



# Moving the California distributed CMS XCache from bare metal into containers using Kubernetes

Edgar Fajardo<sup>1</sup>, presented by Matevž Tadel<sup>1</sup>  
Justas Balcas<sup>2</sup>, Alja Tadel<sup>1</sup>, Frank Würthwein<sup>1</sup>, Diego Davila<sup>1</sup>,  
Jonathan Guiang<sup>1</sup>, Igor Sfiligoi<sup>1</sup>

Caltech<sup>2</sup>, UCSD<sup>1</sup>

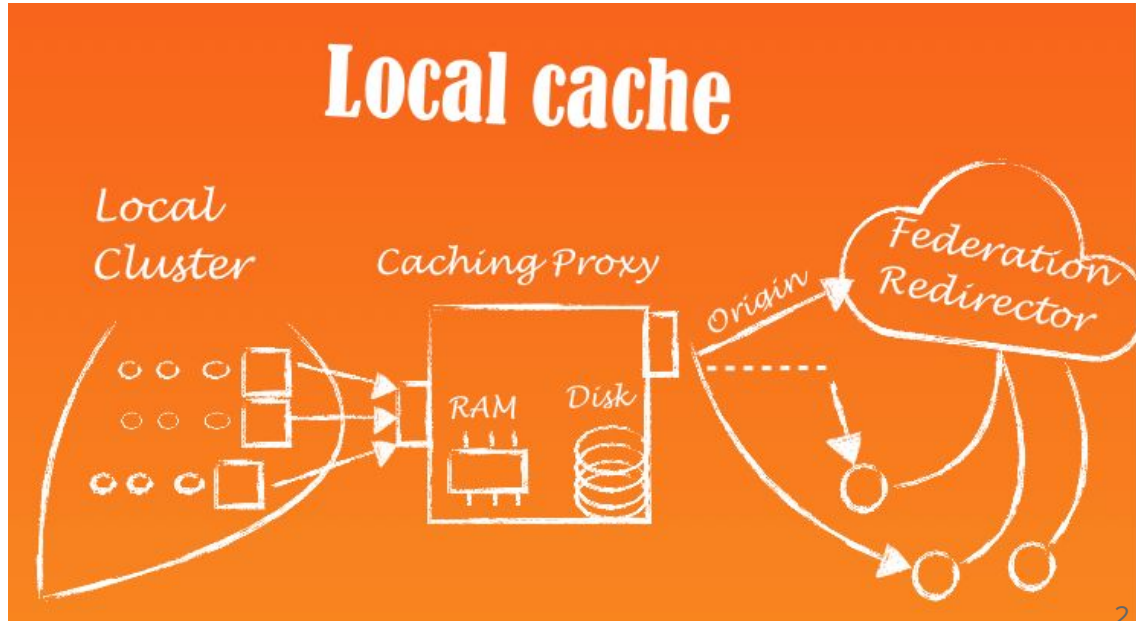
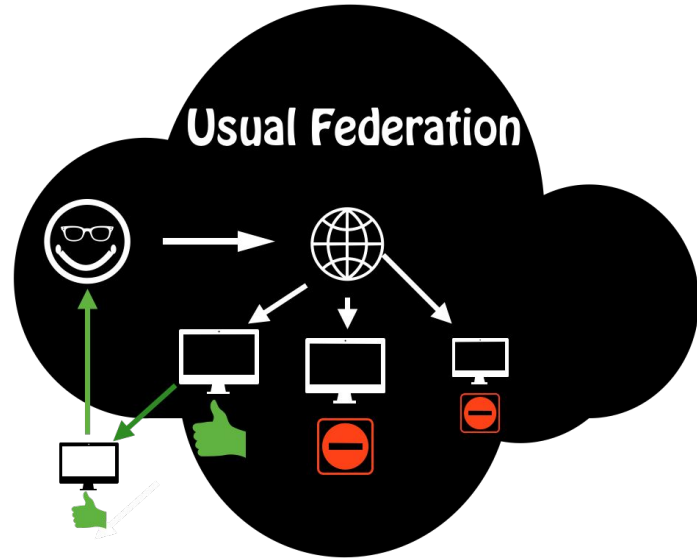
UC San Diego

