Image credit: Marguerite Tonjes
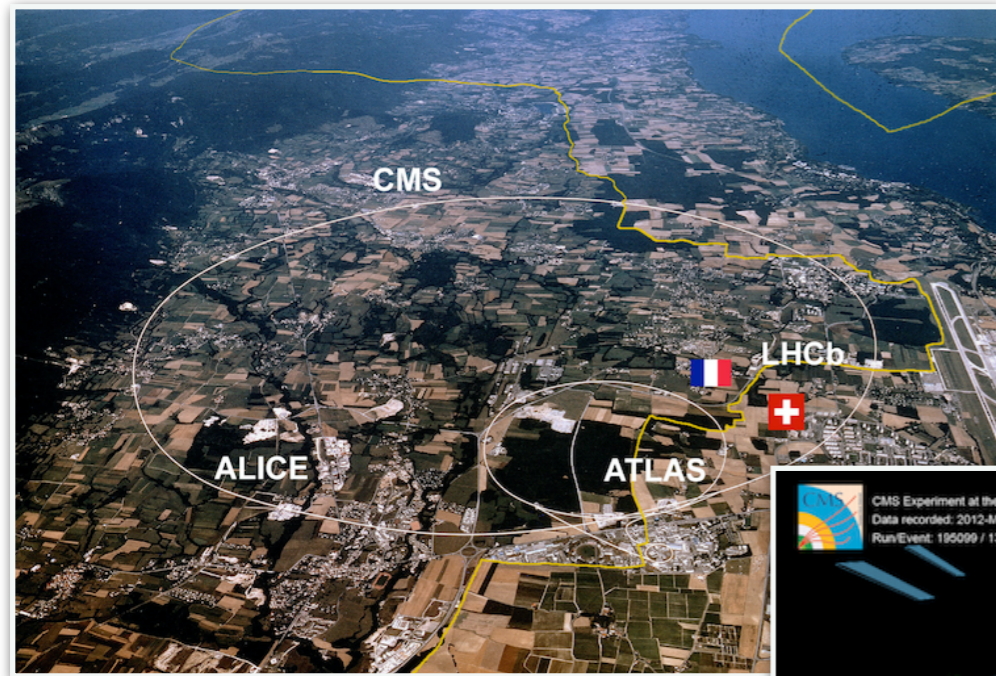
# Storage for HEP

Nick Smith

Computational HEP Traineeship Summer School

23 May 2024

# HEP Experiment: three easy steps
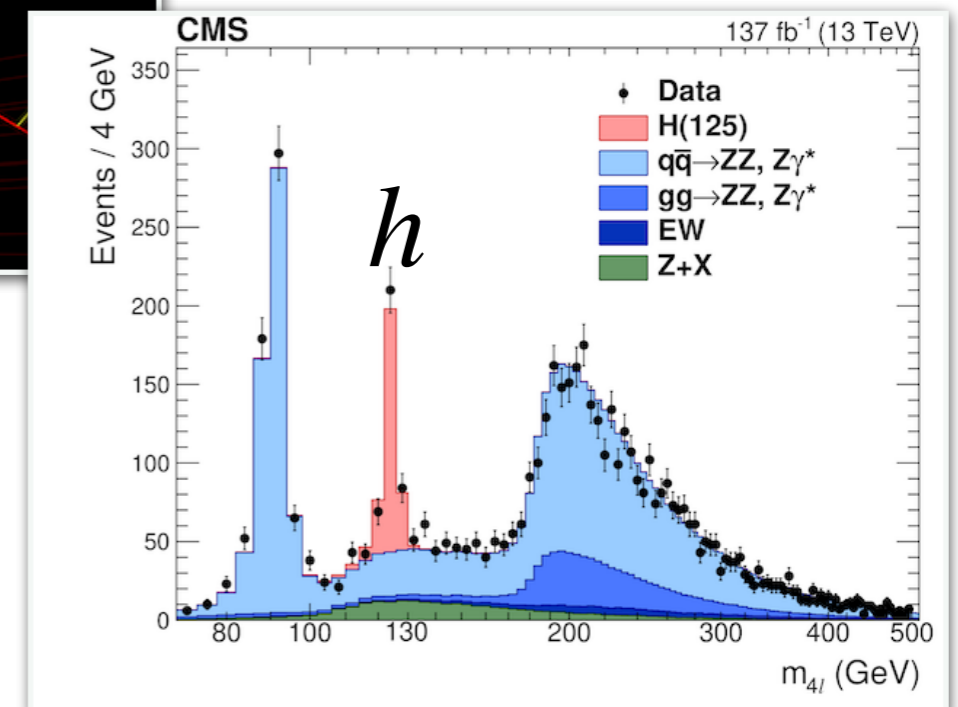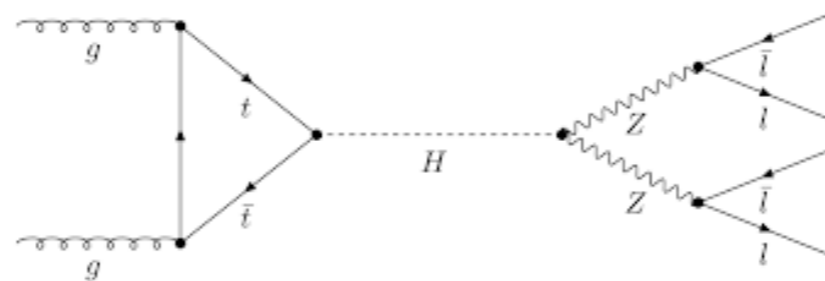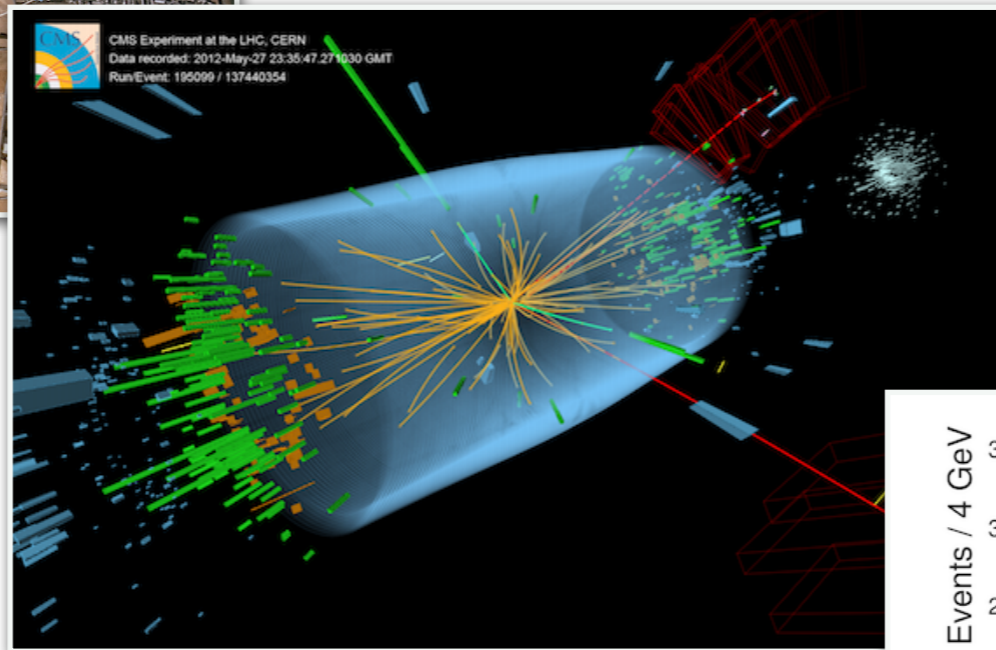
**1. Collide particles**

