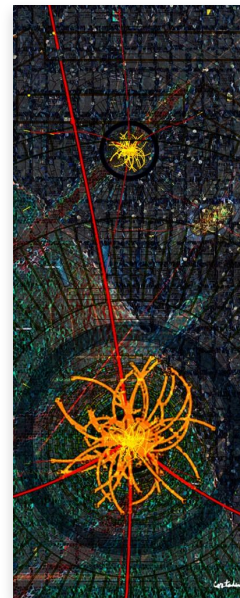
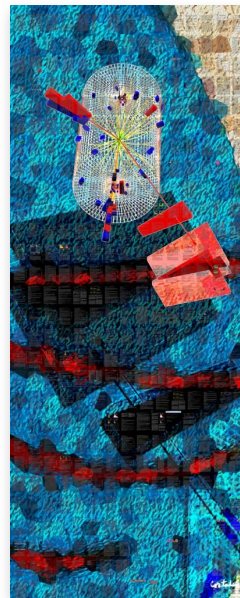
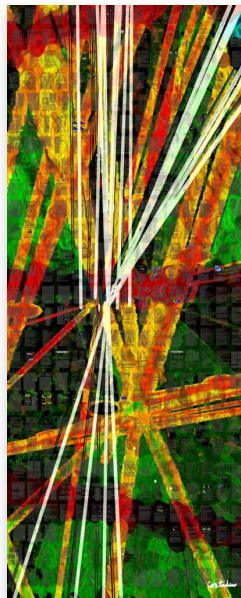
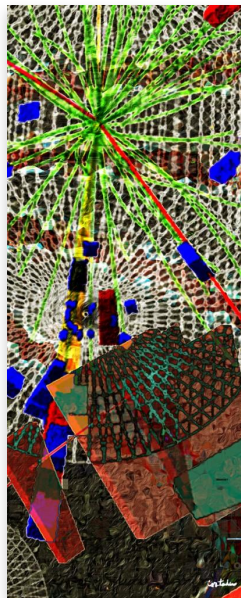
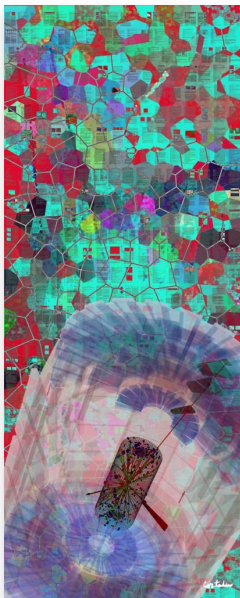


# Computing and Software for Physics

Andrew Melo  
Vanderbilt University  
USCMS Undergraduate Summer Internship  
June 7th, 2024



- Introduction
- Computing's role in CMS
- Challenges
- Accessing data and analyzing it

# About Myself

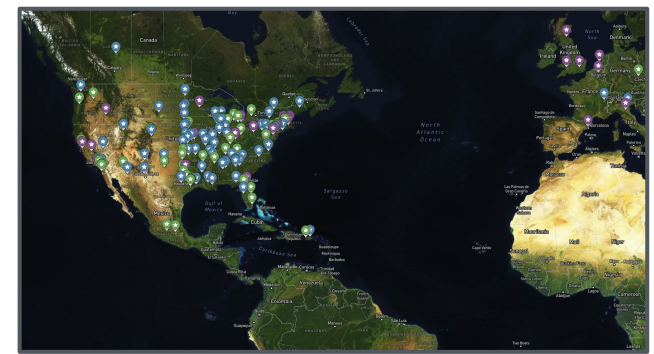
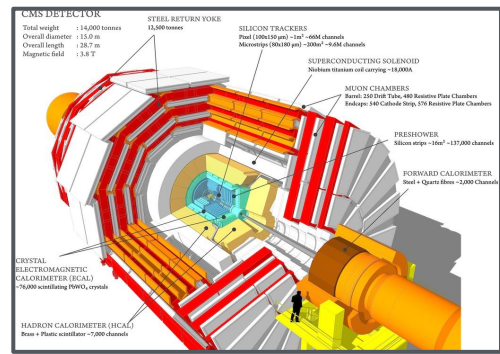
- Research Professor at Vanderbilt University
- Originally wanted to study Computer Science, ended up double-majoring in Physics
- Double-majored in Computer Science and Physics in undergrad
- PhD dissertation a search for evidence for dark matter w/CMS
- Active in CMS computing for nearly 15 years:
  - Contributed to workload management software
  - Maintain the CMS Facility @ Vanderbilt
  - Manages all U.S.-based CMS facilities





