



# Future Facility Plans

Stu Fuess / Scientific Computing Division

2019 ICAC

14 March 2019

# Outline

- [Side note on operations]
- General statement of problem
  - Motivation, complications, solution
- Specifics on current resources, experiment requests – and plans
  - Processing
    - Local, grid, allocations, cloud
    - “HPC”
      - LQCD clusters (new, current, and old)
      - Development systems
  - Storage
    - Disk, tape

## [Side note on Facility operations]

- Local resources are currently specific to **CMS**, **“Public”** (= not CMS, supporting all other experiment activities), or **Lattice QCD**
  - DUNE, Nova, MicroBoone, ICARUS, SBND, Mu2e, Muon g-2, many others... Common funding
- Important to note that **people** operations are (mostly\*) in common
  - Hardware purchasing and provisioning
  - System administration
  - Storage systems
  - Batch systems
  - Supporting services

\* Several services on LQCD clusters traditionally independent, but slowly fixing this

## Motivation for change

- Expect to have limited / insufficient local resources
  - Need to find more elsewhere
- Need to leverage opportunities to utilize new (not traditional HTC) resources
  - Cutting edge technology, accelerators, interconnects
  - Massive size
  - Better economics
- Want to break ties of distinct physical resources (clusters, etc.) that are closely matched to their logical function (support of an experiment or project)
  - Current model of sharing (WLCG, OSG), as pledges or opportunistic, are largely on similar resources

## Complications moving from homogeneous to heterogeneous

- Must understand the importance of **data locality** and networks
- Must support **variety of architectures**
  - Need container build and management infrastructure
- Must understand **local storage limitations** (both on node and on system/cluster)
  - Often optimized for speed/latency, not capacity
- Must deal with In/Out **WAN access limitations**
  - for code (cvmfs), data, workload management, conditions, ...
- Must work with expanded **proposal / allocation / purchase** method
- Need more extensive and complex **monitoring**
- Need more extensive and complex **accounting**
- Need more complex (federated?) **authentication / authorization** infrastructure
- Need to understand impact of **limited support** at remote sites

## Solution: expand the “facility”

- Move to a logical workload description based on characteristics of job, and match to physical resource satisfying those attributes
  - Allows significant expansion of types of jobs and match to heterogeneous resources: HPC sites, commercial clouds
- Supply a “**science gateway**” for workloads, implemented as **HEPCloud**
  - Provisioning based on workload / job characteristics
    - E.g. memory, MPI, architecture, accelerators, allocations, funding, storage...
  - “Best match” made by Decision Engine to resource attributes

- HEPCloud system
  - Have DOE ATO and went “live” this Tuesday, 12-March-2019 !
    - Accessing local clusters, NERSC, Amazon, Google
  - Job submission will look the same, now with additional optional attributes
  - On-boarding of experiments serially to ease transition
    - CMS – interface to global mechanism
    - Nova, Mu2e, DUNE – utilize Fermilab jobsub mechanism
- Initially directing location-agnostic processing (compute cycles)
  - “Low-hanging fruit”
- Matching with storage is more challenging, with continued development
  - Move towards unified data management
  - Co-scheduling as needed / when possible
- Will add more sites in future: LCFs, NSF/XSEDE sites



## Processing: Summary of current resources

- **CMS Tier-1 and LPC**: to meet pledge and provide analysis platform, **~27K cores**, 285 kHS06
- **FermiGrid**: Intensity Frontier and other HTC usage, **~19K cores**, 200 kHS06
- **LQCD clusters**: allocated, high speed interconnect (**IB**), some **GPUs**
  - Existing:
    - pi0 : 5,024 cores --- only ~1/4 allocated to LQCD post 2019
    - pi0G : 512 cores, 128 K40 GPUs --- no allocation to LQCD post 2019
    - Bc : 7,168 cores ---
    - Ds : 6,272 cores | All these are ancient
    - DsG : 320 cores, 80 Tesla M2050 GPUs ---
  - Bid in progress:
    - IC : ~75 nodes (Cascade Lake?) + 5 nodes with dual Voltas --- 92% LQCD allocated
- **Wilson cluster**: development with various accelerators, **small HPC**