**2. Take pictures**

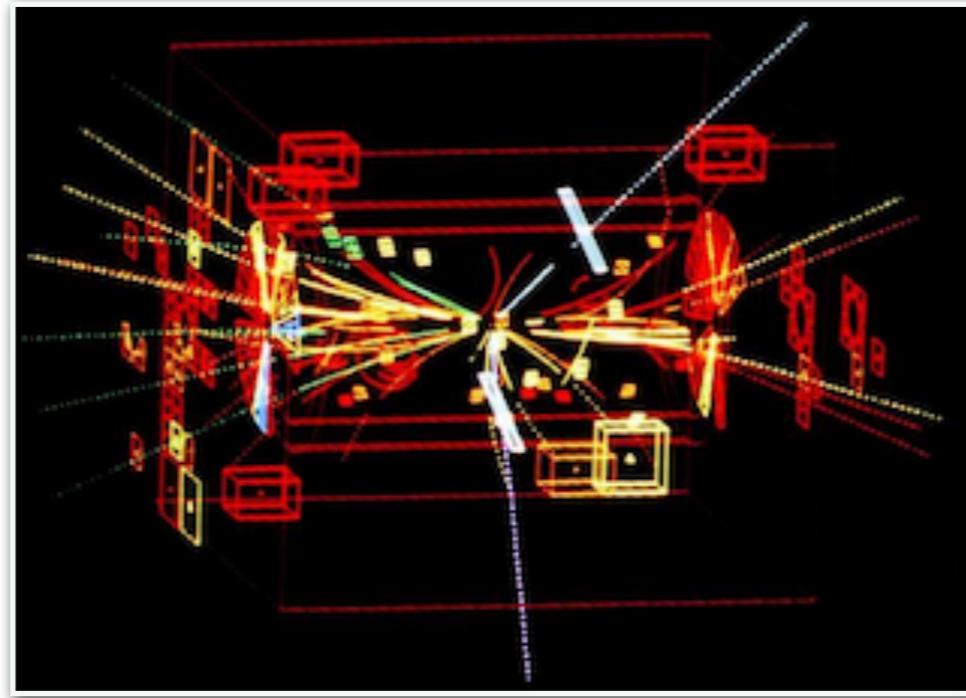**3. Infer parameters**

~1MB / event
~100B events

$h$

# Step 2: take pictures



1 photo / event
~6M events cds:1733654

🧬 **Fermilab**

# Step 2: take pictures

~100kB / event
~10M events cds:182190

1 photo / event
~6M events cds:1733654

🎟 Fermilab

# Step 2: take pictures
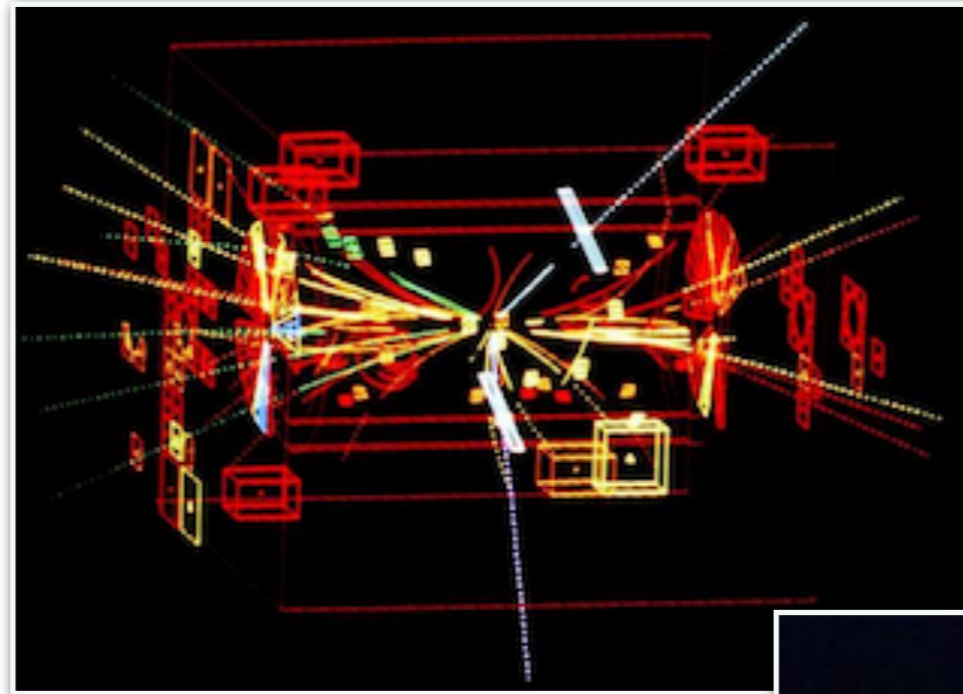
~100kB / event
~1B events 10.1016/j.nima.2017.01.043

~100kB / event
~10M events cds:182190

1 photo / event
~6M events cds:1733654

🔷 Fermilab

# Step 2: take pictures

~1MB / event
~100B events



0kB / event
events [10.1016/j.nima.2017.01.043](https://doi.org/10.1016/j.nima.2017.01.043)

1 photo / event
~6M events [cds:1733654](https://cds.cern.ch/record/1733654)

🔬 Fermilab

# Our inference pipeline

# Our inference pipeline



RAW

# Our inference pipeline



RAW

🔷 Fermilab

# Our inference pipeline

# Our inference pipeline

RAW

Centrally managed

🟦 Fermilab

# Our inference pipeline

Centrally planned, executed

🟦 **Fermilab**

# Our inference pipeline



Centrally planned, executed

RAW

Centrally managed

User-managed

Fermilab

# Our inference pipeline

# Our inference pipeline

Centrally planned, executed



$$P(x|\theta)$$

RAW

Centrally managed

User-managed

🎗 **Fermilab**

# Our inference pipeline



Centrally planned, executed

Analyst / Scientist

$P(x|\theta)$

RAW

Centrally managed

User-managed

🐝 Fermilab

# Our inference pipeline



Centrally planned, executed

**30k CPU-y**  **30k CPU-y**  **30k CPU-y**  **3k CPU-y**

RAW
**100 PB**

Centrally managed
**100 TB-10 PB**

Order of magnitude for **CMS Run-II** (2016-18)
**Processing time** (~100B events)
**Data volume on disk**