# Computing's role in CMS


$$\begin{aligned}
 \mathcal{L}_{3M} = & -\frac{1}{2}\partial_\mu g_\nu^\rho \partial_\mu g_\nu^\rho - g_\nu^{\rho\sigma} \partial_\mu g_\nu^\rho \partial_\mu g_\nu^\rho - \frac{1}{2}g^{\rho\sigma} f^{\mu\nu\alpha} f_{\mu\nu\alpha} g_\nu^\rho g_\nu^\sigma - \partial_\mu W_\nu^+ \partial_\mu W_\nu^- \\
 & - M^2 W_\mu^+ W_\mu^- - \frac{1}{2}\partial_\mu Z_\nu^\rho \partial_\mu Z_\nu^\rho - \frac{1}{2}M^2 Z_\nu^\rho Z_\nu^\rho - \frac{1}{2}\partial_\mu A_\nu \partial_\mu A_\nu - ig_{\text{em}}(\partial_\mu Z_\nu^\rho)(W_\mu^+ W_\nu^- - \\
 & W_\mu^- W_\nu^+) - Z_\nu^\rho(W_\mu^+ \partial_\mu W_\nu^- - W_\mu^- \partial_\mu W_\nu^+) + Z_\nu^\rho(W_\mu^+ \partial_\mu W_\nu^- - W_\mu^- \partial_\mu W_\nu^+) - \\
 & ig_{\text{sw}}(\partial_\mu A_\nu)(W_\mu^+ W_\nu^- - W_\mu^- W_\nu^+) - A_\nu(W_\mu^+ \partial_\mu W_\nu^- - W_\mu^- \partial_\mu W_\nu^+) + A_\nu(W_\mu^+ \partial_\mu W_\nu^- - \\
 & W_\mu^- \partial_\mu W_\nu^+) - \frac{1}{2}g^2 W_\mu^+ W_\nu^- W_\mu^- W_\nu^+ + \frac{1}{2}g^2 W_\mu^- W_\nu^+ W_\mu^+ W_\nu^- + g^2 c_w^2 (Z_\nu^\rho W_\mu^+ Z_\nu^\rho W_\mu^- - \\
 & Z_\nu^\rho Z_\nu^\rho W_\mu^+ W_\mu^-) + g^2 s_w^2 (A_\nu W_\mu^+ A_\nu W_\mu^- - A_\nu A_\nu W_\mu^+ W_\mu^-) + g^2 s_w c_w (A_\nu Z_\nu^\rho)(W_\mu^+ W_\nu^- - \\
 & W_\mu^- W_\nu^+) - 2A_\nu Z_\nu^\rho W_\mu^+ W_\mu^- - \frac{1}{2}\partial_\mu H \partial_\mu H - 2M^2 \alpha_h H^2 - \partial_\mu \phi^+ \partial_\mu \phi^- - \frac{1}{2}\partial_\mu \phi^0 \partial_\mu \phi^0 - \\
 & \beta_h \left( \frac{2M^2}{\Lambda^2} H + \frac{2M}{\Lambda} H + \frac{1}{2}(H^2 + \phi^0 \phi^0 + 2\phi^+ \phi^-) \right) + \frac{2M^2}{\Lambda^2} \alpha_h - \\
 & \frac{1}{8}g^2 \alpha_h (H^4 + (\phi^0)^4 + 4(\phi^+ \phi^-)^2 + 4(\phi^0)^2 \phi^+ \phi^- + 4H^2 \phi^+ \phi^- + 2(\phi^0)^2 H^2) - \\
 & g M W_\mu^+ W_\mu^- H - \frac{1}{2}ig_{\text{sw}} Z_\nu^\rho Z_\nu^\rho H - \\
 & \frac{1}{2}ig(W_\mu^+ (\phi^0 \partial_\mu \phi^- - \phi^- \partial_\mu \phi^0) - W_\mu^- (\phi^0 \partial_\mu \phi^+ - \phi^+ \partial_\mu \phi^0)) + \\
 & \frac{1}{2}g(W_\mu^+ (H \partial_\mu \phi^- - \phi^- \partial_\mu H) + W_\mu^- (H \partial_\mu \phi^+ - \phi^+ \partial_\mu H)) + \frac{1}{2}g \frac{1}{c_w} (Z_\nu^\rho (H \partial_\nu \phi^0 - \phi^0 \partial_\nu H) + \\
 & M \left( \frac{1}{c_w} Z_\nu^\rho \partial_\nu \phi^0 + W_\mu^+ \partial_\mu \phi^- + W_\mu^- \partial_\mu \phi^+ \right) - ig_{\text{em}} M Z_\nu^\rho (W_\mu^+ \phi^- - W_\mu^- \phi^+) + ig_{\text{sw}} M A_\nu (W_\mu^+ \phi^- - \\
 & W_\mu^- \phi^+) - ig \frac{1-2s_w^2}{2c_w} Z_\nu^\rho (\phi^+ \partial_\nu \phi^- - \phi^- \partial_\nu \phi^+) + ig_{\text{sw}} A_\nu (\phi^+ \partial_\nu \phi^- - \phi^- \partial_\nu \phi^+) - \\
 & \frac{1}{2}g^2 W_\mu^+ W_\mu^- (H^2 + (\phi^0)^2 + 2\phi^+ \phi^-) - \frac{1}{2}g^2 \frac{1}{c_w^2} Z_\nu^\rho Z_\nu^\rho (H^2 + (\phi^0)^2 + 2(2s_w^2 - 1)\phi^+ \phi^-) - \\
 & \frac{1}{2}g^2 \frac{1}{c_w} Z_\nu^\rho \partial_\nu \phi^0 (W_\mu^+ \phi^- + W_\mu^- \phi^+) - \frac{1}{2}ig^2 \frac{1}{c_w} Z_\nu^\rho H (W_\mu^+ \phi^- - W_\mu^- \phi^+) + \frac{1}{2}g^2 s_w A_\nu \phi^0 (W_\mu^+ \phi^- + \\
 & W_\mu^- \phi^+) + \frac{1}{2}ig^2 s_w A_\nu H (W_\mu^+ \phi^- - W_\mu^- \phi^+) - g^2 \frac{1}{c_w} (2c_w^2 - 1) Z_\nu^\rho A_\nu \phi^+ \phi^- - \\
 & g^2 s_w^2 A_\nu A_\nu \phi^+ \phi^- + \frac{1}{2}ig_s \lambda_3 (\bar{q}^i \gamma^\mu q^j) g_\mu^2 - e^4 (\gamma \partial + m_2) e^2 - \nu^4 (\gamma \partial + m_2) \nu^4 - u_2^4 (\gamma \partial +
 \end{aligned}$$




After the collisions occur and are recorded by the detector, it's still not in a format usable by analyzers (you!). The data needs to be:

- Reconstructed
- Slimmed down for a format useful for analysis
- Transmitted to our "grid" of ~80 facilities in 50 countries
- Be processed by analysis software who transmits the final results to the end user

EUROPEAN ORGANIZATION FOR NUCLEAR RESEARCH (CERN)





CERN-PH-EP/2012-220  
2013/01/29

---

CMS-HIG-12-028

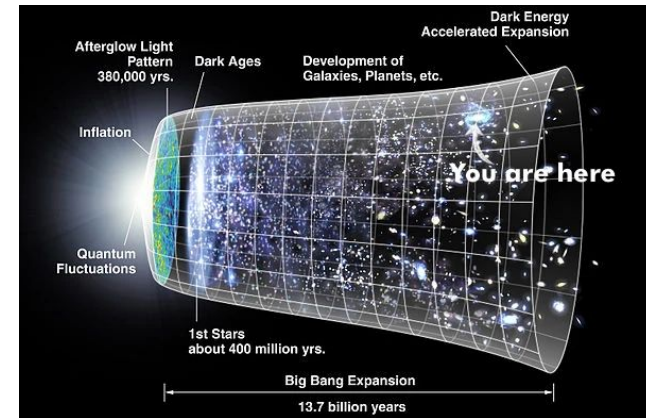
Observation of a new boson at a mass of 125 GeV with the CMS experiment at the LHC

The CMS Collaboration

# Scale of CMS Computing

“Space is big. You just won't believe how vastly, hugely, mind-bogglingly big it is. I mean, you may think it's a long way down the road to the chemist's, but that's just peanuts to space.”

