command line  
analysis



register data  
(production)



# Dynamo ecosystem



**SmartDataProjects**

Repositories 6 | People 5 | Teams 0 | Projects 0 | Settings

Search repositories... | Type: All | Language: All | Customize pinned repositories | [New](#)

**dynamo**  
CMS next-generation dynamic data management  
Python 6 Updated 9 hours ago

**dynamo-policies**  
Repository for policies defining dynamo application behaviors.  
2 Updated 2 days ago

**dynamo-consistency**  
Python 1 Updated 22 days ago

**dynamo-monitoring**  
Python 1 MIT Updated on Feb 8

**Top languages**  
Python

**People** 5 >  
[Profile pictures of contributors]  
[Invite someone](#)

## Ongoing work & next steps

- Porting dynamo from bash-script-philosophy to server mode that is able to handle tasks asynchronously was the big achievement late last year.
  - Much more stable
  - Much faster
  - Much easier to interface dynamo with other services
- Implementation of REST APIs #1 priority to query inventory status and allow any tool to “dock on”
- Phasing out PhEDex as underlying tool for handling transfers & deletions
- Working on documentation to make it easy for other experiments to install and use dynamo
- Continuous integration and unit testing will come soon (right now development model consists of developing and testing on dedicated server and git-reviewing among collaborators)