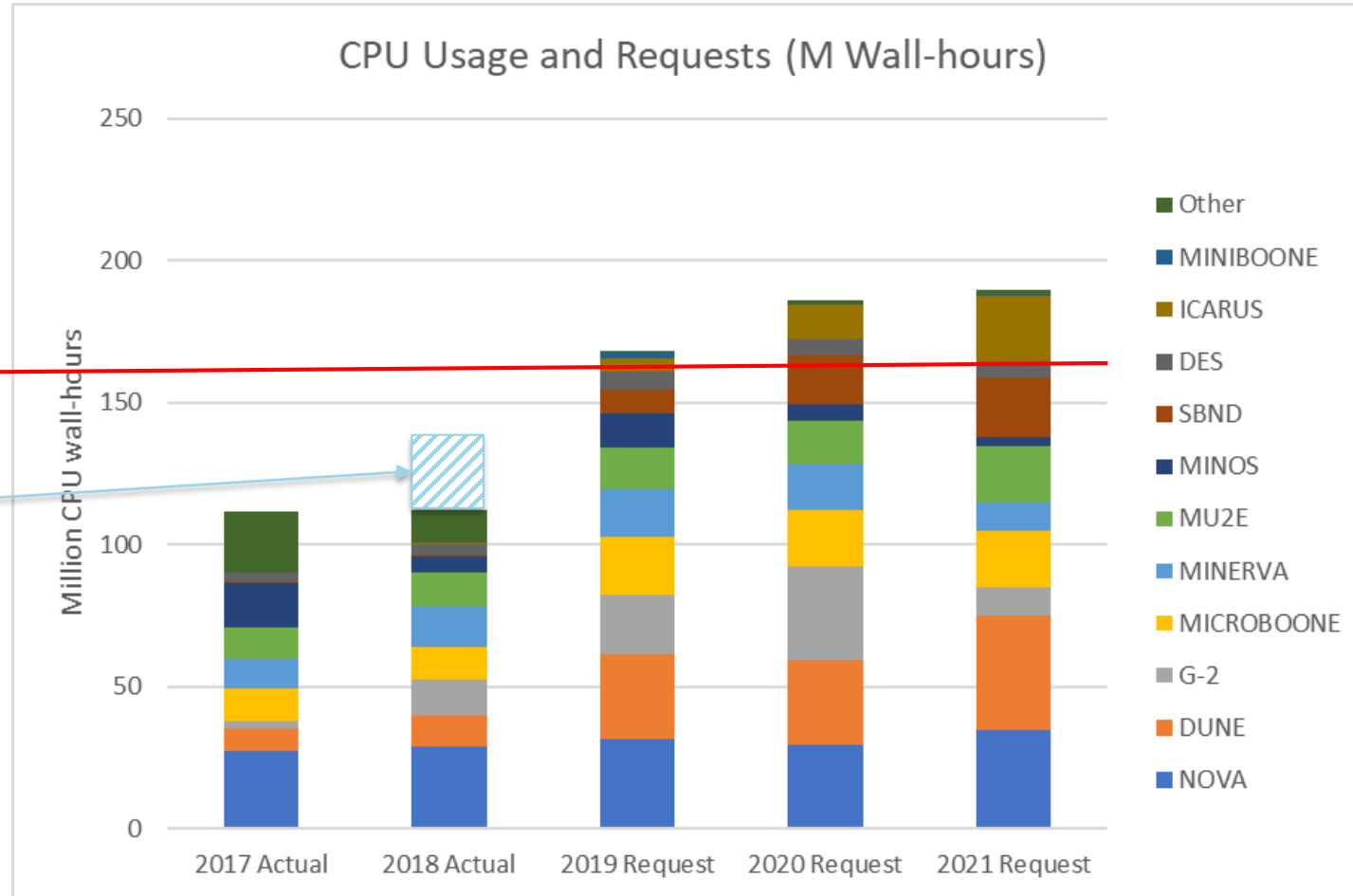


## Processing future: CMS use of HEPCloud

- 2019 Tier-1 pledge: 260 kHS06 (285 kHS06 currently available)  
2020-2021 pledge: 338 kHS06 (need to replace retirements, add some)
- 2019 CMS HPC allocations (requested annually)
  - DOE
    - NERSC (82M hours Cori)
    - ALCF (0.5M hours Theta)
  - NSF/XSEDE
    - SDCS (Comet), PSC (Bridges), TACC (Stampede)
- Eventually expand T1\_US\_FNAL to include all HPC allocations
  - Map workflow characteristics to resource capabilities
  - Meet some of the pledge with external resources
  - Discussion started if and how some part of the pledge can be met with external resources

# Processing future: Public HTC Requests

- Summary of processing history and current requests from all experiments participating in SCPMT:



Current capacity  
160 M hours/year

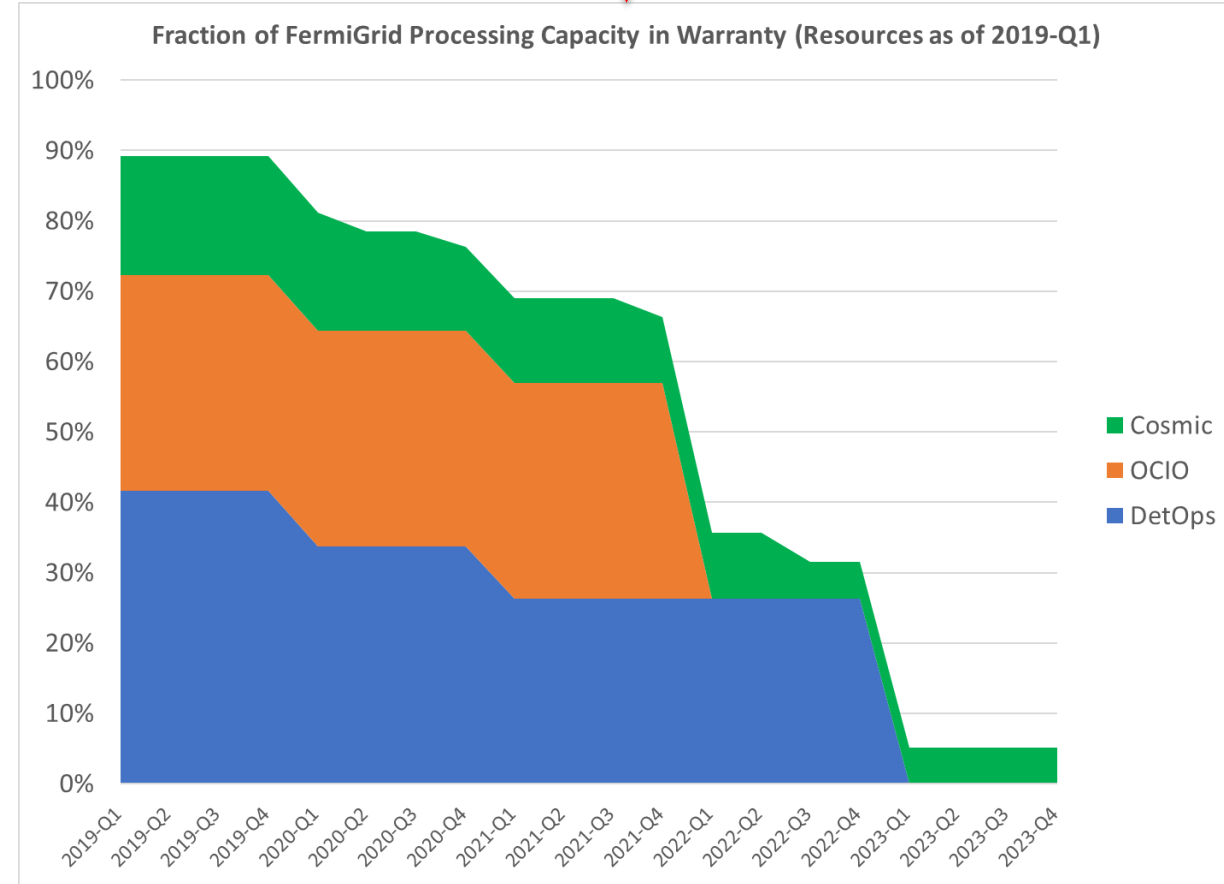
Opportunistic use  
from OSG ~ 24M

Add ~ 5M hours/year  
to requests for other  
local usage

Bottom line:  
HTC need is to  
sustain at approx.  
current level

# Processing future: Public HTC resources

- FermiGrid: shared (all except CMS) worker nodes
  - Approximately **19,000 cores** of various vintage
    - Availability of ~ **160M core-hours per year (200 kHS06 units)**
    - Last purchase using Computing and Detector Operations funds was in FY17
    - No funds for additions in FY19
      - ~ \$2M purchase price
      - To replenish 20%/year need ~ \$400K
  - At least 2 GB per core
    - some (for DES) have ~ 5-6 GB per core (256 GB/node)