**SDSC** SAN DIEGO  
SUPERCOMPUTER CENTER

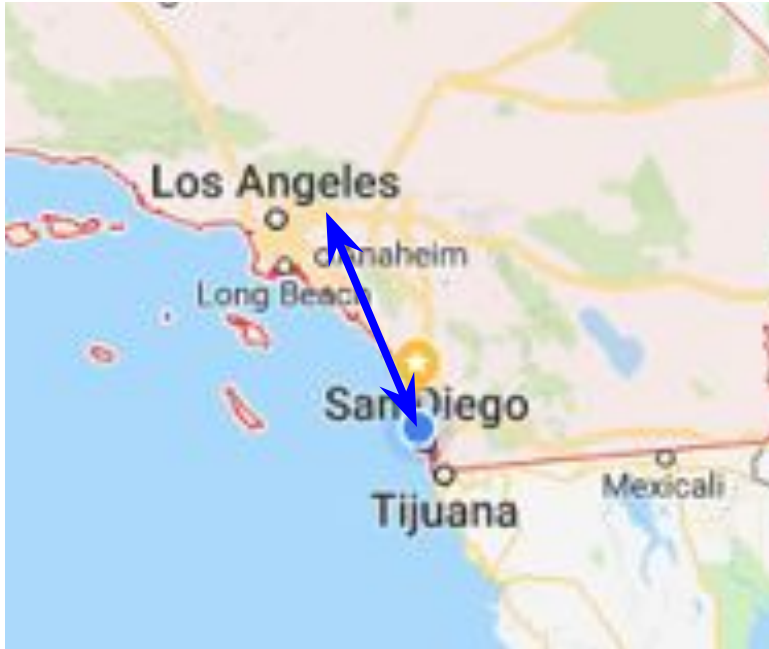


# What is an XRootD Cache? Caching Cluster!

Each node acts as an independent cache. But they can be clustered through cmsd to scale horizontally.



# Opportunity for two sites to merge some namespace and profit from closeness



UCSD ↔ Caltech link:

- 120 miles
- 100 Gbit/sec
- below 3 ms

# Recipe: How to setup an XCache for CMS at a given site(s)?

1. Decide on the namespace to cache. (We used /MINIAOD\*)
2. Calculate the working set (the set of unique files of the namespace that are accessed on a given period).
3. Provision the storage for some fraction of data from step #2.
4. Install and configure XCache on top of disk from step #3.



↑  
This is where kubernetes kicks in



# 1. Decide on the namespace

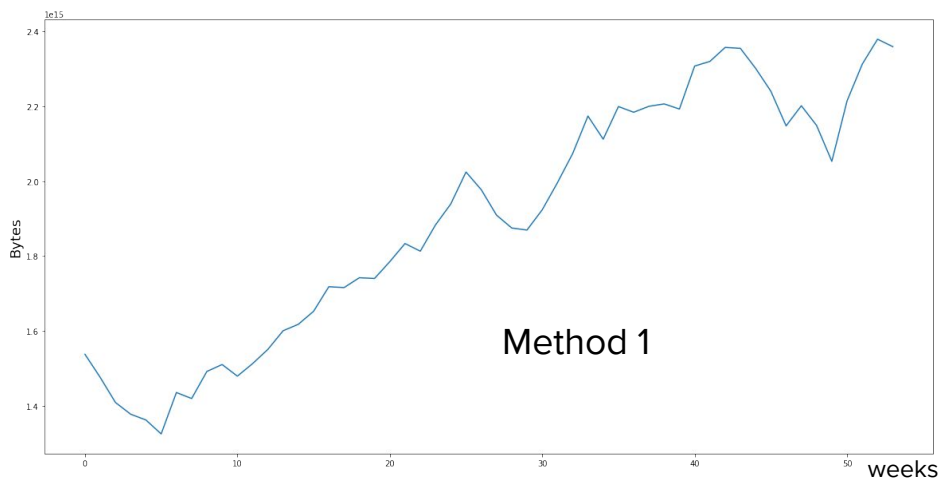
Dataset	Size (PB)
/*Run2016*-03Feb2017*/MINIAOD	0.182
/*RunII Summer16MiniAODv2-PUMoriond17_80X_*/ MINIAODSIM	0.5
/*/*RunIIFall17MiniAODv2*/MINIAODSIM	0.2
/*/*-31Mar2018*/MINIAOD	0.14
Total	1.04
/*/*/MINIAOD	2.92
/*/*/MINIAODSIM	4.6
Total	7.52

Until end of  
september 2019  
we had a more  
restricted  
namespace.