🎇 **Fermilab**

# Our inference pipeline



Centrally planned, executed

30k CPU-y        30k CPU-y        30k CPU-y

RAW
**100 PB**

Centrally managed
**100 TB-10 PB**

Sustained high throughput for
Run-II data reprocessing
(If we did this on AWS: ~$40M)

Order of magnitude for **CMS Run-II** (2016-18)
**Processing time** (~100B events)
**Data volume on disk**

🔷 **Fermilab**

# Our inference pipeline

Centrally planned, executed



WLCG: cloud of the 2000s

**Running cores**

Sustained high throughput for
Run-II data reprocessing
(If we did this on AWS: ~$40M)

...rally managed
...0 TB-10 PB

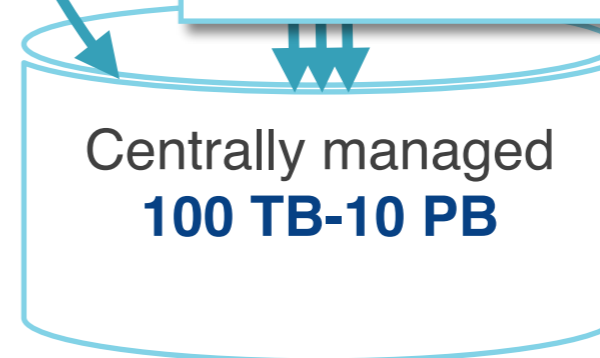Order of magnitude for **CMS Run-II** (2016-18)
**Processing time** (~100B events)
**Data volume on disk**