## Processing future: HPC/accelerator

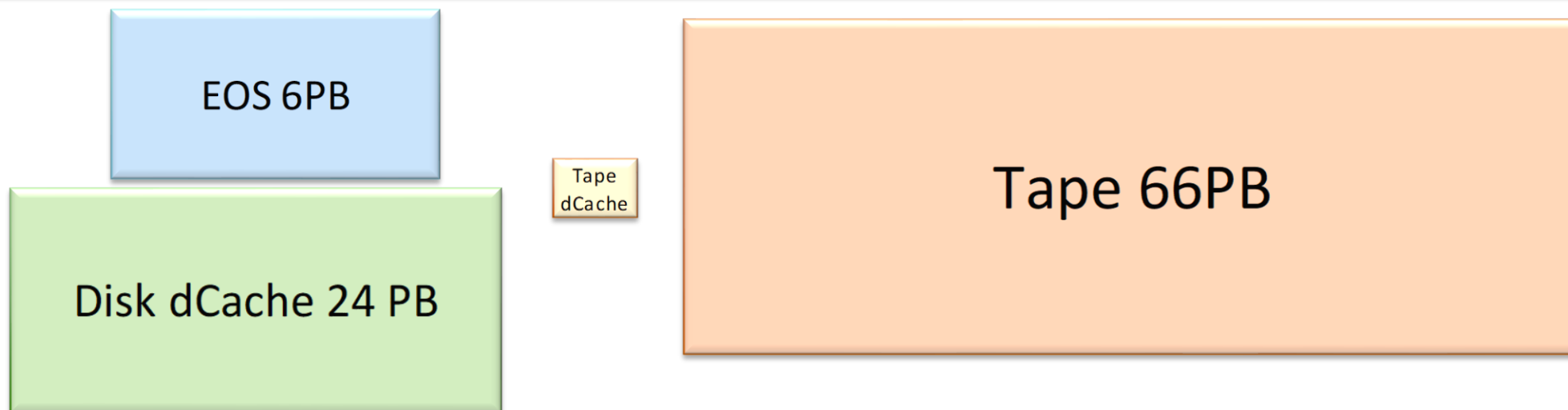
- Existing resources
    - pi0G cluster (512 cores, 128 K40 GPUs) will be available for general use in 2020
      - “HPC like” in that nodes have no external connectivity
      - Limited cluster storage (~1PB Lustre)
    - Wilson cluster
      - Currently available, small, but very ancient HPC cluster
      - Also home of various development platforms:
        - 5 GPU enabled hosts, 1 KNL host, 1 “Summit” Power9 node (these will move to IC, below)
  - New/pending resources
    - “Institutional Cluster” (\*) RFP in progress
      - ~75 nodes + 5 nodes with Voltas, IB, ~1PB Lustre
      - Operated as a service, with LQCD “purchasing” hours (promised ~92% of available)
- \* The “processing as a service” model will be applied to all local resources  
With access via HEPCloud

## Processing future: Summary

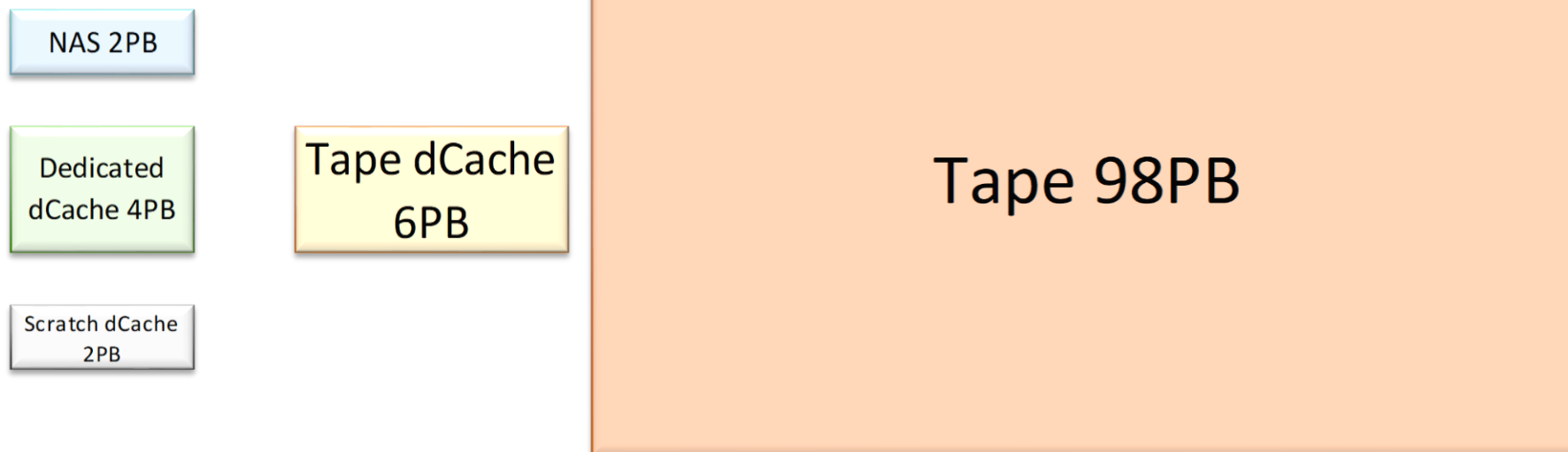
- HEPCloud will be the gateway to both local and external resources
- In aggregate, local resources will follow the “Institutional Cluster” model
  - “Processing as a service”
  - With allocations and “cost” accounting
- Local HPC resources provided at a level enabling:
  - Code development
  - Container development
  - Testing at small-to-mid scale

