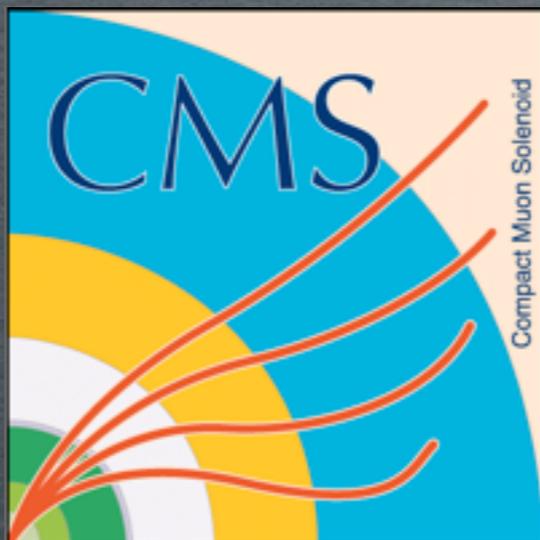


SCALING UP CMS TIER-3 DATA PROCESSING

KEVIN LANNON, MIKE HILDRETH
ON BEHALF OF ND CMS GROUP

Last year's Workshop: Data Preservation

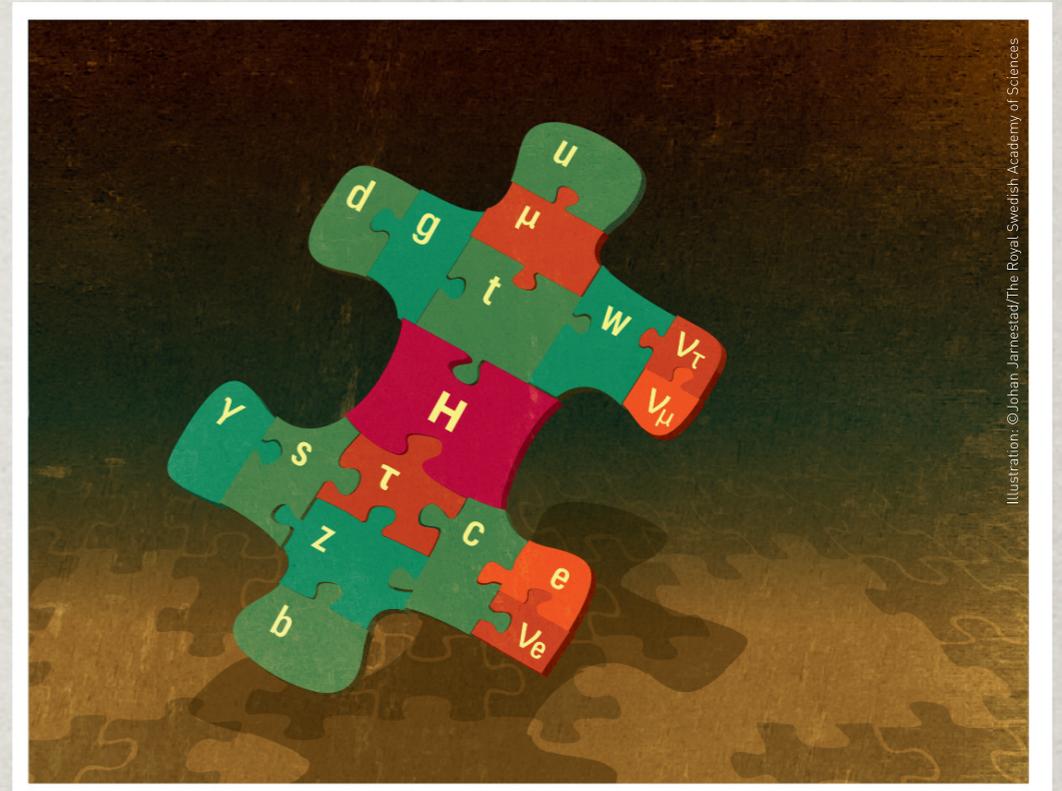
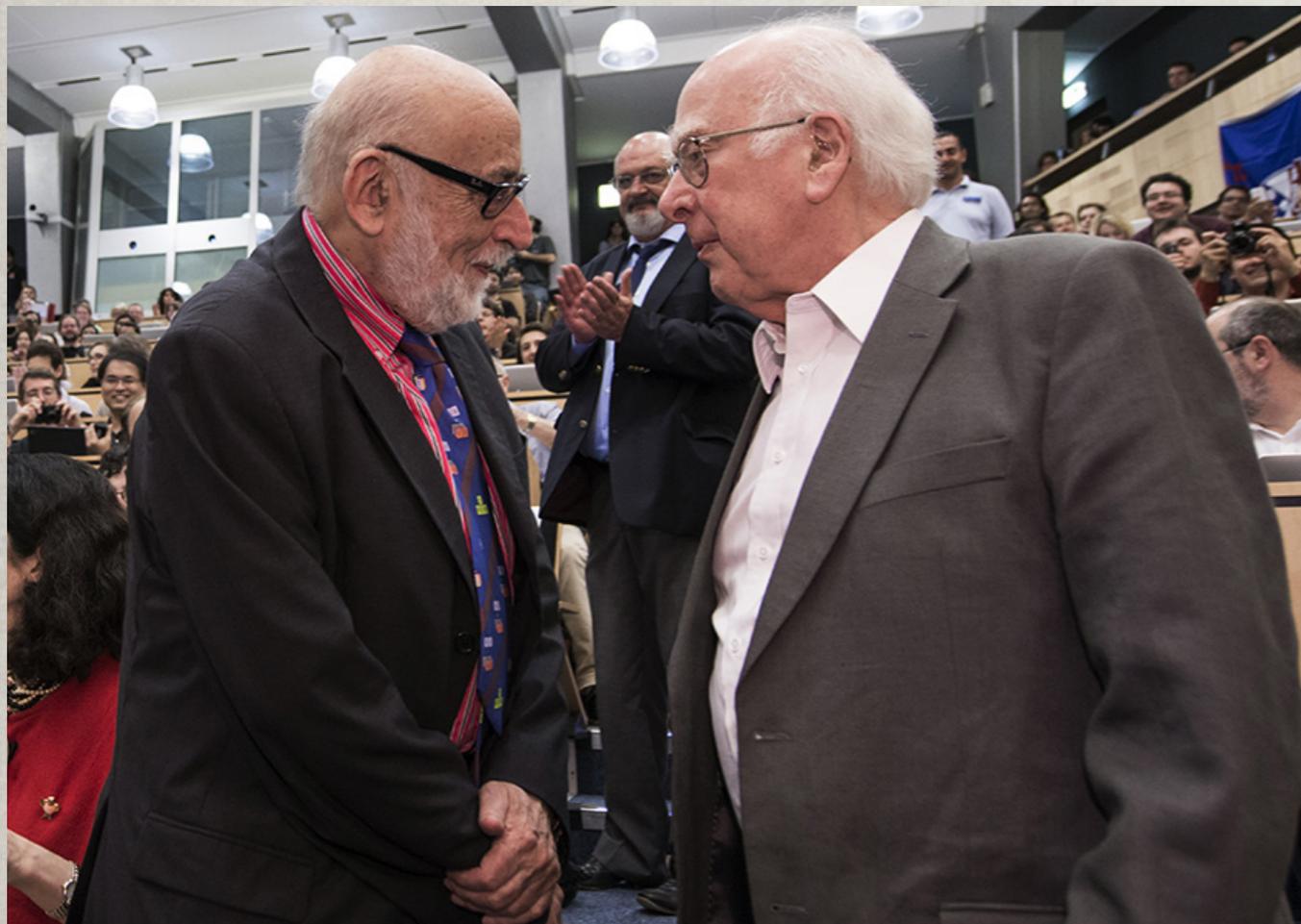