**Fermilab**

# Our inference pipeline

Centrally planned, executed

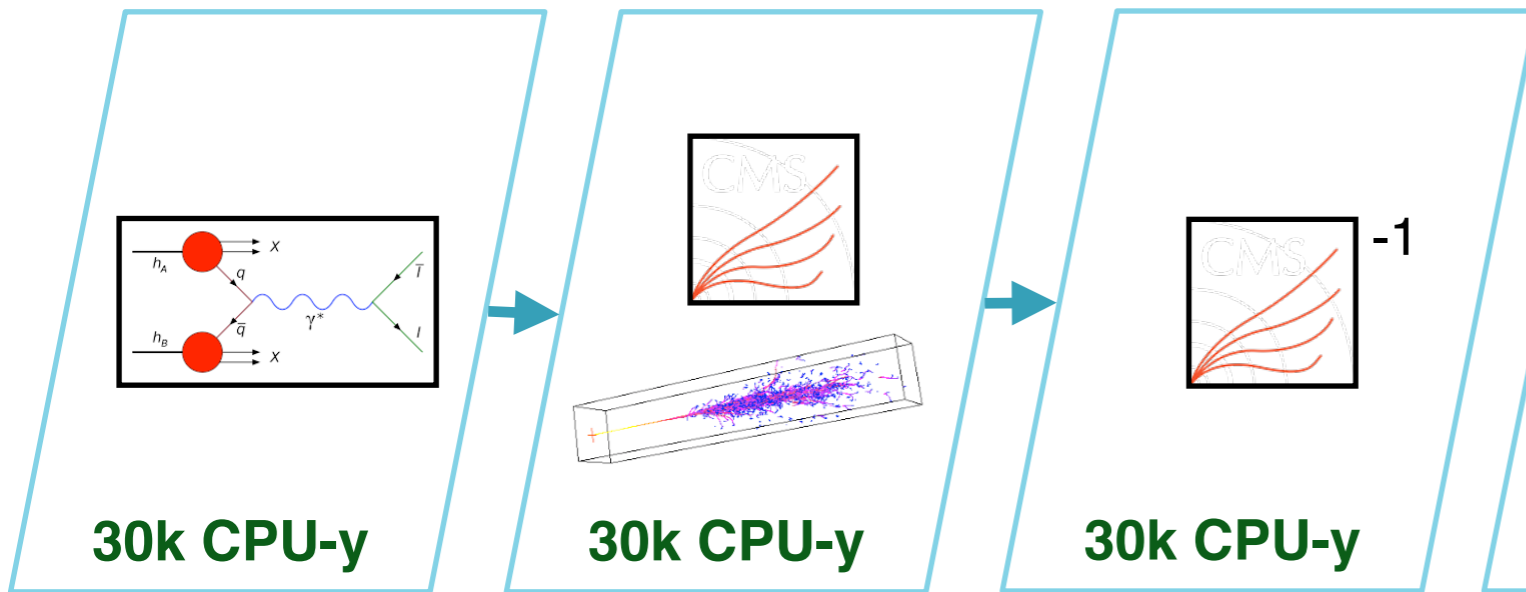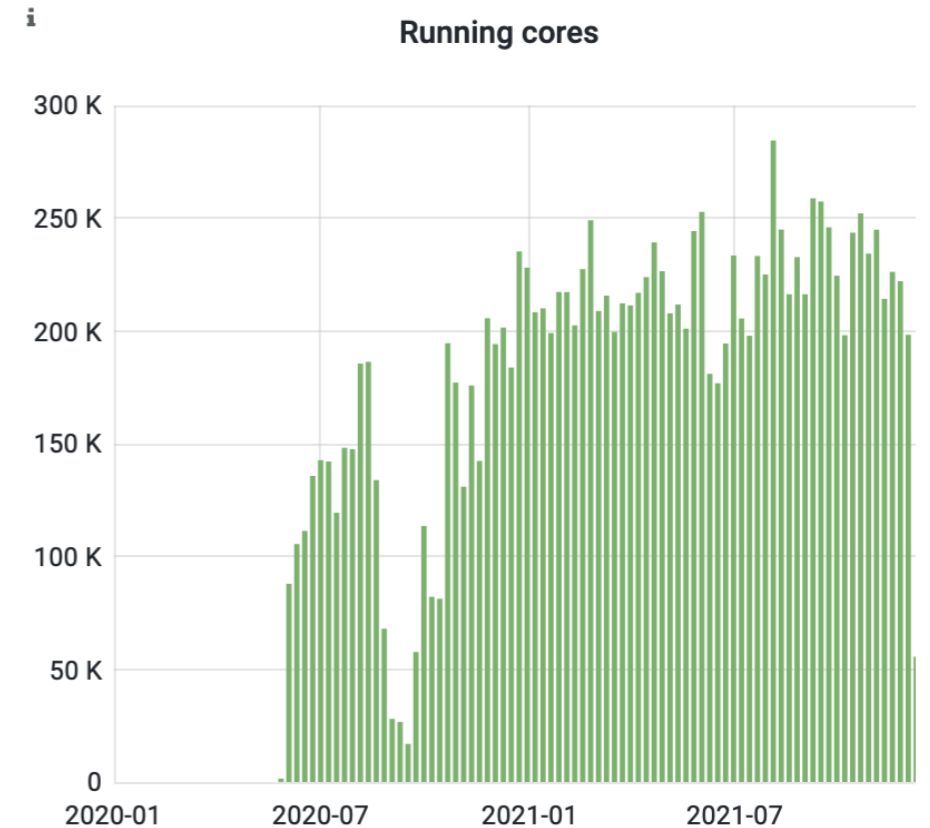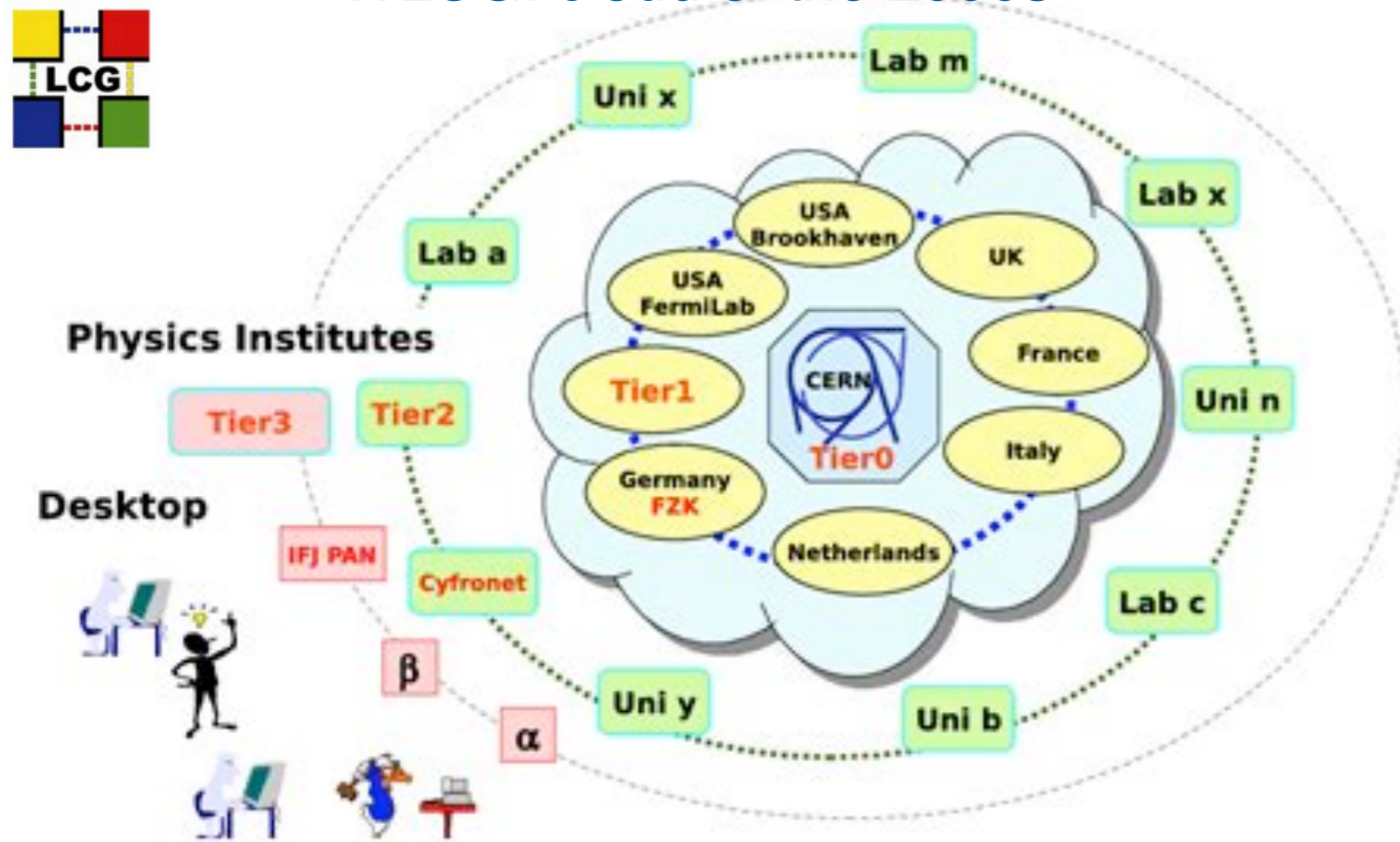WLCG: cloud of the 2000s



**Running cores**

**LCG**

Physics Institutes

Desktop

Uni x
Lab a
USA FermiLab
Tier3  Tier2
Tier1
Germany FZK
IFJ PAN
Cyfronet
β
α
Uni y

Running jobs: 365644
Active CPU cores: 807139
Transfer rate: 21.54 GiB/sec

**☰ Fermilab**

# Centrally managed data



## Primary dataset

Abstract, "what kind of events."
e.g. hard scatter process for simulation, trigger filter for data

### Data tiers

#### AOD
1e5/event

Data columns pertaining to
low-level reconstruction

.root

1e9/file

#### MiniAOD
1e4/event

Calibrated physics objects
Particle-flow candidates

.root

1e9/file

**Data volume**
order of magnitude
[bytes]

🎜 Fermilab

# Centrally managed data



Primary dataset

Abstract, "what kind of events."
e.g. hard scatter process for simulation, trigger filter for data

Data tiers

## AOD
1e5/event

Data columns pertaining to low-level reconstruction

.root
1e9/file

## MiniAOD
1e4/event

Calibrated physics objects
Particle-flow candidates

.root
1e9/file

. . .