## Storage: Current usage

- CMS

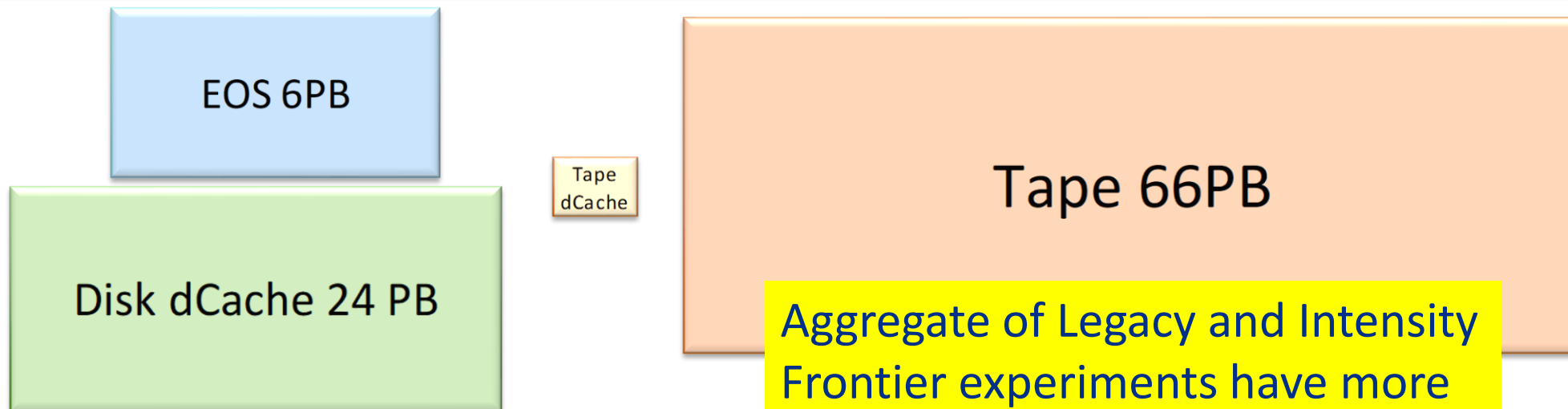


- Public



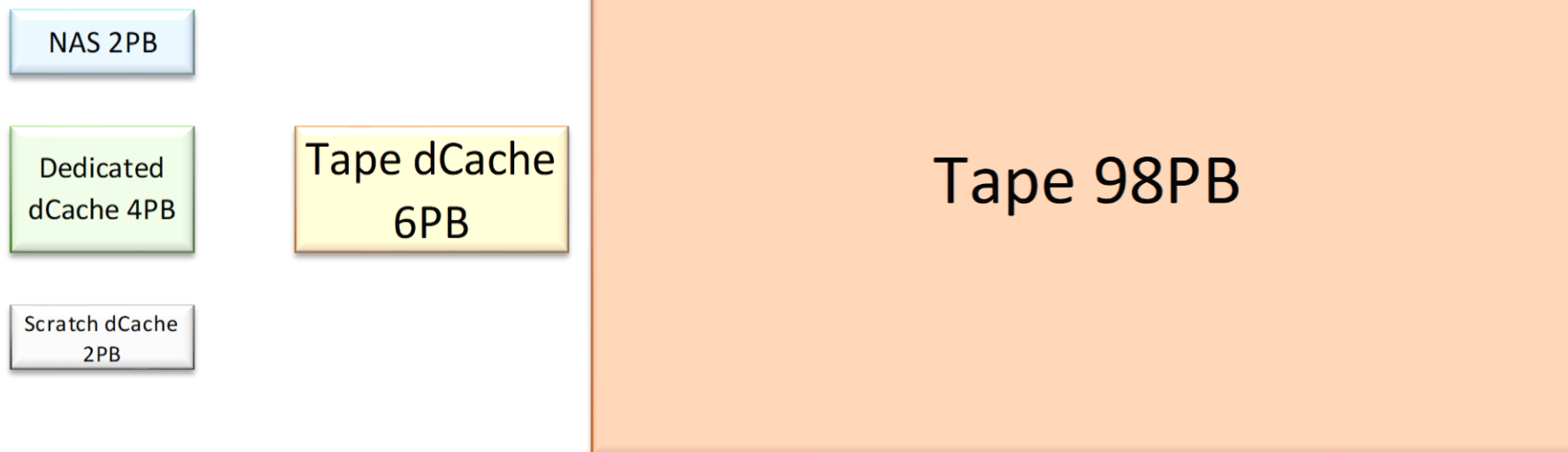
## Storage: Current usage

- CMS



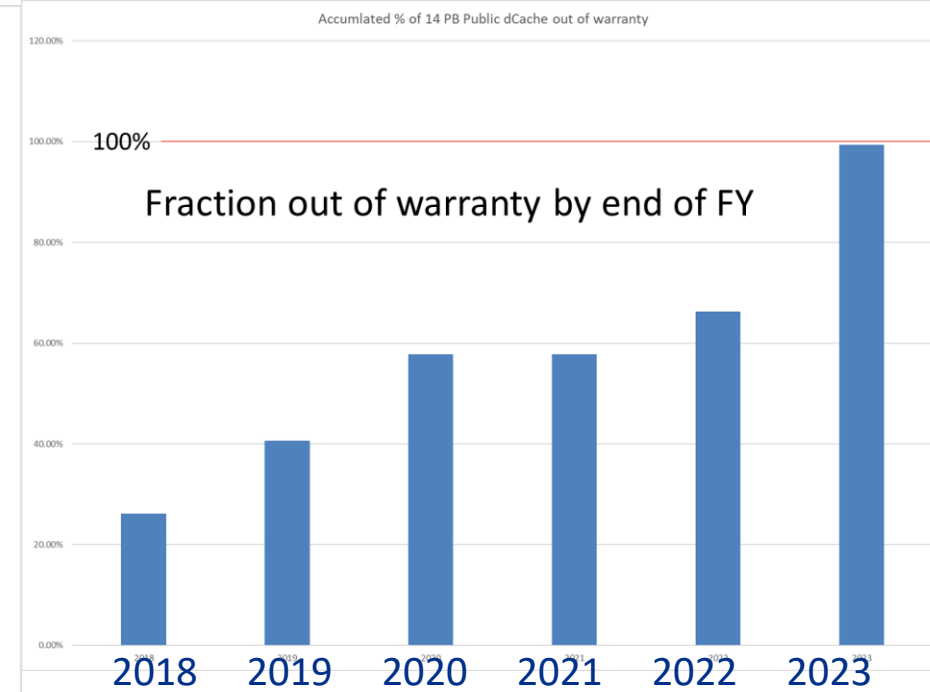
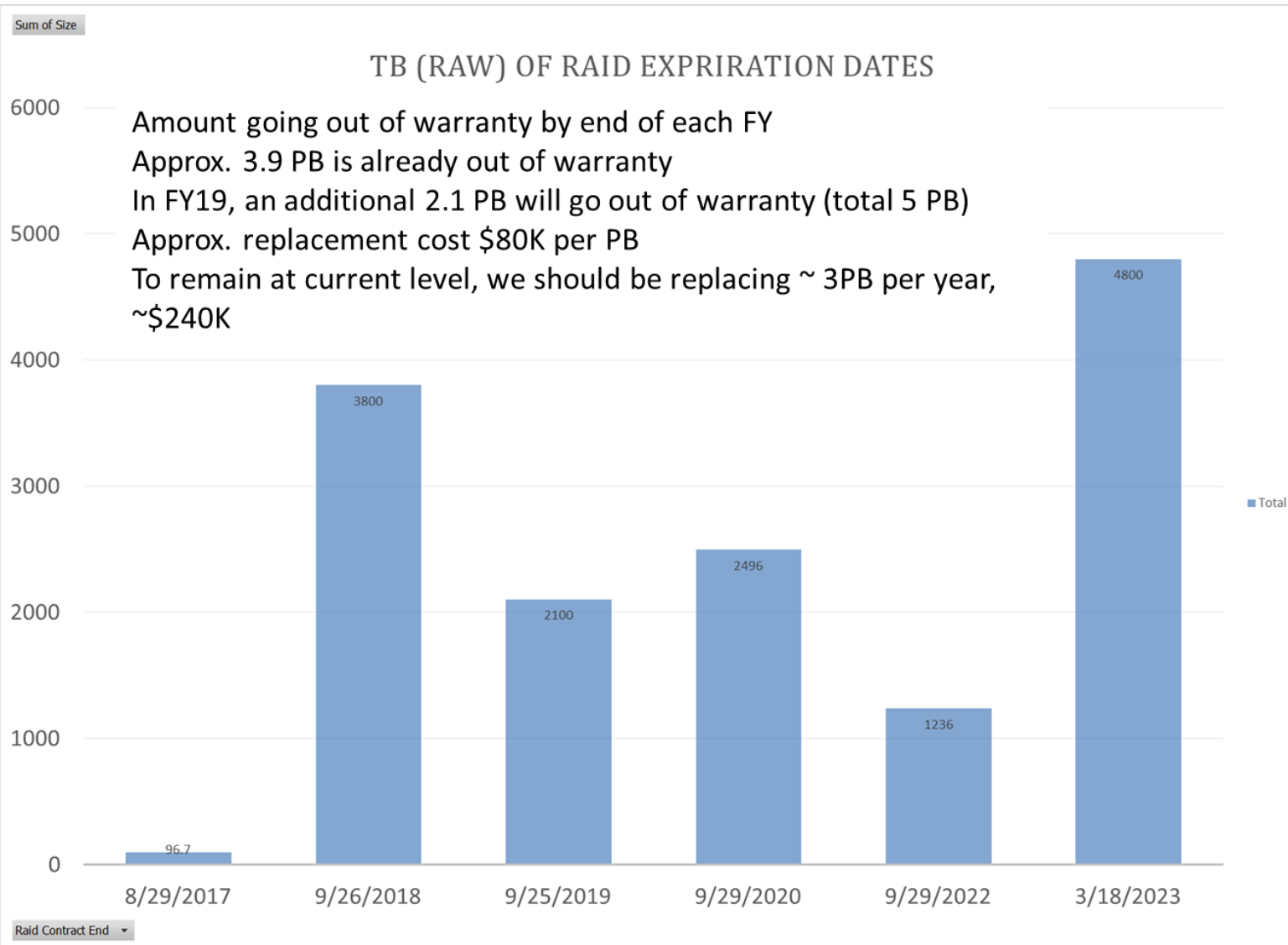
- Public