-Douglas Adams, The Hitchhiker's Guide to the Galaxy



CMS produces, stores and processes an enormous amount of data:

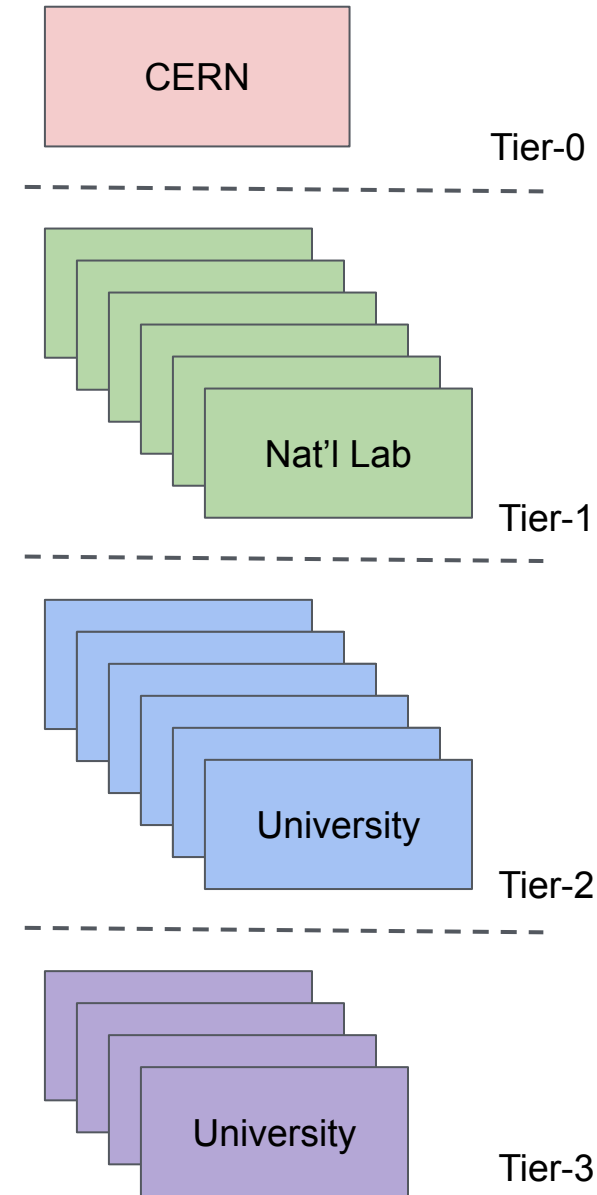
- The detector has 100s of millions of channels, and collisions happen 40 million times per second – In theory, this is  $10^{17}$  bytes every second
- The trigger filters down these 40 million events/sec into a more manageable 2,000 events/sec – but even after this reduction, the amount of data remaining is still several gigabytes a second, approximately a DVDs worth of data for every second the detector runs
  - CMS has recorded 100,000,000 gigabytes of data, which is all securely archived on two separate tape backup systems
- Each collision can take up to a full minute to process on a CPU. In total, we use 4.8 billion CPU-hours each year
- And this doesn't even consider all the simulated collisions we need to produce

# CMS Computing Facilities

It's not feasible to have one large facility for all of CMS' computing needs. We instead have a large number of globally-distributed and interconnected facilities, called "The Grid". We divide the grid into four different tiers based on size:

- Tier-0: At CERN, has trigger and a tape copy of all RAW data produced by the detector
- Tier-1: At 7 national laboratories. The Tier-1s store an additional tape copy of all RAW data
- Tier-2: 60 large university facilities. Where most user analysis is run
- Tier-3: 20 Smaller facilities usually purpose-built by the institution. Fermilab is the exception with a very large Tier-3 for users

The sites are all connected by fast networking links, to allow the data to be quickly moved between sites



# Moving Data Between Sites

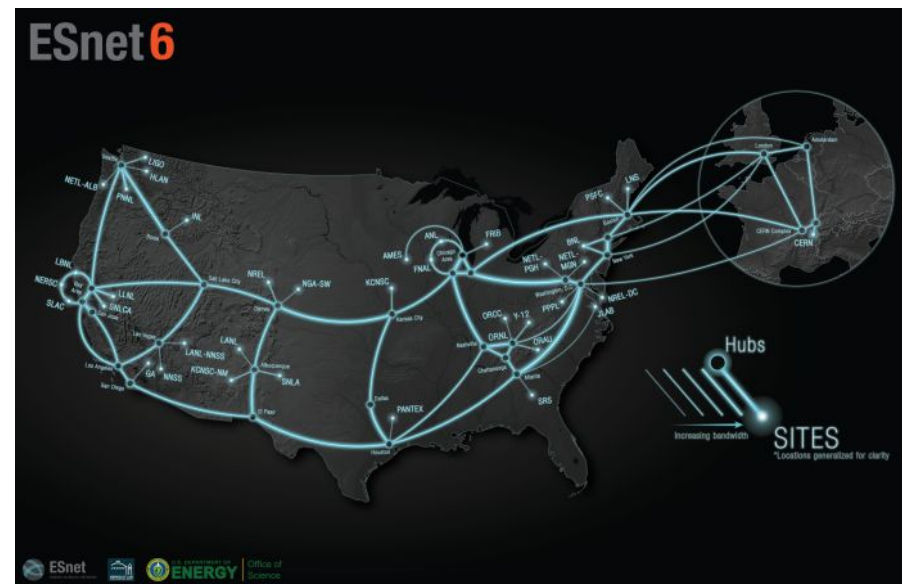
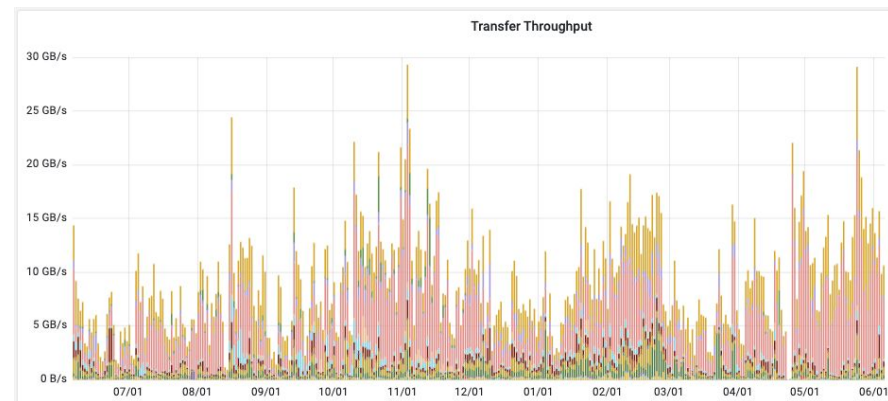
With the sites being distributed globally, a key problem being able to move the data around globally.