## Data volume
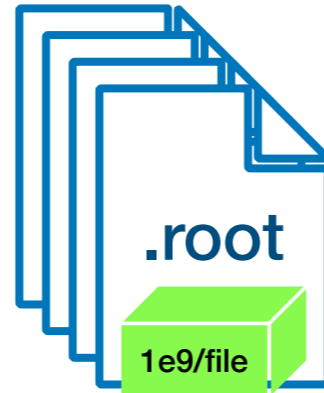order of magnitude [bytes]

🟦 Fermilab

# Primary d

Abstract, "what kind o
e.g. hard scatter proc

## Data tiers

### AOD

Data columns pe
low-level reconst

.root

1e9/file

Accessed

Not accessed



lume

gnitude
s]

🔶 Fermilab

# File format

- Event Data Model (TTree)
- Branch: metadata about C++ data type, basket positions
- Basket: serialized C++ objects stored contiguously*



MiniAODSIM data products

# Simulation processing

🎟 **Fermilab**

# Simulation processing



- Generate events at many sites

🔶 **Fermilab**

# Simulation processing



- Generate events at many sites
- Collect output on disk & tape

🔷 Fermilab

# Data reconstruction



**30k CPU-y**

T1 site

T2 site

T2 site

🔷 **Fermilab**

# Data reconstruction

**30k CPU-y**

- Copy raw input from tape to many sites



T1 site

T2 site

T2 site

🔷 Fermilab

# Data reconstruction



**30k CPU-y**

- Copy raw input from tape to many sites
- Process data



T1 site

T2 site

T2 site

Fermilab

# Data reconstruction

**30k CPU-y**

- Copy raw input from tape to many sites
- Process data
- Collect output on disk & tape

T1 site

T2 site

T2 site

🎗 **Fermilab**

# Data reconstruction

- Copy raw input from tape to many sites
- Process data
- Collect output on disk & tape

**30k CPU-y**

Transfers via File Transfer Service "third-party copy" between storage; orchestrated via Rucio data management



T1 site

T2 site

T2 site

🔷 Fermilab

# Data cataloging

Mario Lassnig

## Rucio main functionalities

**Provides many features that can be enabled selectively**

Findable  Accessible  Interoperable  Reusable

More advanced features →

**Horizontally scalable catalog** for files, collections, and metadata

Transfers between facilities including **disk, tapes, clouds, HPCs**

**Authentication and authorisation** for users and groups

**Many interfaces** available, including CLI, web, FUSE, and REST API

**Extensive monitoring** for all dataflows

Expressive **policy engine** with rules, subscriptions, and quotas

Automated **corruption identification and recovery**

Transparent support for **multihop, caches, and CDN dataflows**

**Data-analytics based flow control**

**Rucio is not a distributed file system, it connects existing storage infrastructure over the network**

No Rucio software needs to run at the data centres **(!)**

Data centres are free to choose which storage system suits them best - **No Vendor Lock-In (!)**

🐿 **Fermilab**

# Data cataloging



2024-05-23 18:13:13 UTC

🐾 Fermilab

# Disk storage at CERN

Vladimir Bayhl

Tape services

Physics services

Sync & Share services

local batch cluster
O($10^5$) cores

CERN Tape Archive

SSD

SSD

EOS

CERNBox

openstack

30-40 GB/s

LHC-B Detector

## EOS

**Total Space**
**600 PB**

**Files Stored**
**~7 Billion**

**# Storage Nodes**
**~1600**

**# Disks**
**~80000**

🔷 **Fermilab**

# Disk storage elsewhere

- Facebook Tectonic FS: one disk cluster per datacenter, two basic workloads:
  - Blob storage: pictures/videos
    - Steady-state IOPS, random access
  - Warehouse: engagement data (clicks/likes)
    - Bursty, more sequential access
- Potential analog:
  - Blob storage: pileup mixing in generation
  - Warehouse: analysis queries
- Many spindles!
  - Load-balance➡performance
  - Scalability via indirection
    - 3 (!) metadata queries /access

- https://www.usenix.org/system/files/fast21-pan.pdf

| Capacity | Used bytes | Files | Blocks | Storage Nodes |
|----------|-----------|-------|--------|---------------|
| 1590 PB | 1250 PB | 10.7 B | 15 B | 4208 |

Table 2: Statistics from a multitenant Tectonic production cluster. File and block counts are in billions.



(a) Aggregate cluster IOPS

(b) Aggregate cluster bandwidth

# Tape storage

Michael Davis



1974: Tape Storage Vault – when robots were human

춘 Fermilab

# Tape storage

🔷 **Fermilab**

# Tape storage



Michael Davis

- Advantages:
  - Reliability
  - Cost
  - Energy efficiency

🟦 **Fermilab**

# Tape storage

Michael Davis

- Advantages:
  - Reliability
  - Cost
  - Energy efficiency

## The Outlook for Tape Technology

### New Advanced Materials

- Very fine magnetic particles
- Smooth surfaces with low friction
- 3D stacking of magnetic particles

Disk technologies are pushing the limits of storage density. Tapes have plenty of room to improve capacity.

- **The cost advantages of tape will increase over time**

Lubricant
Protective layer
Magnetic layer
Substrate
Back coating

50nm    500nm

**Fermilab**

# Projected usage

‡ Fermilab

# Projected usage



**CMS** *Public*
Total Disk
*2022 Estimates*
- No R&D improvements
- Weighted probable scenario
- 10 to 20% annual resource increase

**CMS** *Public*
Total Tape
*2022 Estimates*
- No R&D improvements
- Weighted probable scenario
- 10 to 20% annual resource increase

**CMS** *Public*
Total Disk HL-LHC (2031/No R&D Improvements) fractions
*2022 Estimates*

CACHE: 13%
AODSim: 11%
MINIAOD: 13%
AOD: 12%
ALCARECO: 4%
USER: 4%
SKIM: 7%
RECOSim: 2%
RECO: 5%
RAWSim: 4%
MINIAODSim: 23%
PREMIX: 3%
Other: 5%
NANOAODSim: 3%
OPERATIONS: 10%

**900 PB**

**CMS** *Public*
Total Tape usage HL-LHC (2031/No R&D Improvements) fractions
*2022 Estimates*

HIAOD: 8%
AODSim: 12%
HIRAW: 12%
AOD: 11%
MINIAOD: 2%
ALCARECO: 4%
MINIAODSim: 3%
SKIM: 6%
Other: 4%
RAW: 39%

**2400 PB**

🎗 **Fermilab**

# Disk as a cache

- Disk is expensive (vs. tape)
  - Only MiniAOD, NanoAOD data tiers reliably on disk now
    - Ok because of 10+y experience with detector to know what we need
    - For HL-LHC, new detectors may require more time with low-level information
- Best cache: all the columns you need, none you don't
  - Different set of columns needed for different PDs, analyses
  - Not all rows read if filtering (skimming)
  - **How much can we reduce disk use from PD*tier granularity we have now?**



**CMS** *Public*
Total Disk HL-LHC (2031/No R&D Improvements) fractions
*2022 Estimates*

- CACHE: 13%
- AODSim: 11%
- MINIAOD: 13%
- AOD: 12%
- MINIAODSim: 23%
- ALCARECO: 4%
- USER: 4%
- SKIM: 7%
- NANOAODSim: 3%
- RECOSim: 2%
- RECO: 5%
- OPERATIONS: 10%
- RAWSim: 4%
- Other: 5%
- PREMIX: 3%