Paucity of disk means far greater use of tape by average user





# Public dCache disk: Warranty expiration dates



**Bottom line:**  
 Funding constraints unlikely to allow little expansion of Public disk

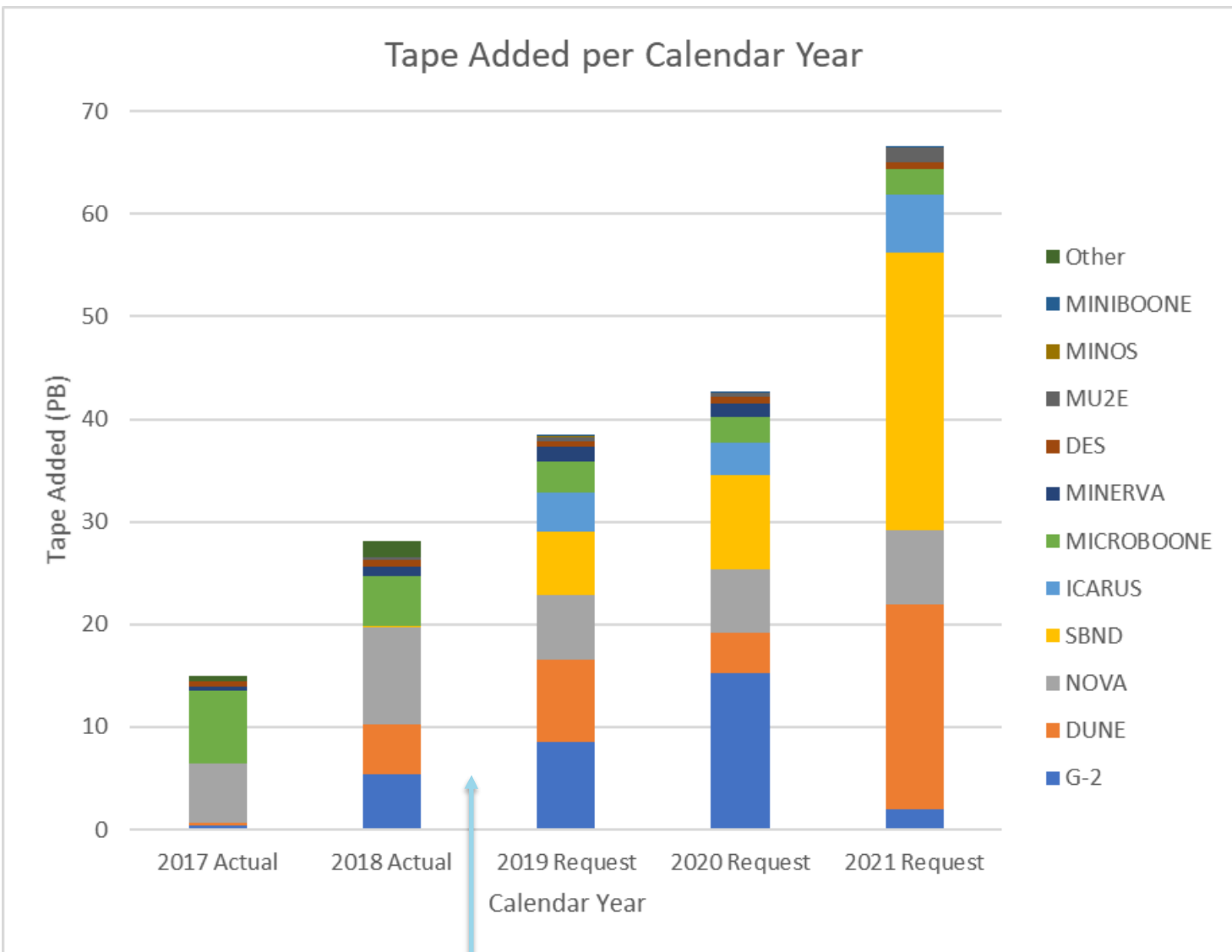
## Tape: Hardware status

- We see no near-term alternative hardware technology for archival storage
- Technology change (from Oracle to...):
  - At start of 2018 we had 7 10K-slot SL8500 libraries with ~80 enterprise drives
  - Have retired 2 libraries, purchased 2 new 8.5K slot IBM libraries (will do 3<sup>rd</sup> this year)
  - Moving to (~100) LTO8 drives with M8/LTO8 media
    - With LTO8, each new IBM library is ~ 100PB
- Need to both ingest new data and migrate legacy data
  - ~140 PB (+20PB CDF, D0) of existing data to potentially migrate

## Tape: Software status, plans

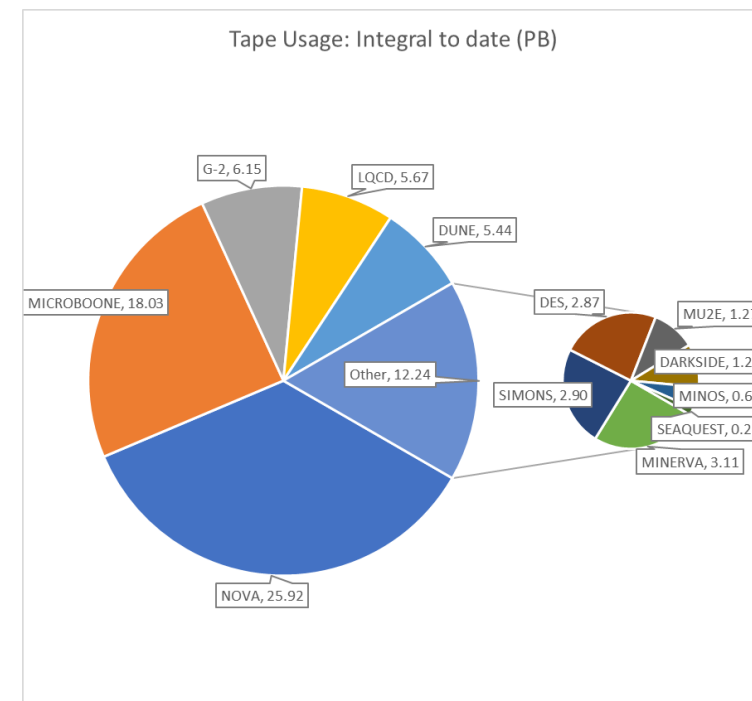
- Fermilab uses enstore for all tape storage
  - Closely connected as HSM to dCache
  - enstore also used by another CMS Tier-1 (PIC) and several Tier-2s
  - But limited personnel with enstore expertise
- CERN has used Castor, moving to CTA
- Fermilab will evaluate CTA as future option
  - Tape format is a complication
    - CERN uses “CERN format” for both Castor and CTA, so can physically “move” tapes to CTA
    - enstore uses CPIO format, which would require copying files (so best done at a migration)
  - Need to evaluate effort in all surrounding utilities