- Where possible, we “move computation to the data”

Within the USA, ESNet runs a private, Research and Education-only portion of the internet to connect all DOE labs and many Universities together

Even when the experiment isn't running, CMS is moving several Gigabytes/sec across the network

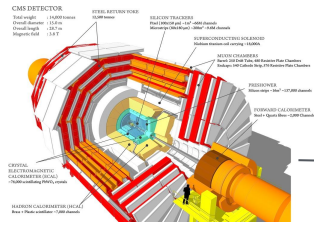
Within the next decade, ESNet will have a Terabit/sec of bandwidth Across the atlantic



More Info:  
<https://www.es.net/welcome-esnet6/>



# Preparing for Analysis: The Data Processing Workflow



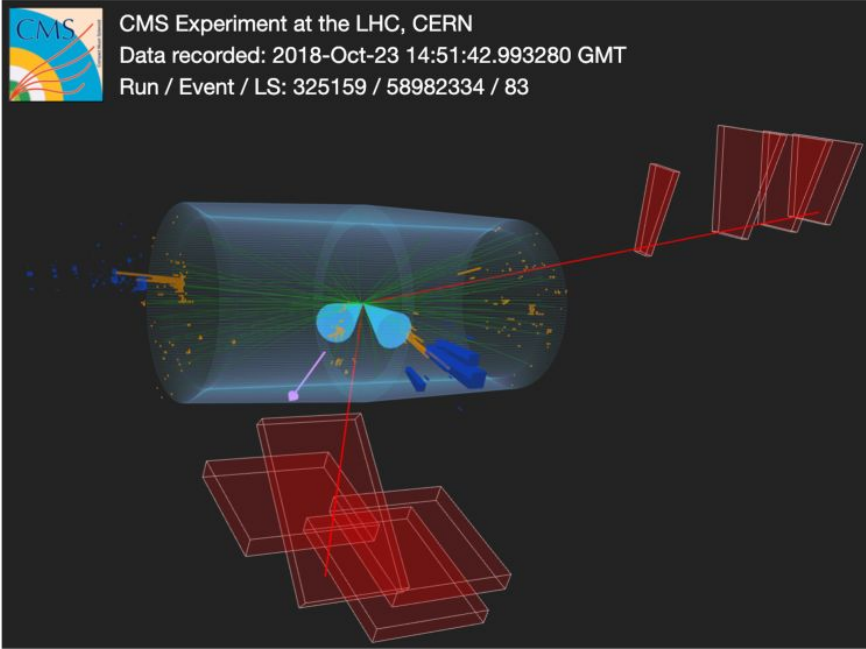
GEN-SIM-DIGI

You are here!

The raw data (simulated or real) is not in a format very useful for analysis.

- The detector data is points in 3D-space where particles interacted w/sensors
- We want to study *particles*, so we have to do *reconstruction* to “connect the dots”

Visit <https://tinyurl.com/cmsfnal23> for an interactive event browser



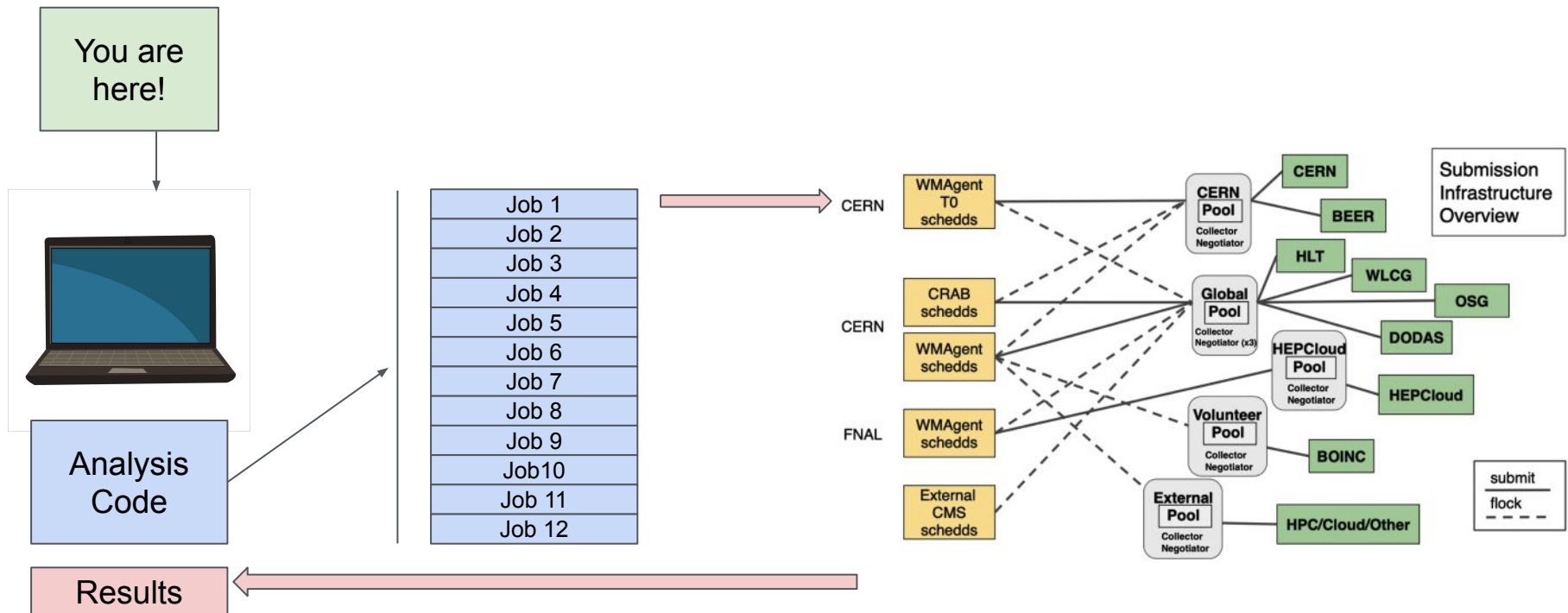


# Processing Data at Scale

There is far too much data to reasonably process on a single computer. To process data at scale, we have to

- Split the data into smaller chunks
- Process each chunk on a separate machine
- Combine all the results together

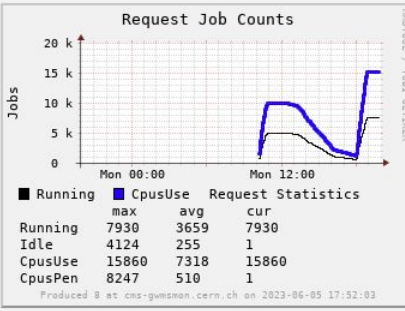
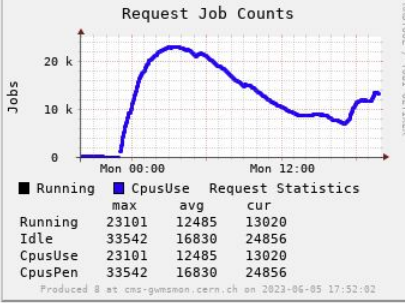
To enable this, we developed The Global Pool, which combines all CMS' compute hardware into a single "virtual computer"