900 PB

**CMS** *Public*
Total Tape usage HL-LHC (2031/No R&D Improvements) fractions
*2022 Estimates*

- HIAOD: 8%
- AODSim: 12%
- HIRAW: 12%
- AOD: 11%
- MINIAOD: 2%
- MINIAODSim: 3%
- ALCARECO: 4%
- Other: 4%
- SKIM: 6%
- RAW: 39%

2400 PB

**Fermilab**

# Our inference pipeline



Centrally planned, executed

Analyst / Scientist

RAW

Centrally managed

User-managed

$P(x|\theta)$

x1000

🎇 Fermilab

# Facility view

# Facility view

# Facility view



- Small data (kB-GB)
  - User code, calibration payloads, output histograms, …
- Medium data (GB-TB)
  - Intermediate datasets (skims)
- Large data (TB-PB)
  - Input datasets

🔷 **Fermilab**

# Small data

🔷 **Fermilab**

# Small data

# Small data

🔶 **Fermilab**

# Small data

🔷 **Fermilab**

# Small data

🟰 **Fermilab**

# Small data



- Authorization: unix (token?)
- Logical organization via directories
  - Directory quotas
- Performance: IOPS
  - Read-only mount on workers?

**Mounted**

**Less often mounted**
**Xrootd/https**

Distributed Filesystem

CVM...

REST APIs

🎴 **Fermilab**

# Medium data

# Medium data



- Authorization: token
- Logical organization via provenance
  - Derived dataset catalog?
- Performance: IOPS & Bandwidth
- Lateral movement is non-trivial
  - TPC across facilities?

Xrootd/https

Distributed Filesystem

🔷 **Fermilab**

# Medium data



- Authorization: token
- Logical organization via provenance
  - Derived dataset catalog?
- Performance: IOPS & Bandwidth
- Lateral movement is non-trivial
  - TPC across facilities?

Xrootd/https

Distributed Filesystem

Object Store

**Fermilab**

# Medium data

🔵 **Fermilab**

# Medium data

🔵 **Fermilab**

# Medium data



- A few pros:
  - Cloud-friendly: support industry query platforms
  - More flexible authorization & QoS
- A few cons:
  - Fighting 50y of unix knowledge
  - Existing infrastructure built on top of POSIX(-ish) base layer

🔬 **Fermilab**

# Large data



Xrootd/https access only

Scalable Storage

**Fermilab**

# Large data

🔷 Fermilab

# Large data



- Authorization: token
- Logical organization via DM service
  - User *requests* (subset of) datasets
- Performance: IOPS & Bandwidth
- How to interact with the lake?
  - Data management philosophy

Xrootd/https access only

Cache

Scalable Storage

Data lake Streaming/CDN

JupyterHUB

SSH

USER UI

Jupyter

dask

Batch cluster

HTCondor  slurm  Spark

Local (cloud) resources

AF "seed" of resources

**Fermilab**

# Data management philosophy



Primary dataset

Abstract, "what kind of events."
e.g. hard scatter process for simulation, trigger filter for data

Data tiers

**AOD** — 1e5/event

Data columns pertaining to
low-level reconstruction

.root

1e9/file

**MiniAOD** — 1e4/event

Calibrated physics objects
Particle-flow candidates

.root

1e9/file

. . .

**Data volume**

order of magnitude
[bytes]

**Similar layout for xAOD / PHYSLITE**

🔷 **Fermilab**

# Can we align unit of data access to unit of physics content?

- Dataset = list of 2-4 GB files, totaling 10GB-1PB. Why?



- 1kB
- 100kB
- 10MB
- 1GB
- 100GB

🔷 Fermilab

# Can we align unit of data access to unit of physics content?

- Dataset = list of 2-4 GB files, totaling 10GB-1PB. Why?

- Lower limits:
  - IOPS!!
  - Erasure-code block size ~ 4-16kB
  - Catalog/filesystem overhead ~ 10MB-1GB?

- Upper limits:
  - Third party copy timeout ~ 20 GB
  - Tape cartridge ~ 10 TB

- 1kB    - 100kB    - 10MB    - 1GB    - 100GB

**Fermilab**

# Can we align unit of data access to unit of physics content?

- Dataset = list of 2-4 GB files, totaling 10GB-1PB. Why?

> - TBasket / Page sizes ~ 10-100kB
>   - This is physics-relevant
>   - One float column for O(100k) events
>   - One ragged column for O(10k) events
> - Motivation for byte-range xcache

- 1kB   - 100kB   - 10MB   - 1GB   - 100GB

**辈 Fermilab**

# Can we align unit of data access to unit of physics content?

- Dataset = list of 2-4 GB files, totaling 10GB-1PB. Why?

- Cluster size ~ 10MB
  - All columns pertaining to same group of events
- Good target for read-ahead buffer size
- Do we want to cluster *all* columns though?
  - Typical analysis accesses 10-50%
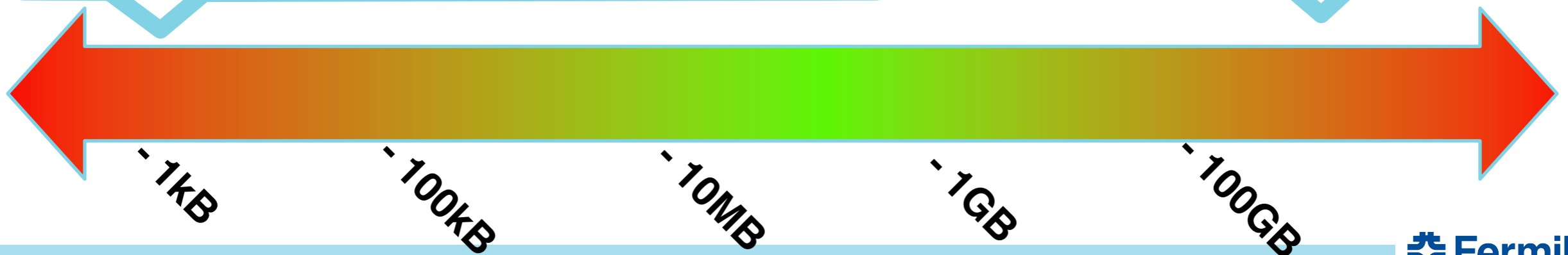  - How will column joins be performed?

- 1kB
- 100kB
- 10MB
- 1GB
- 100GB

🔷 **Fermilab**

# Can we align unit of data access to unit of physics content?

- Dataset = list of 2-4 GB files, totaling 10GB-1PB. Why?