Now with more  
hardware we  
relaxed the  
constraints.

## 2. Estimating the working set for SoCal

Working set size MINI\* Window = 4 weeks, Time period: 2018-06-21 - 2019-06-27



### Method 1:

1. Look at the unique MINI\* **data-sets** accessed globally (at all sites) within a four week window and calculate their size.
2. Move the window 1 week at a time for a year worth of data from the Global pool ClassAds
3. Results: The **monthly working set** is somewhere between 1.6PB and 2.4 PB

### Method 2:

1. Look at the unique MINIAOD\* files accessed in SoCAL during the month of October.
2. Estimate the month working set as the size of all the unique files accessed during October
3. Results: 451TB

### 3. Provision disk infrastructure based on the needs in step #2.

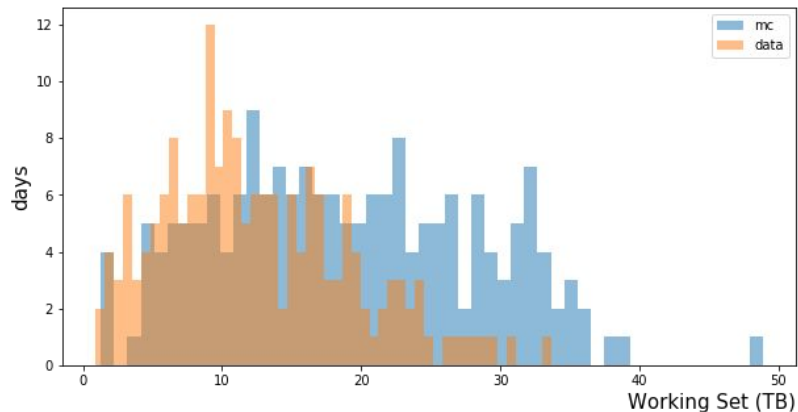
	UCSD	Caltech
<b>Nodes</b>	11 (+1 JBOD)	2*
<b>Disk Capacity per node</b>	12 x 2TB = 24TB (+ 48 x 11TB)	30 x 6TB (HGST Ultrastar 7K6000)
<b>Network Card per node</b>	10 Gbps (+ 40 Gbps)	40 Gbps
<b>Total Disk Capacity</b>	264 TB (+ 528 TB) = <b>792 TB</b>	<b>360 TB*</b>
<b>TOTAL</b>	792 TB + 360 TB* = <b>1,152 TB</b>	

\* Caltech has 1 JBOD (440TB) ready to be added anytime (currently in HDFS managed space)



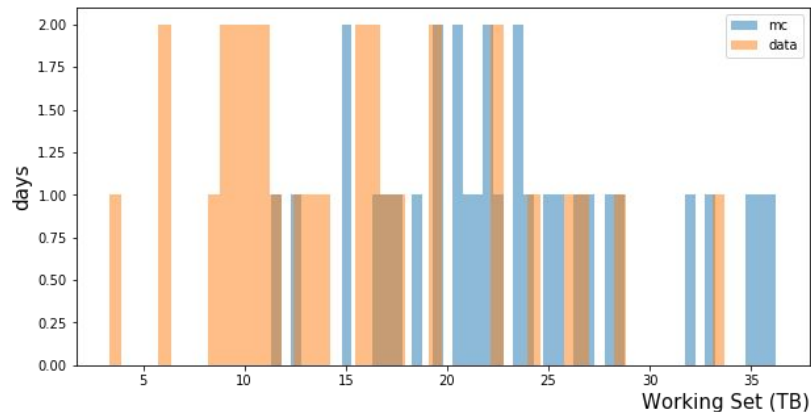
# 3. Provision disk infrastructure (Working set for SoCal)

Daily Working Set (Apr - Sept, 2019)



- The daily working set for SoCal was at most 80 TB = 48 TB MC + 32 TB Data.