2013 PHYSICS PRIZE



Francois Englert

Peter Higgs



The Higgs Boson (aka “God Particle”)

“...for the theoretical discovery of a mechanism that contributes to our understanding of the origin of mass of subatomic particles, and which recently was confirmed through the discovery of the predicted fundamental particle, by the ATLAS and CMS experiments at CERN's Large Hadron Collider”

2013 PHYSICS PRIZE

Brought to you by the power of the grid...

Francois Englert

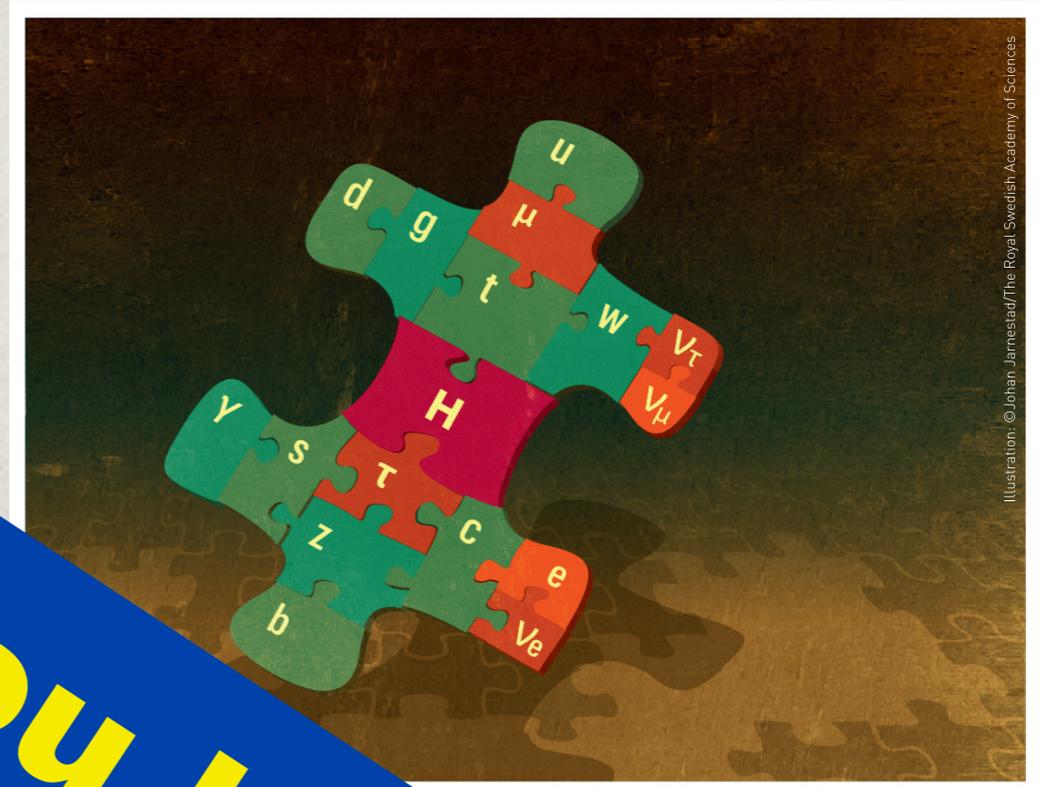
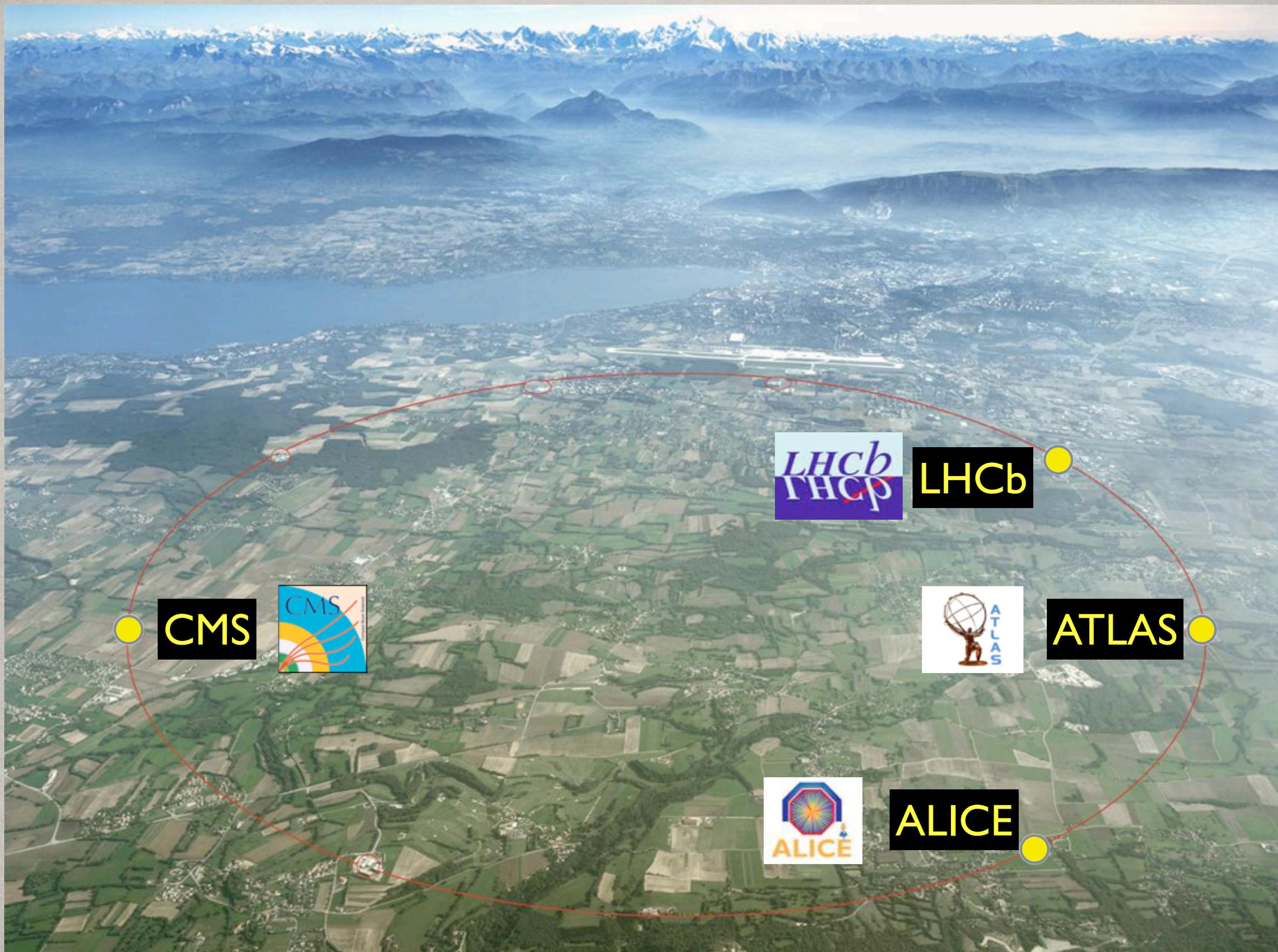


Illustration: ©Johan Jarnestad/The Royal Swedish Academy of Sciences

Higgs Boson (aka

“... mechanism
that con... ending of the
origin of m... ticles, and which
recently was co... h the discovery of
the predicted fundam... particle, by the
ATLAS and CMS experiments at CERN's
Large Hadron Collider”