- Sweet spot for access ~ 1MB
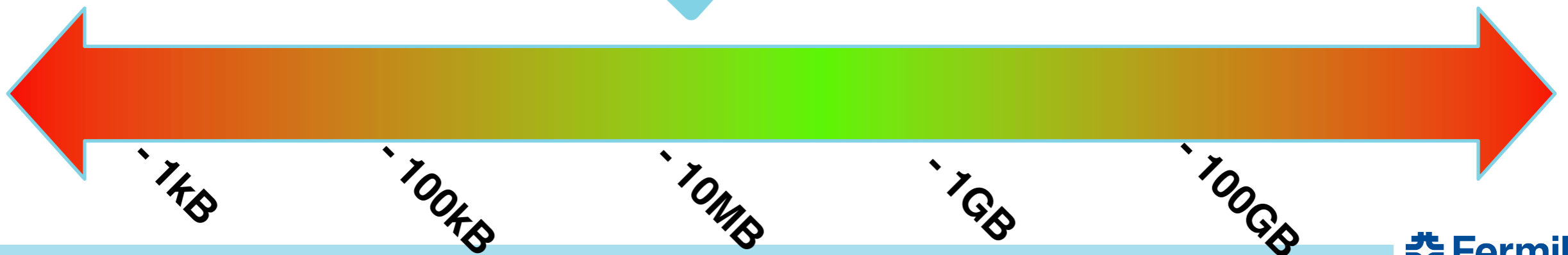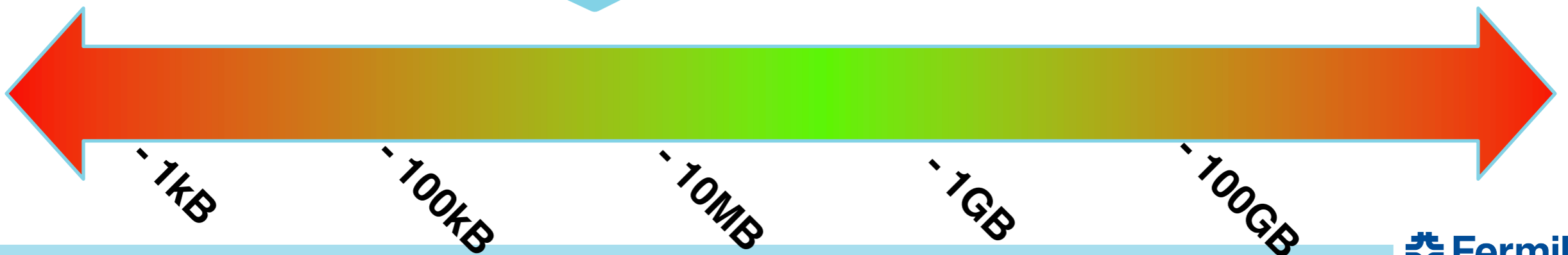  - Few ragged columns for O(10k) events?
  - Many columns for O(1k) events?
  - Do we want small # events per unit?
- Whole-unit cache ➜ off-shelf solutions
- Catalog challenge: need indirection

1kB    100kB    10MB    1GB    100GB

**Fermilab**

# Object store vs. filesystem

- Traditional data storage technology: distributed filesystem
  - e.g. NFS, EOS, dCache, Lustre, HDFS*, …
  - Often with remote access protocol (xrootd)
  - Files are concurrently read/writeable
- Popular new-ish technology: object store
  - Native remote access (http)
  - Objects are immutable (overwrite possible)



[attrib](#)

# Breaking down the ROOT file

- Essentially storing (+ moving) smaller units
  - This is usually a bad thing



Intermodal container

?



Break-bulk cargo

# Breaking down the ROOT file

- Essentially storing (+ moving) smaller units
  - This is usually a bad thing
- Calculated placement
  - Like a hash, client-side
  - Downside: cluster state change causes reshuffle
    - Consistent hashing to minimize movement



Intermodal container

?



Break-bulk cargo



**C**ontrolled
**R**eplication
**U**nder
**S**calable
**H**ashing

OBJECTS

cluster state

OBJECT NAME → PG ID → [OSD.185, OSD.67]

CLUSTER

🟦 Fermilab

# Higher levels of indirection



For intermediate data, (ab)using POSIX filesystem as an implicit data catalog.

Bring Rucio to intermediate data? Does Rucio have sufficient indirection layers?

How do we enable a "facility grid" (cross-facility namespace for intermediate data)
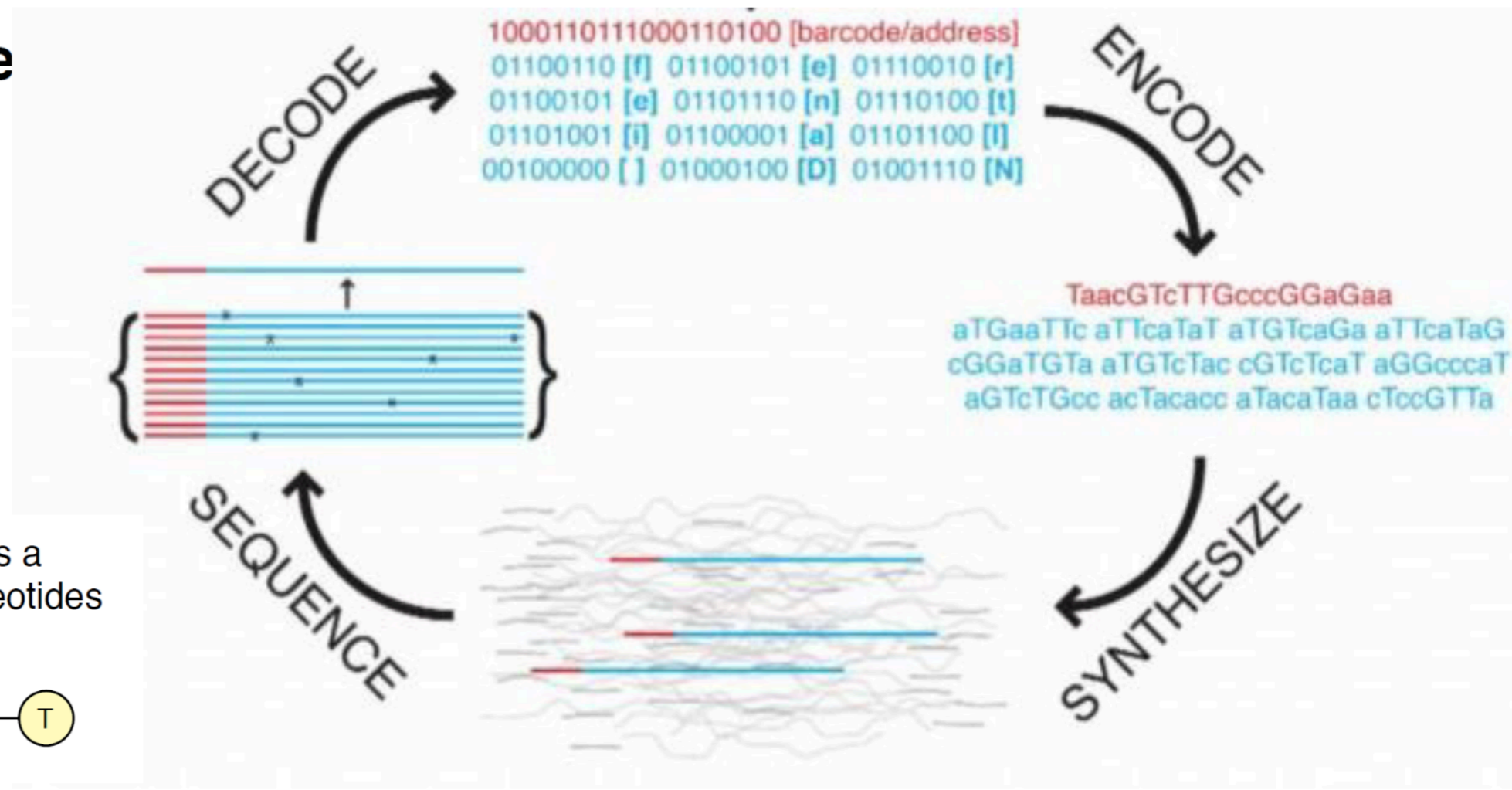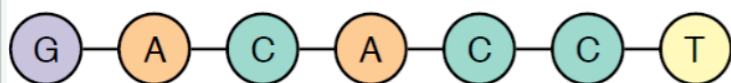
Fermilab