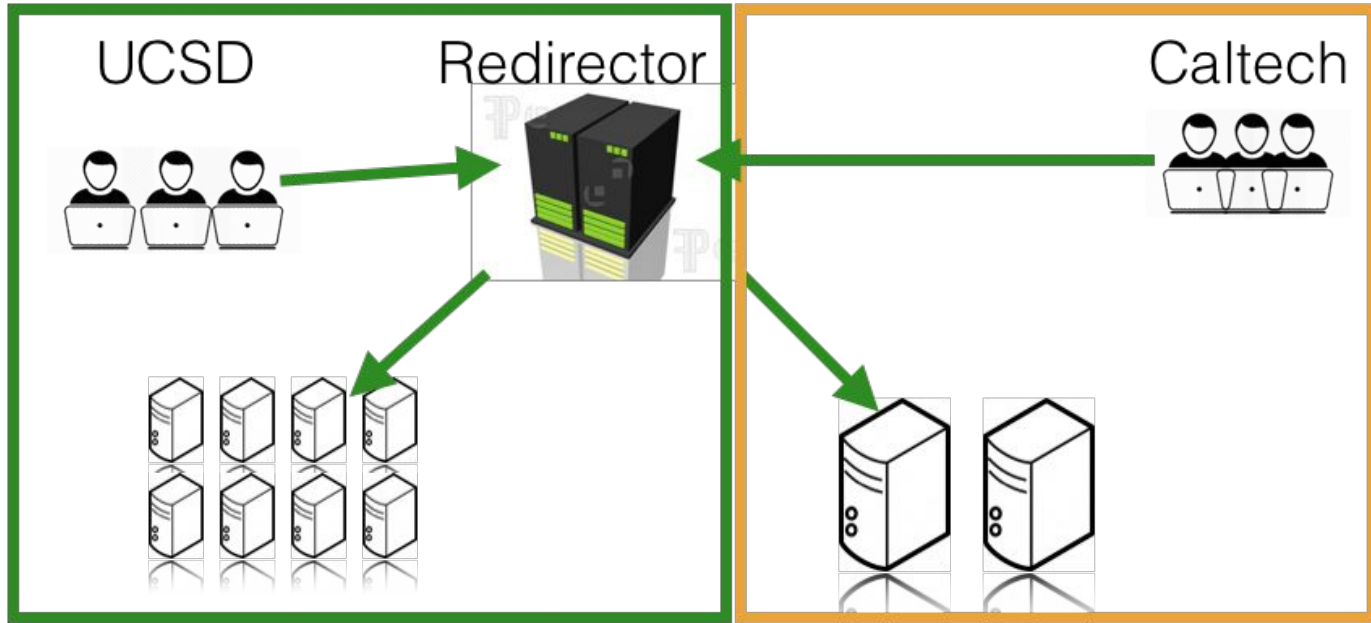
Daily Working Set (October, 2019)



- The daily working set for SoCal for October was at most:  
70 TB = 35 TB Data + 35 TB MC.