# Processing Data at Scale

- Through the global pool, users can access 10s of thousands of CPUs
- HTCondor is used as our job scheduling system and is the entry point for all jobs to execute on CMS resources
- To be successful, you will become familiar with HTCondor

Idle Jobs/CpusPen	207995 / 396591
Unique Pressure	139623
Running Jobs/CpusInUse	51978 / 78985
Not Queued Jobs	103019
User Count	65
Tasks Running	1261
Tasks Held	20136
Update Time	6/5/2023, 12:51:46 PM

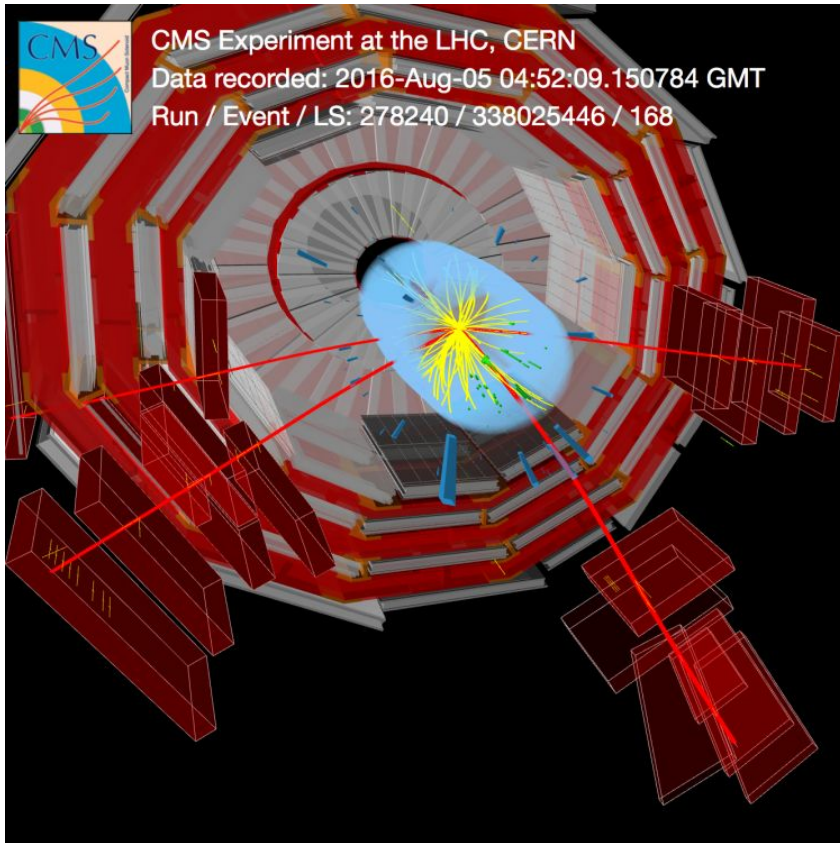
Username	Total	Running	Idle	CpusInUse	CpusPen	Tasks Running	Tasks Idle	Matching Site Count	User History	All ta																				
srosenzw	7,932	7,932	0	15,864	0	14	0	67	 <p>Request Job Counts</p> <p>Jobs</p> <p>Mon 00:00 Mon 12:00</p> <p>Running CpusUse Request Statistics</p> <table border="1"> <tr> <td></td> <td>max</td> <td>avg</td> <td>cur</td> </tr> <tr> <td>Running</td> <td>7930</td> <td>3659</td> <td>7930</td> </tr> <tr> <td>Idle</td> <td>4124</td> <td>255</td> <td>1</td> </tr> <tr> <td>CpusUse</td> <td>15860</td> <td>7318</td> <td>15860</td> </tr> <tr> <td>CpusPen</td> <td>8247</td> <td>510</td> <td>1</td> </tr> </table> <p>Produced @ cms-gwmsmon.cern.ch on 2023-06-05 17:52:03</p>		max	avg	cur	Running	7930	3659	7930	Idle	4124	255	1	CpusUse	15860	7318	15860	CpusPen	8247	510	1	
	max	avg	cur																											
Running	7930	3659	7930																											
Idle	4124	255	1																											
CpusUse	15860	7318	15860																											
CpusPen	8247	510	1																											
iaandree	37,875	13,020	24,855	13,020	24,855	44	0	39	 <p>Request Job Counts</p> <p>Jobs</p> <p>Mon 00:00 Mon 12:00</p> <p>Running CpusUse Request Statistics</p> <table border="1"> <tr> <td></td> <td>max</td> <td>avg</td> <td>cur</td> </tr> <tr> <td>Running</td> <td>23101</td> <td>12485</td> <td>13020</td> </tr> <tr> <td>Idle</td> <td>33542</td> <td>16830</td> <td>24856</td> </tr> <tr> <td>CpusUse</td> <td>23101</td> <td>12485</td> <td>13020</td> </tr> <tr> <td>CpusPen</td> <td>33542</td> <td>16830</td> <td>24856</td> </tr> </table> <p>Produced @ cms-gwmsmon.cern.ch on 2023-06-05 17:52:02</p>		max	avg	cur	Running	23101	12485	13020	Idle	33542	16830	24856	CpusUse	23101	12485	13020	CpusPen	33542	16830	24856	
	max	avg	cur																											
Running	23101	12485	13020																											
Idle	33542	16830	24856																											
CpusUse	23101	12485	13020																											
CpusPen	33542	16830	24856																											