# Tape: Volume of "Public" (=not CMS) new tape requests



For reference, the net tape usage to date:

Experiment	Net to date (PB)
NOVA	25.92
MICROBOONE	18.03
G-2	6.15
LQCD	5.67
DUNE	5.44
MINERVA	3.11
SIMONS	2.90
DES	2.87
MU2E	1.27
DARKSIDE	1.25
MINOS	0.63
SEAQUEST	0.21
Other	0.81
<b>TOTAL Public</b>	<b>74.25</b>

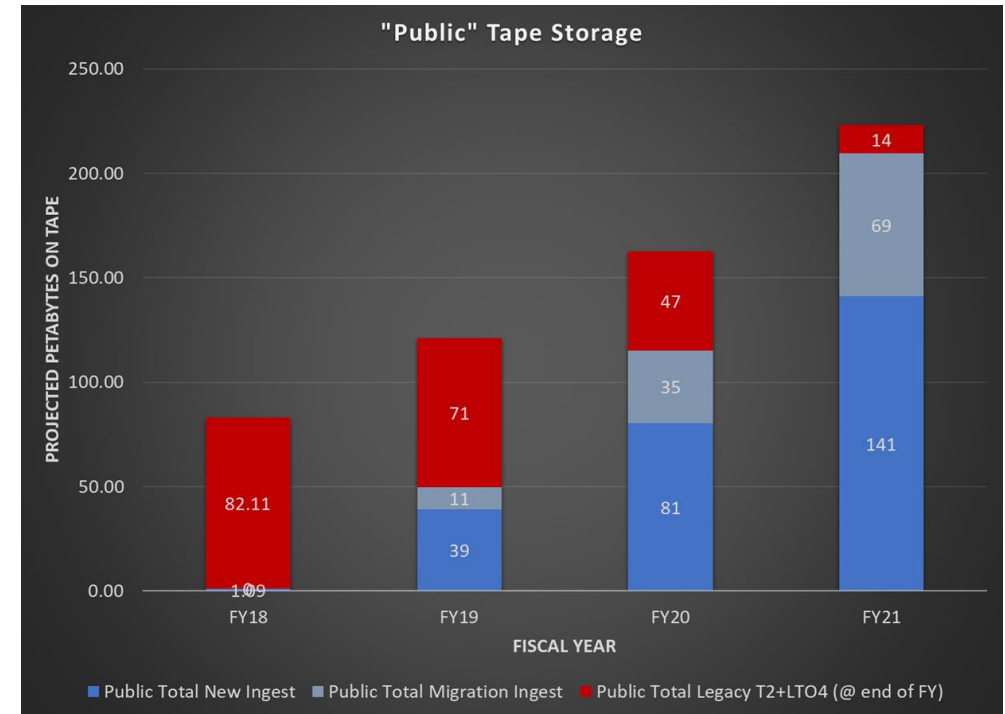
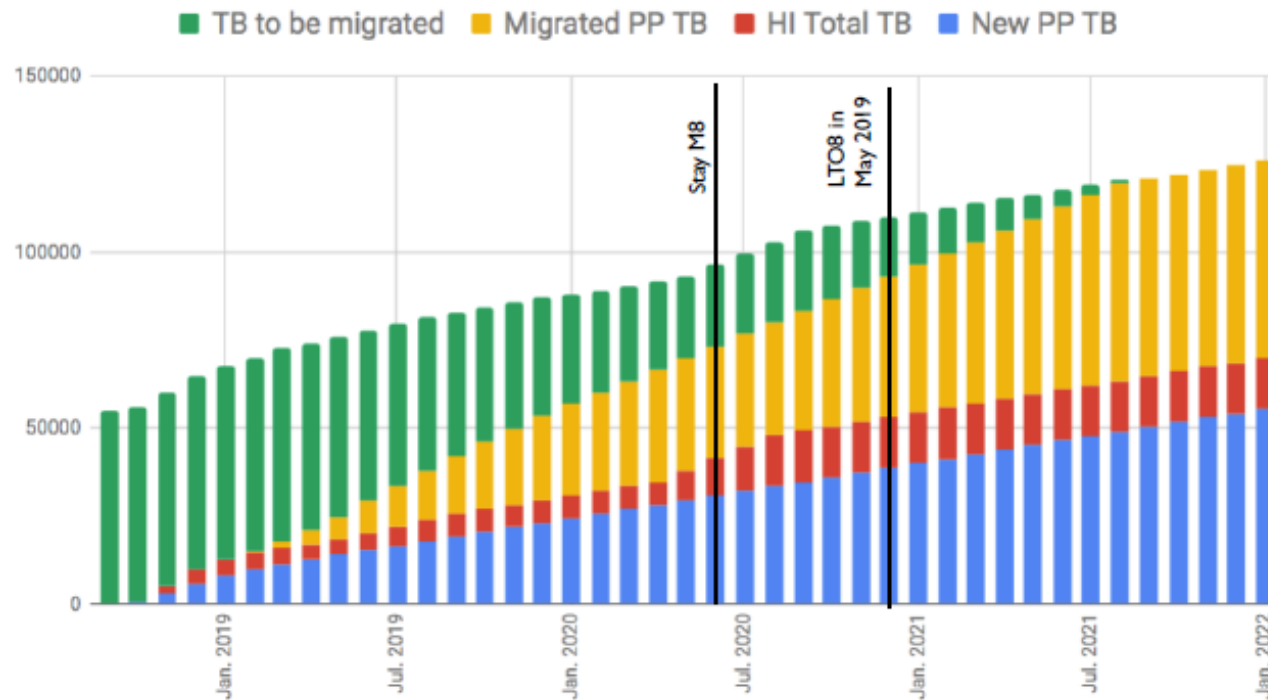


# Tape: Integral

CMS (125PB by 2022)

Public (225PB by 2022)

Facility Tape Written



## Storage future: Summary

- There is a discrepancy between CMS and Public storage architectures and disk/tape balance
  - Would like greater coherence of methodologies
- Storage architecture decisions will be greatly influenced by plans emerging from HSF etc.
- Concern that funding will constrain options for Public systems

## Conclusions

- HEPCloud is seen as the path for uniform access to heterogeneous processing
  - Long path to incorporating more resources, attributes, storage...
- Local resources will appear as a “processing service” to which allocations and cost accounting will apply (the “Institutional Cluster” model)
- The path of storage architecture evolution is not yet clear



# Backup slides

## Disk: numbers

Use	Type	Capacity
CMS	dCache disk only	24 PB
CMS	EOS	6 PB
CMS	dCache tape	1 PB
Public	dCache tape	6 PB
Public	dCache scratch	2 PB
Public	dCache dedicated	4 PB
Public	NAS	2 PB