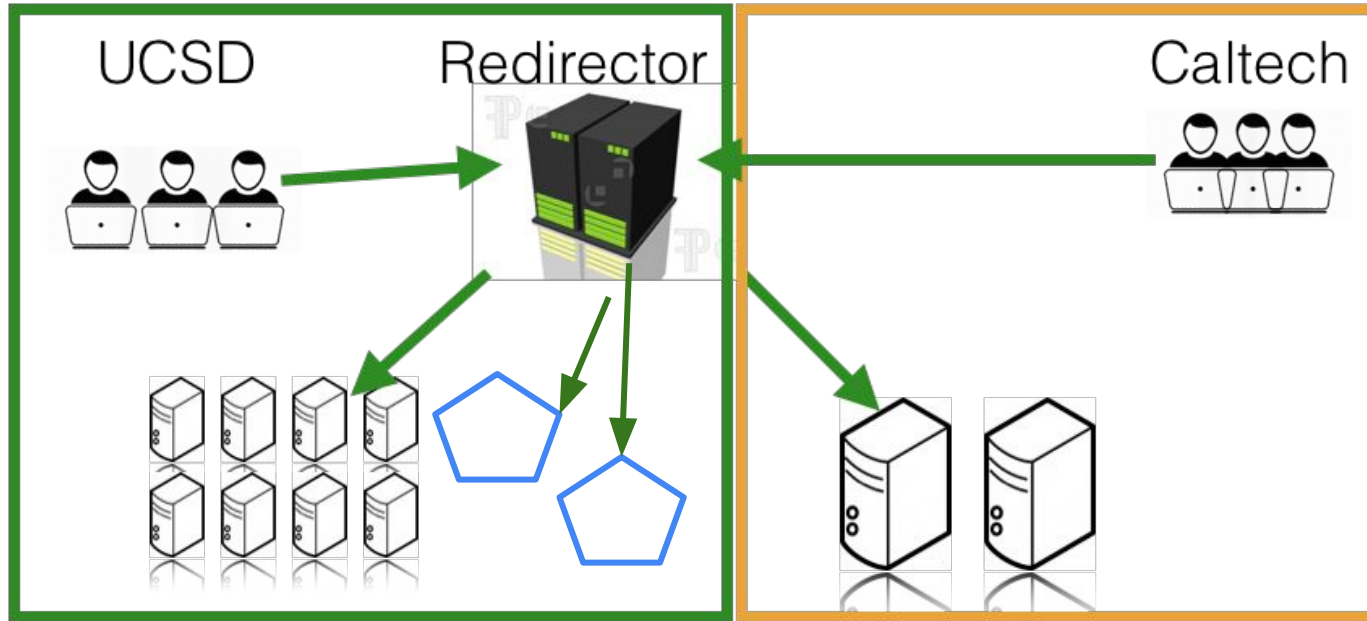
# 4. Install and configure XCache

2018 setup



Each server acts as an independent cache and through the redirector they all work as a logical cache

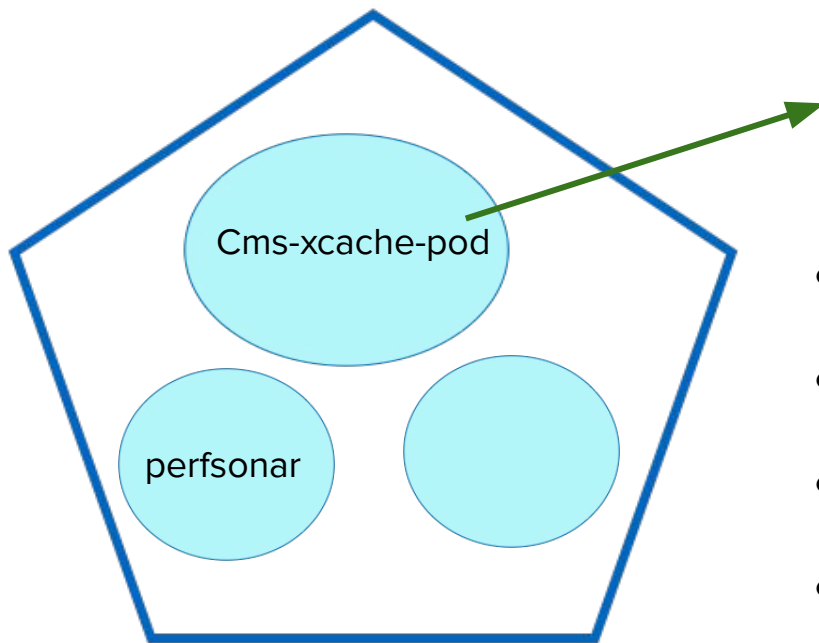
## 2019 and forward Setup



All new XCaches are installed via kubernetes. Two of them currently installed this way.



# How do the pods look inside?



Physical machine with kubernetes installed

Redirector



- CMS XCache [image provided by OSG](#)
  - Based on OSG [CMS XCache RPM](#)
- We added the extra configuration via kubernetes knobs
  - Our github repo with k8s yaml is [here](#)
- We are not using virtualized network capabilities: we bind the pod to the network interface.
- Authentication comes from the RPM and is done via LCMAPS
- We use the PRP kubernetes federation ([See Igor's talk](#))

# Recipe: How to setup an XCache for CMS at a given site(s)?

1. Decide on the namespace to cache. (We used /MINIAOD\*)
2. Calculate the working set (the set of unique files of the namespace that are accessed on a given period).
3. Provision the storage for some fraction of data from step #2.
4. Install and configure XCache mentioned on step #3.
5. Monitor
6. Experiment with some changes
7. Compare
8. Conclude