# Analyzing CMS Data



# Expectation vs Reality

Many think that looking at CMS data is exactly that – visualizing events and drawing conclusions. There’s way too many events to look at them manually. Instead we write software to examine the data for us and look at those results



```

~/projects/xrootd-multiuser/src — ssh brazil.accre.vanderbilt.edu — 84x51
71 0xcc2b1d17, 0xc8ea00a0, 0xd6ad50a5, 0xd26c4d12, 0xdf2f6bcb,
72 0xdbee767c, 0xe3a1cbc1, 0xe760d676, 0xea23f0af, 0xee2ed18,
73 0xf0a5bd1d, 0xf464a0aa, 0xf9278673, 0xfde69bc4, 0x89b8fd09,
74 0x8d79e0be, 0x803ac667, 0x84fbd0b0, 0x9abc8bd5, 0x9e7d9662,
75 0x933eb0bb, 0x97ffad0c, 0xafb010b1, 0xab710d06, 0xa6322bdf,
76 0xa2f33668, 0xbcb4666d, 0xb8757bda, 0xb5365d03, 0xb1f740b4
77 };
78
79
80 static std::string
81 human_readable EVP(const unsigned char *evp, size_t length)
82 {
83     unsigned int idx;
84     std::string result; result.reserve(length*2);
85     for (idx = 0; idx < length; idx++)
86     {
87         char encoded[3];
88         sprintf(encoded, "%02x", evp[idx]);
89         result += encoded;
90     }
91     return result;
92 }
93
94
95 ChecksumState::ChecksumState(unsigned digests)
96 : m_digests(digests),
97   m_cksum(0),
98   m_crc32(crc32(0, NULL, 0)),
99   m_adler32(adler32(0, NULL, 0)),
100  m_md5_length(0),
101  m_cur_chunk_bytes(0),
102  m_offset(0),
103  m_md5(NULL),
104  m_file_sha1(NULL),
105  m_chunk_sha1(NULL)
106 {
107     if (digests & ChecksumManager::MD5)
108     {
109         m_md5 = EVP_MD_CTX_create();
110         EVP_DigestInit_ex(m_md5, EVP_md5(), NULL);
111     }
112     if (digests & ChecksumManager::CVMFS)
113     {
114         m_file_sha1 = EVP_MD_CTX_create();
115         EVP_DigestInit_ex(m_file_sha1, EVP_sha1(), NULL);
116         m_chunk_sha1 = EVP_MD_CTX_create();
117         EVP_DigestInit_ex(m_chunk_sha1, EVP_sha1(), NULL);
118     }
119 }
XrdChecksumCalc.cc 116,1 29%

```



To analyze CMS data one needs to

1. Access a CMS compute cluster via SSH
2. Develop and test code to process each event
3. Make HTCondor configuration files
4. Submit jobs to the grid

There will be more detailed talks/tutorials on each of these steps in the coming days!

# Accessing Fermilab LPC Cluster via SSH

First step to using CMS' compute resources is to connect to them via SSH

SSH stands for “Secure SHell” and allows you to connect to a remote computer and type into it like you were physically there

CMS interactive machines come preconfigured with all the necessary software to write code and process data

```
1 → [meloam][melo-mbp][~]
      ssh lxplus.cern.ch
* *****
* Welcome to lxplus711.cern.ch, CentOS Linux release 7.9.2009 (Core)
* Archive of news is available in /etc/motd-archive
* Reminder: you have agreed to the CERN
*   computing rules, in particular OC5. CERN implements
*   the measures necessary to ensure compliance.
*   https://cern.ch/ComputingRules
* Puppet environment: production, Roger state: production
* Foreman hostgroup: lxplus/nodes/login
* Availability zone: cern-geneva-a
* LXPLUS Public Login Service - http://lxplusdoc.web.cern.ch/
* An AlmaLinux8 based lxplus8.cern.ch is now available
* An AlmaLinux9 based lxplus9.cern.ch is now available
* Please read LXPLUS Privacy Notice in http://cern.ch/go/TpV7
* *****
-bash: cite: command not found
-bash: cite: command not found
-bash: cite: command not found
2 → [meloam][lxplus711][~]
      hostname
lxplus711.cern.ch
```

[https://uscms.org/uscms\\_at\\_work/computing/getstarted/uaf.shtml](https://uscms.org/uscms_at_work/computing/getstarted/uaf.shtml)



# Develop code to process each event

To process each event, a custom hand-written program with the analysis is run over all the data.

This code is unique to each analysis (though groups with similar analyses might share large parts)

Often, this analysis code is written in C++ which means that a typical day involves repeating:

- The code is written in a text editor (like VIM)
- The code is *compiled* into an executable
- The executable is then run over the data

```
void MyClass::Loop() {  
    size_t nEvents;  
    // load...  
  
    for (Long64_t iEvent=0; iEvent<nEvents; iEvent++) {  
        double MET_pt;  
        int nElectron;  
        double * Electron_pt;  
        double * Electron_eta;  
        // load...  
  
        if ( MET_pt > 100. ) continue;  
  
        for(size_t iEl=0; iEl<nElectron; ++iEl) {  
            if ( Electron_pt[iEl] > 30. ) {  
                hist->Fill(Electron_eta[iEl]);  
            }  
        }  
    }  
}
```

VIM Text Editor:

<https://linuxconfig.org/vim-tutorial>

# Version Control with Git

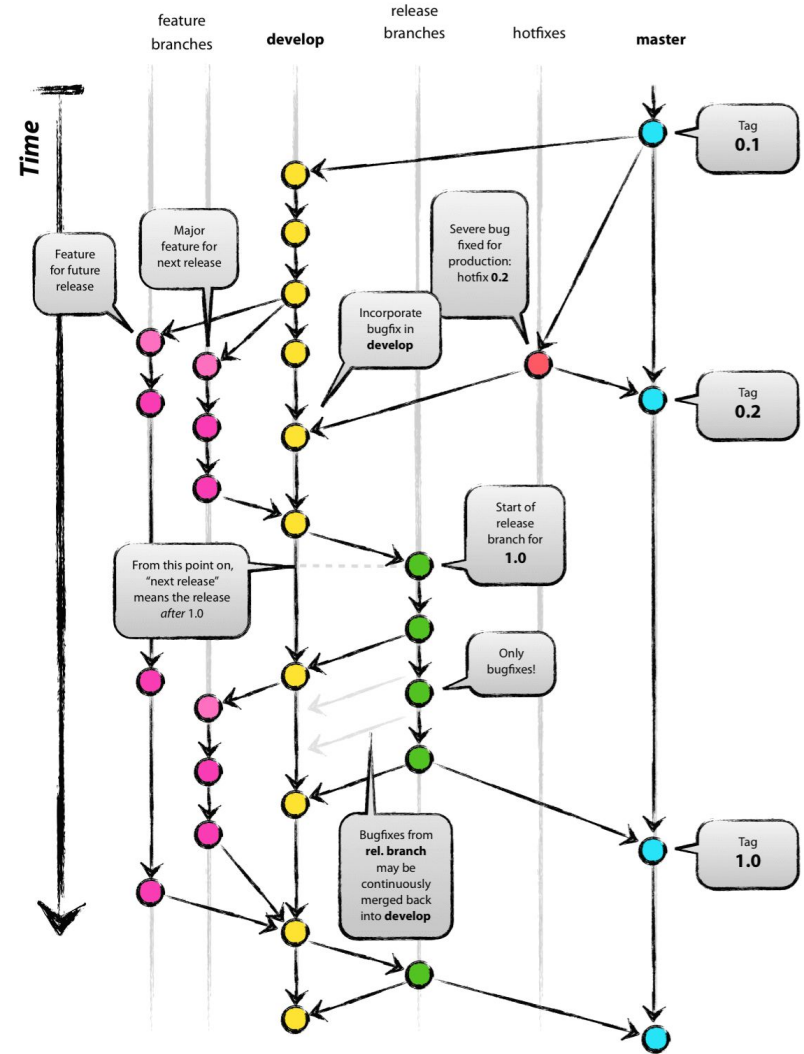
When developing an analysis, its important to:

- Collaborate with others
- Document the history of the code

*Version Control* enables both of these. Within CMS, we typically use the *Git* version control system

A typical Git workflow is:

- **checkout** the code from Git
- Write a new feature
- **commit** the code to Git
- **push** the code to a central location



See also: <https://git-scm.com/video/what-is-version-control>

# Make HTCondor Config Files

With the code running functioning on a single machine, the next step is to write a configuration to tell HTCondor how to run our executable.