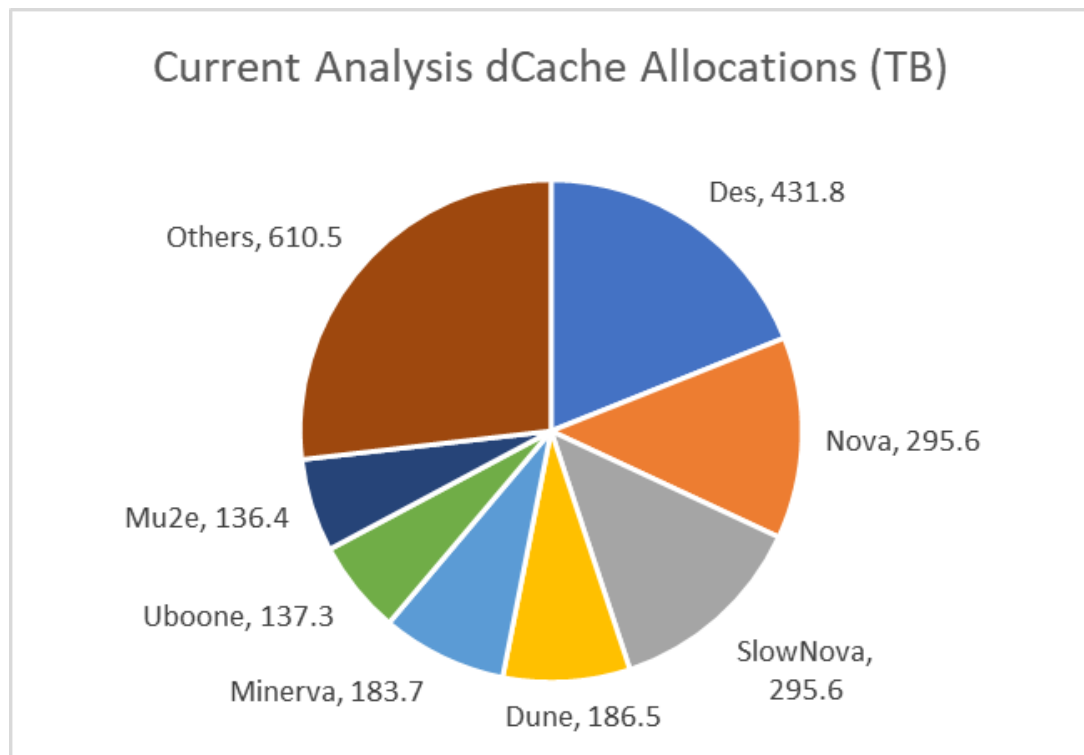
## dCache disk: Resources

- dCache is split into a number of pool groups, some for general use and others dedicated to specific experiment or project use

Pool Type	Number of Pools	Available Space (TB)
Read/Write Cache	2	5,695
Scratch Cache	2	2,122
Analysis / Persistent	32	2,277
Expt. Dedicated	13	2,145
Utility	6	438
<b>TOTAL</b>	<b>55</b>	<b>12,677</b>

## dCache disk: Analysis / Persistent

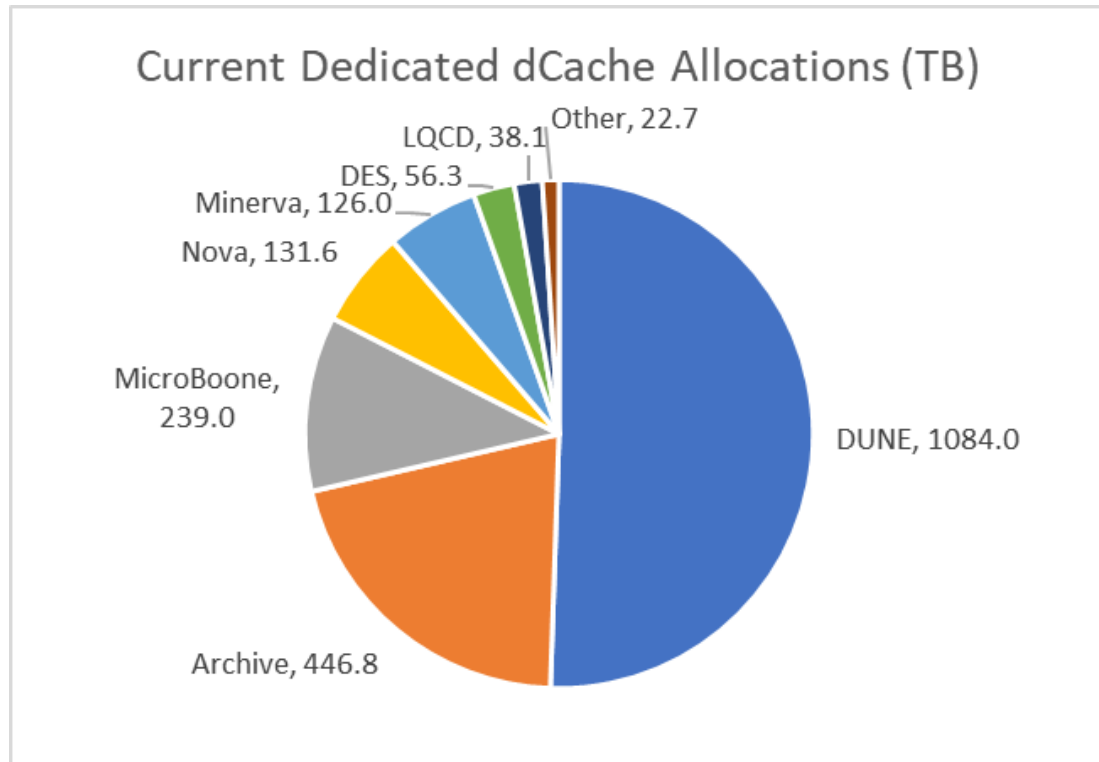
- This is disk space that is permanently resident but with no backup
  - Allocated via SCPMT / SPPM process
  - Management under experiment control
  - 2.3 PB split across 32 experiment/project users



Experiment	2019 Request	2020 Request	2021 Request
DES	400	500	500
DUNE	400	400	800
ICARUS	100	150	200
MicroBoone	300	300	300
Mu2e	150	200	300
g-2	150	300	300
Nova	450	450	450
SBND	100	125	150
Minerva	250	250	250
Others	450	450	450
<b>TOTAL</b>	<b>2,750</b>	<b>3,125</b>	<b>3,700</b>