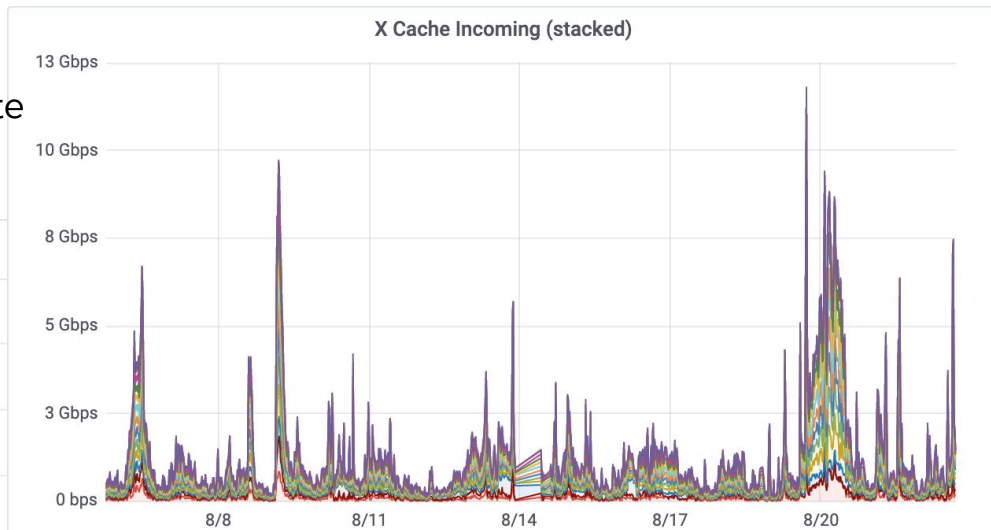
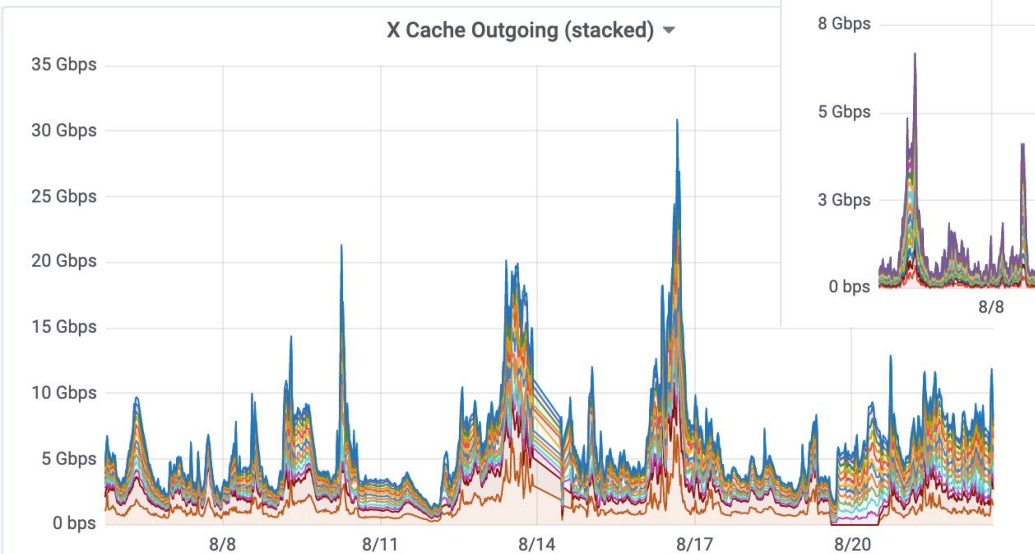
# Monitoring

- Right now monitoring comes from several sources:
  - Telegraf + influxDB + Grafana on the bare metal
  - Monit job information (failure rates, IO times, etc.)
  - Kubernetes own network monitoring (Prometheus)
  - XRootD monitoring dashboards going through monit.