This configuration, called a JDL specifies:

- The executable to run
- Where to store the log files
- What files to move to/from the grid
- What arguments to pass to the executable
- How many copies of the job to run

HTCondor JDL

```
universe = vanilla
Executable = sleep.sh
should_transfer_files = YES
when_to_transfer_output = ON_EXIT
Output = sleep_$(Cluster)_$(Process).stdout
Error = sleep_$(Cluster)_$(Process).stderr
Log = sleep_$(Cluster)_$(Process).log
Arguments = 60
Queue 2
```

sleep.sh

```
#!/bin/bash
set -x
# Sleep
sleep $1
echo "##### HOST DETAILS #####"
echo "I slept for $1 seconds on:"
hostname
date
```

More Info:

[https://uscms.org/uscms\\_at\\_work/computing/setup/batch\\_systems.shtml](https://uscms.org/uscms_at_work/computing/setup/batch_systems.shtml)



# Submit Jobs to the Grid

With the executable and the JDL, it's time to submit to the Grid.

There are a number of HTCondor commands, but the most important two are:

- `condor_submit`
- `condor_q`

We hope that the Grid works 100% every time, but unfortunately there are several ways for the jobs to fail, so some “hand holding” has to be done to recover from failed jobs.

```
[username@cmslpc132 ~]$ condor_submit sleep-condor.jdl
Querying the CMS LPC pool and trying to find an
available schedd...
Attempting to submit jobs to lpcschedd3.fnal.gov
Submitting job(s)..
2 job(s) submitted to cluster 76596545.
```

```
[username@cmslpc132 condor]$ condor_q
```

```
-- Schedd: lpcschedd3.fnal.gov : <131.225.188.235:9618?... @ 09/08/22 15:12:37
```

ID	OWNER	SUBMITTED	RUN_TIME	ST	PRI	SIZE	CMD
76596545.0	tonjes	9/8 15:12	0+00:00:00	I	0	0.0	RunAN.sh 2018 MC

```
Total for query: 1 jobs; 0 completed, 0 removed, 1 idle, 0 running, 0 held, 0 suspended
```

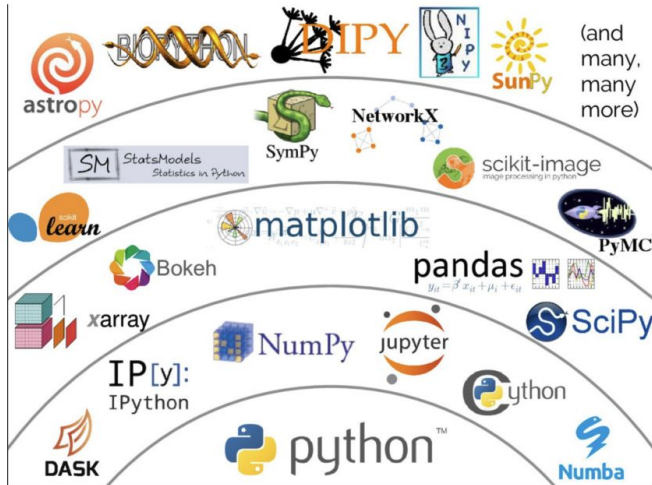
```
Total for tonjes: 1 jobs; 0 completed, 0 removed, 1 idle, 0 running, 0 held, 0 suspended
```

```
Total for all users: 1208 jobs; 0 completed, 41 removed, 898 idle, 203 running, 66 held, 0 suspended
```

# Notebooks and Python-based Analysis

In addition to the “traditional” analysis tools, CMS has begun to embrace the python-based Data Science tools embraced by industry.

Several of these “Analysis Facilities” are being developed, which use industry-standard tools like Pandas, Numpy and Dask



[img](#)

File Edit View Run Kernel Tabs Settings Help

20-z-peak.ipynb hats-spark-2022

```
sys.stdout.write('.')
print("done")
```

Loading DataFrames .....done

## Combine DataFrames

The previous step yields a DataFrame per-dataset, which can be unwieldy to manage. Spark's `union` function concatenates two DataFrames together. Let's loop over everything and make a DataFrame for MC and another for data.

```
[4]: # Monte-Carlo
if mcDFList:
    mcDF = mcDFList[0][0]
    for x in mcDFList[1:]:
        mcDF = mcDF.union(x[0])
    print("MC Partitions: %d" % mcDF.rdd.getNumPartitions())

# Data
if dataDFList:
    dataDF = dataDFList[0][0]
    for x in dataDFList[1:]:
        dataDF = dataDF.union(x[0])
    print("Data Partitions: %d" % dataDF.rdd.getNumPartitions())
```

MC Partitions: 7576  
Data Partitions: 10515

## Drop unneeded columns

Spark is intelligent about how it executes queries, so any unneeded columns are simply not read (similar to how ROOT handles reading branches). That being said, these ntuples have 300 columns, so functions like `printSchema()` and `show()` are unintelligible. Execute the following cell to remove branches that won't be useful for this tutorial.

```
[5]: reqCol = ['Muon_pt', 'Muon_eta', 'Muon_phi', 'Muon_energy',
             'Muon_charge', 'Muon_loose', 'Muon_tight', 'Muon_isGlobal',
             'Jet_pt', 'Jet_eta', 'Jet_phi', 'Jet_energy',
             'Met_type1PF_pt', 'Met_type1PF_px', 'Met_type1PF_py',
             'Met_type1PF_pz', 'Met_type1PF_phi',
             'eff', 'xsec', 'Category', 'SubCategory']
trimMCDF = mcDF.select(reqCol)
```

Simple 0 3 hats-spark-2022 | Idle Mode: Command Ln 1, Col 1 20-z-peak.ipynb



Python notebooks have several advantages:

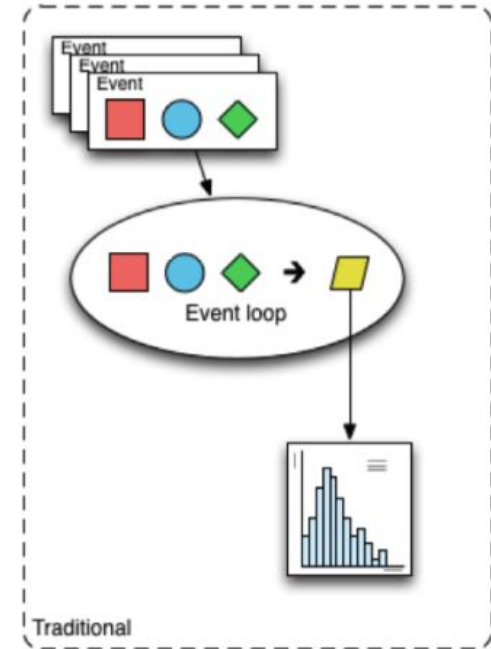
- Plots and histograms can be displayed in-line with the code, making it obvious what code generated which plots
- The Python syntax can be much simpler to work with than the equivalent C++ code
- Because we're using libraries and packages from other fields, we take advantage of the significant effort others are putting into Data Science

The python code can also be used outside of notebooks, if one is accessing a more “traditional” cluster

# Columnar Analysis

In a “traditional” analysis, each event is processed one-at-a-time, meaning:

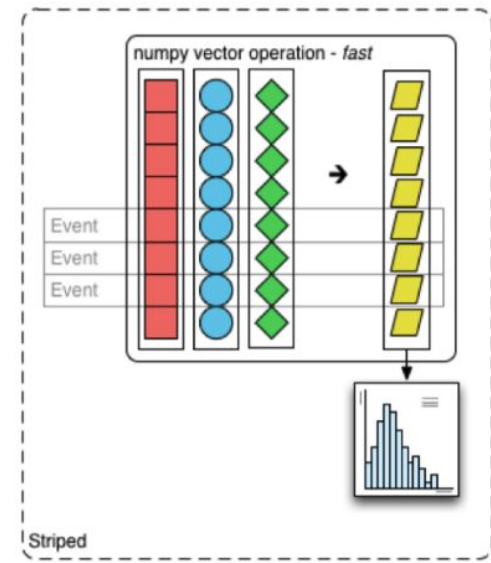
- An event is loaded and the appropriate values are extracted
- Some computation happens over this single event
- The temporary space is cleaned and the next event is loaded



Traditional

In a “columnar” analysis, whole *batches* of events are processed at once, meaning

- 100s or even 1000s of events are loaded at once
- A bulk computation is done over the whole batch
- The temporary space is cleaned and the next event is loaded



Striped

The lower overhead for this method is particularly attractive for Data Scientists and is being embraced by CMS

# GPUs and Accelerators

# Why GPUs?

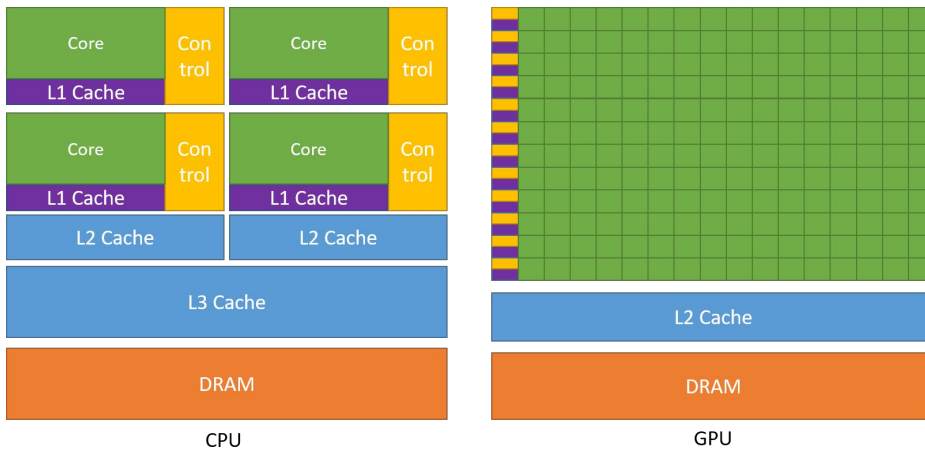
- We discussed rates and pileup before
- The time to process an event increases (roughly) by the square of the pileup
- During HL-LHC, the pileup will increase by a factor of two
- This means (roughly), the time to process each event increases by a factor of four!!
- Due to other affects, the worst-case estimate for HL-LHC is a need for 10x the amount of CPUs
- Will be very hard to make this work financially, need to speed things up!



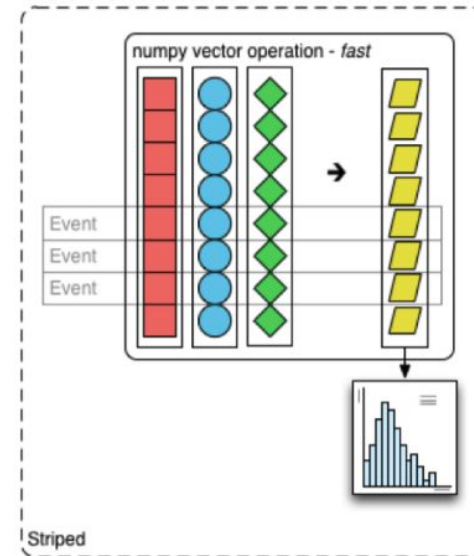
One path CMS is exploring to accelerate our workflows is to exploit GPUs

GPUs are much faster both per-dollar and per-watt of performance, however you can't simply take code written for CPUs and run them on a GPU, you need to rewrite them

Several people are targeting the most-time consuming portions of our workflow to port to GPUs, and they have reached 30% of the modules, with more improvements to come



Accelerator	Architecture	Socket	FP32 CUDA Cores	FP64 Cores (excl. Tensor)	M INT: C
H100	Hopper	SXM5	16896	4608	168
A100 80GB	Ampere	SXM4	6912	3456	691
A100 40GB	Ampere	SXM4	6912	3456	691



[img](#)

Recently, the use of GPUs for AI/ML has exploded

- ChatGPT, etc are in the news constantly

There are a number of efforts within CMS to use Machine Learning (ML) to improve all phases of data processing, including:

- Anomaly detection for detector performance
- “Tagging” physics objects in the detector
- Full end-to-end ML-based reconstruction

There are a lot of opportunities to exploit ML within CMS, but a lot of questions remain from how to support them at the sites to how do we validate the results that come out of these algorithms.

# The Future is Bright!

...come join us!

# Questions?