## dCache disk: Dedicated

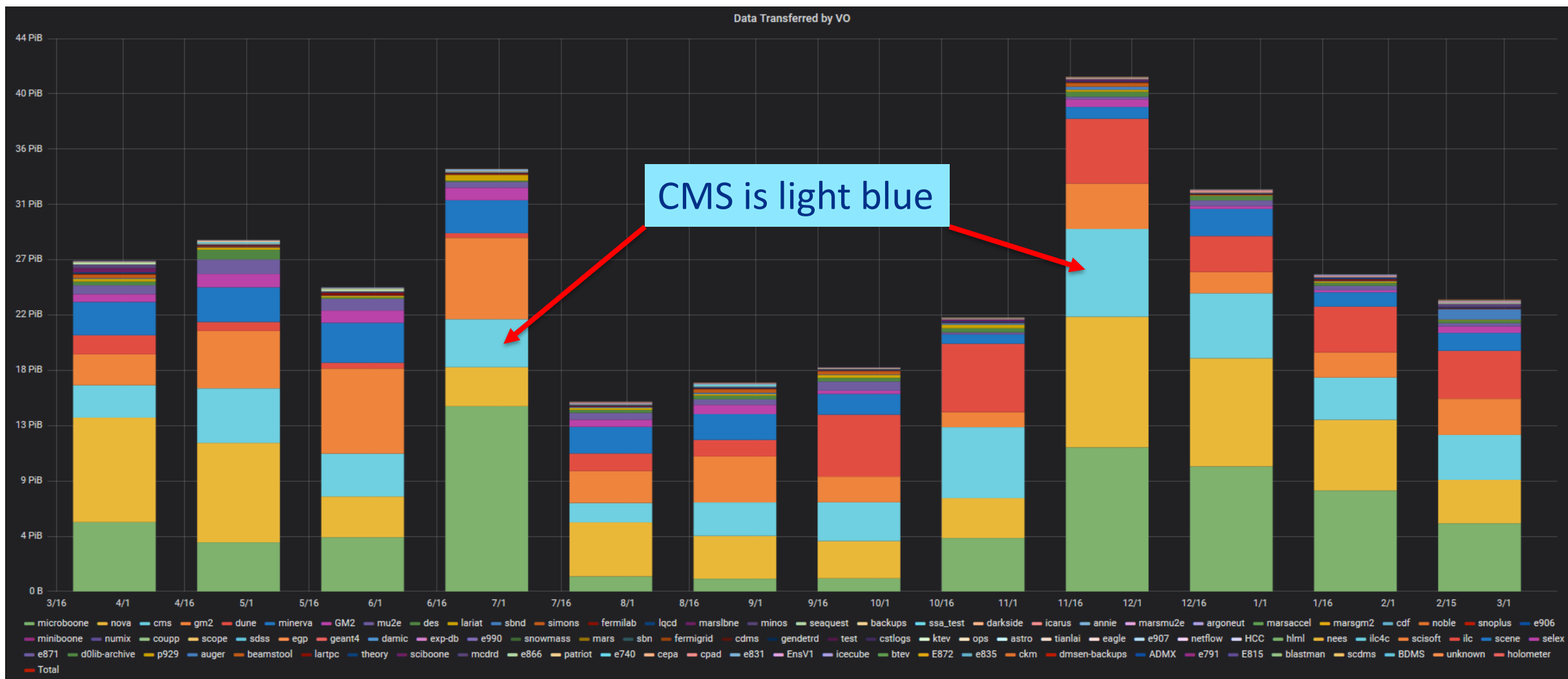
- This is “tape backed” disk space that is dedicated to a specific experiment
  - Allocated via SCPMT / SPPM process
  - Typically for raw data ingest or pre-staging
  - 2.1 PB split across 13 functions



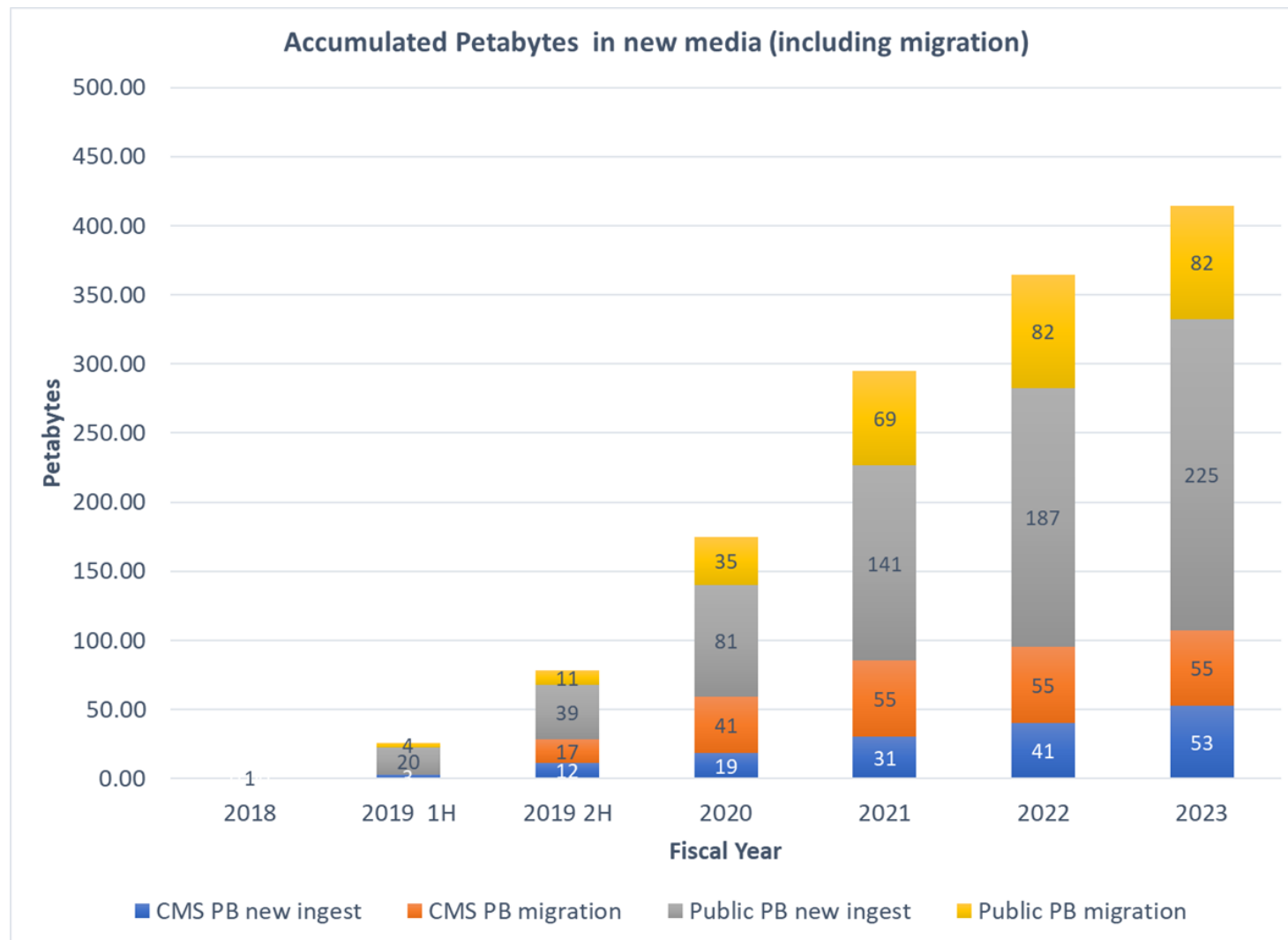
Experiment	2019 Request	2020 Request	2021 Request
DUNE	1,100	1,100	1,500
MicroBoone	?	?	?
Mu2e	0	0	60
Nova	132	132	132
SBND	2	2	2
Minerva	126	126	125
Others	132	132	132
<b>TOTAL</b>	<b>1,234</b>	<b>1,234</b>	<b>1,694</b>

Requests not substantially different than current allocations

# Disk: dCache Transfers by VO (per month)



# Tape: Integral, CMS & Public on new media





# Tape: Transfers by VO (writes, reads per month)