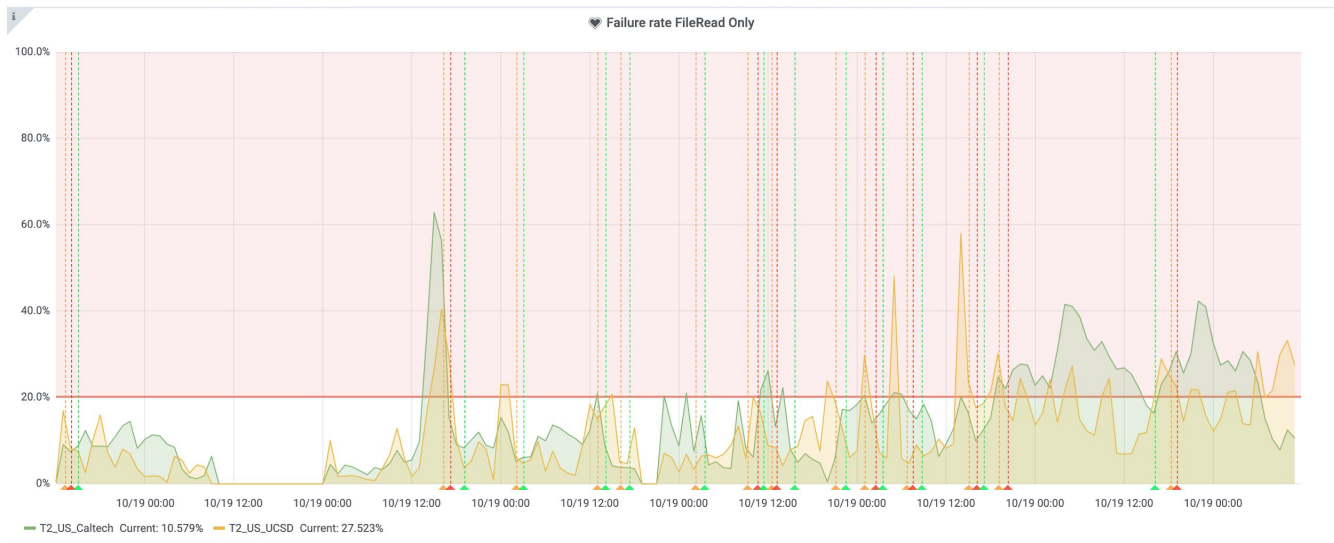


# Telegraf + InfluxDB + Grafana (Bare Metal)

Allows us to approximate instantaneous hit rate



# Monit job information (failure rate)



Failure rate for CMS Global pool jobs for only MINI\* for SoCal sites. But also: CPU Efficiency, AvgReadTime, AvgInputSize, AvgOutputSize.



# Experiment with some changes

At the end of October we made a switch to have MINI\* jobs running at **Caltech** read from **AAA** and **local hadoop** while **UCSD** jobs use the **XCACHE**.