# DNA as a digital storage media

**DNA molecule**

Four nucleotides:

- (A) Adenine
- (C) Cytosine
- (G) Guanine
- (T) Thymine

DNA strand (oligonucleotide) is a linear sequence of these nucleotides

G — A — C — A — C — C — T

# A totally different storage technology

## DNA as a digital storage media

**DNA molecule**

Four nucleotides:

- (A) Adenine
- (C) Cytosine
- (G) Guanine
- (T) Thymine

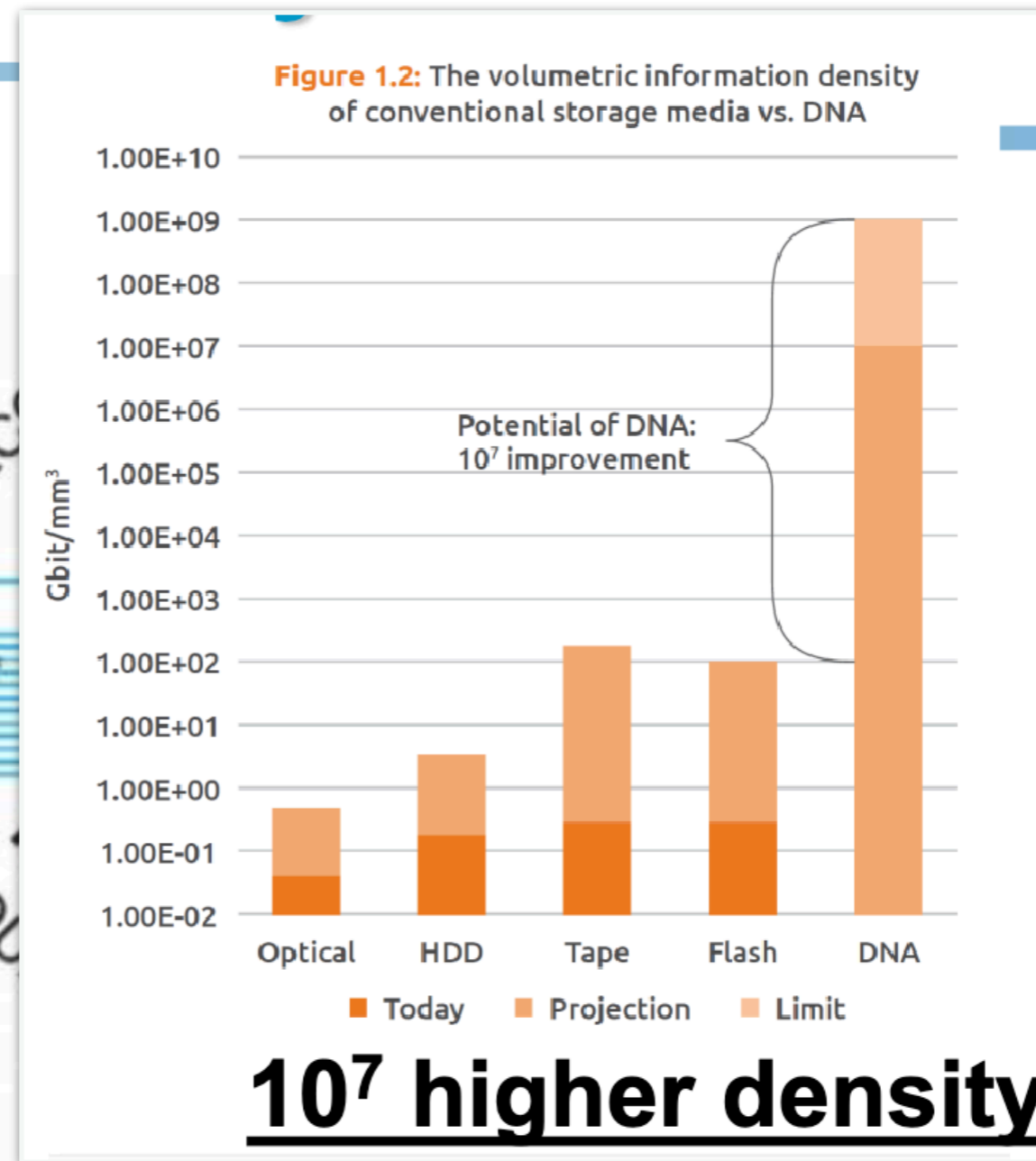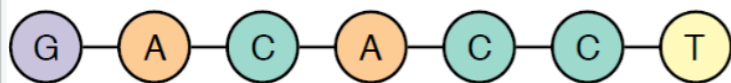DNA strand (oligonucleotide) is a linear sequence of these nucleotides

G - A - C - A - C - C - T

Figure 1.2: The volumetric information density of conventional storage media vs. DNA

Potential of DNA: $10^7$ improvement

Gbit/mm³

| | Optical | HDD | Tape | Flash | DNA |
|---|---|---|---|---|---|

■ Today  ■ Projection  ■ Limit

## $10^7$ higher density

Fermilab

# Conclusion

- Data management for large HEP experiments is a complex topic
- Physics workflows drive requirements
- Always keep an eye out for new technologies

**Fermilab**