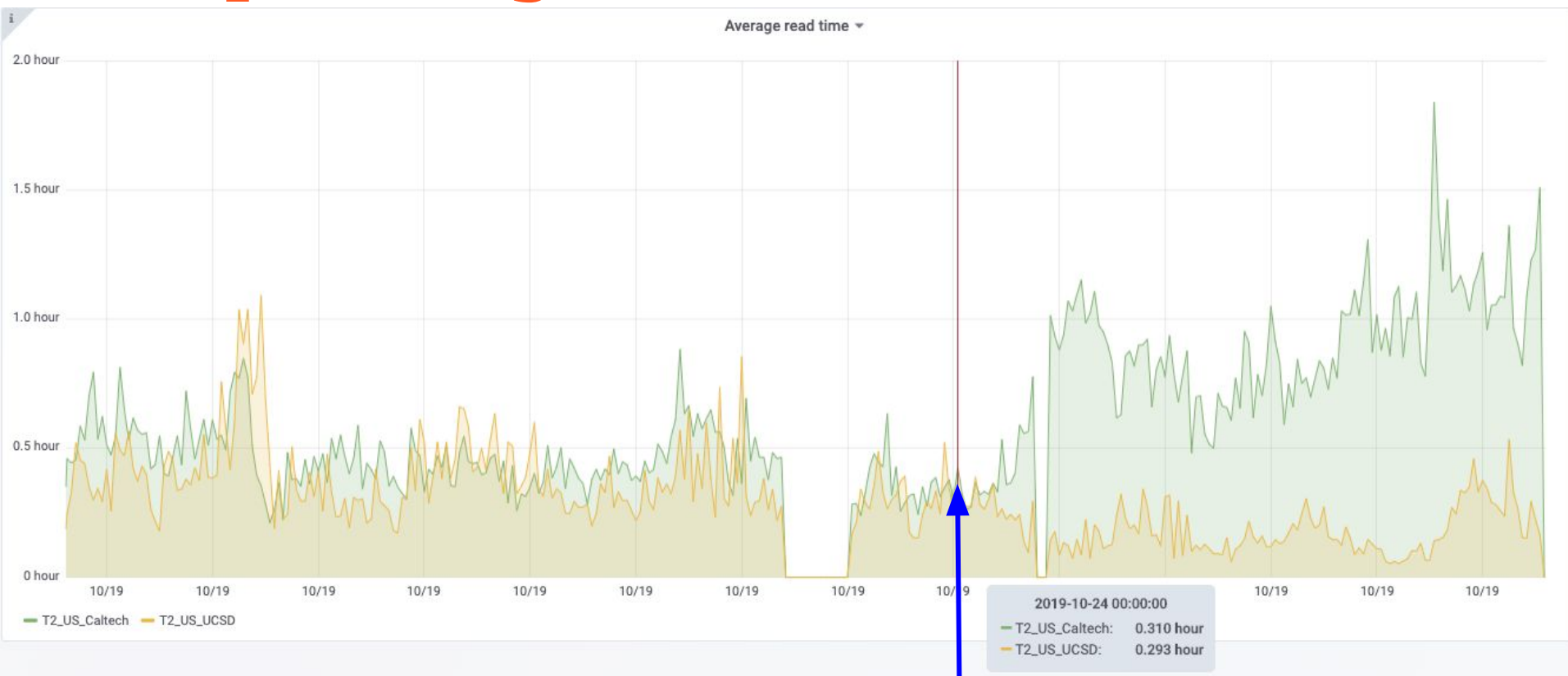
## UCSD

```
# Portions of /store in xcache
<!fn-to-pfn protocol="direct" destination-match=".*"
  path-match="/+store/(data/.*/.*MINIAOD/.*)"
  result="root://xrootd.t2.ucsd.edu:2040//store/$1"/>
<!fn-to-pfn protocol="direct" destination-match=".*"
  path-match="/+store/(mc/.*/.*MINIAODSIM/.*)"
  result="root://xrootd.t2.ucsd.edu:2040//store/$1"/
```

## Caltech

```
<!fn-to-pfn protocol="hadoop" destination-match=".*"
  path-match="(.)" result="$1"/>
<!-- Xrootd fallback rules -->
<!fn-to-pfn protocol="xrootd" destination-match=".*"
  path-match="/+store/(.)"
  result="root://cmsxrootd.fnal.gov//store/$1"/>
```

# Compare: Avg Read Time



Caltech made the change here (Oct 24 - 2019)

# Future Work

- Investigate other deletion paradigms beyond LRU.
- Have caches serve as data origins to other caches.
- XCache to advertise to DDM's their capabilities rather than their state:  
*I can gather /foo rather than I have /foo*
- We will be installing a second set of caches for NanoAOD with a joint project with ESNNet.
- We will be installing NanoAOD (100TB) caches in some US Tier 2 sites.
- Combine all monitoring sources in a sensible way so one can look only at a few plots and alerts and know the status of the cache.

# Conclusions

- The current SoCal cache has been successful at efficiently using the disk space allocated to only store datasets that are actually used.
  - Drastically reduces Average Read Time per job
- UCSD and Caltech successfully operated a service for the region that can grow horizontally based on needs.
- Working closely with other efforts in CMS:
  - INFN regional caching initiative
  - NANOAOOD caches on all US CMS sites
- Influences data-storage & distribution strategies for future computing models and data lakes.

# Acknowledgements

- The authors would like to thank the funding agencies for this work, in particular the following grants:
  - PRP: NSF OAC-1541349
  - OSG: NSF MPS-1148698
  - IRIS-HEP: NSF OAC-1836650
  - US CMS Ops: NSF MPS-1624356
  - TNRP: NSF OAC-1826967
- To the CMS experiment and central submission infrastructure team
- To the PRP kubernetes federation admins
- XRootD Developers