CMS



LHCb

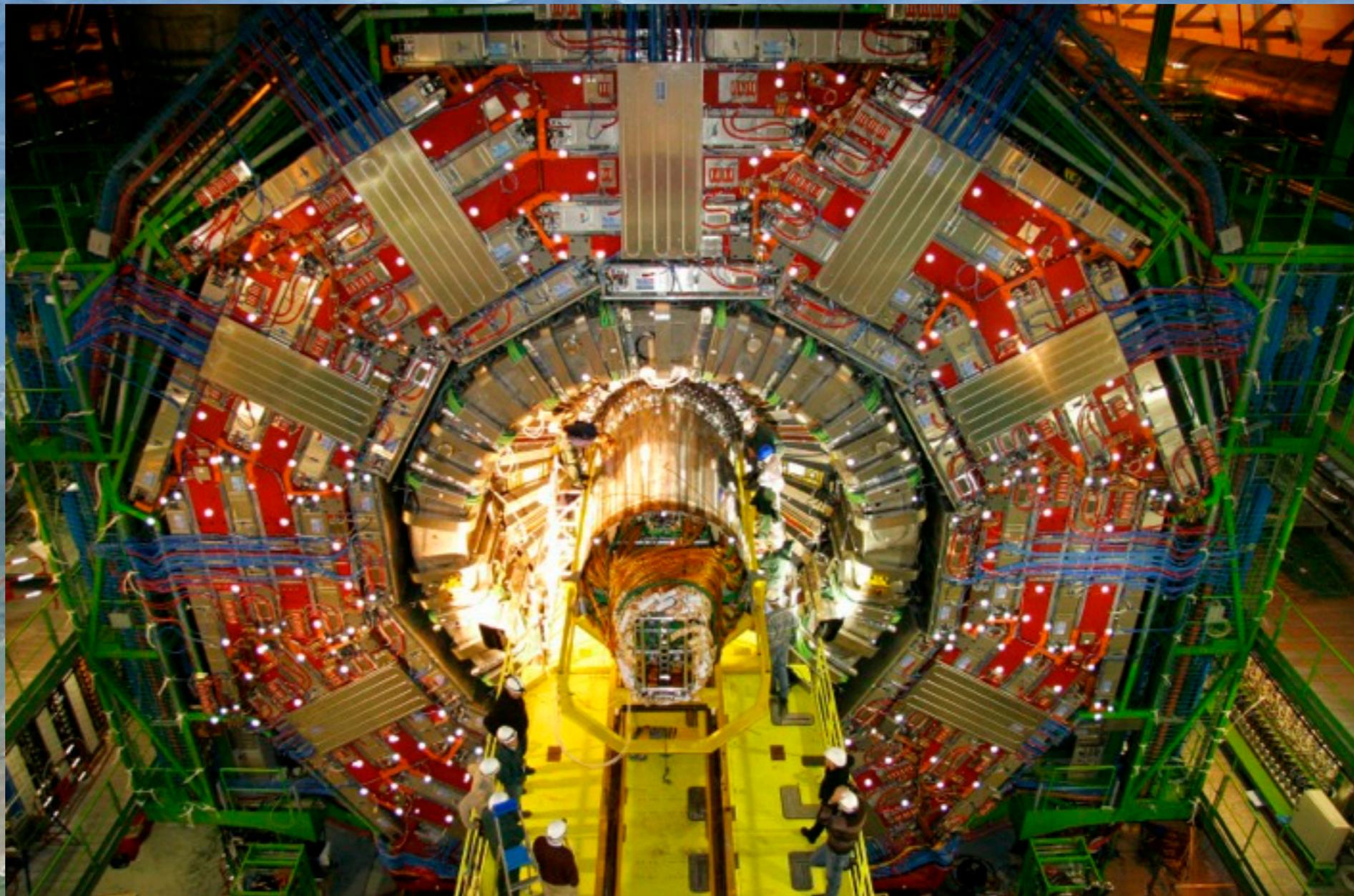


ATLAS

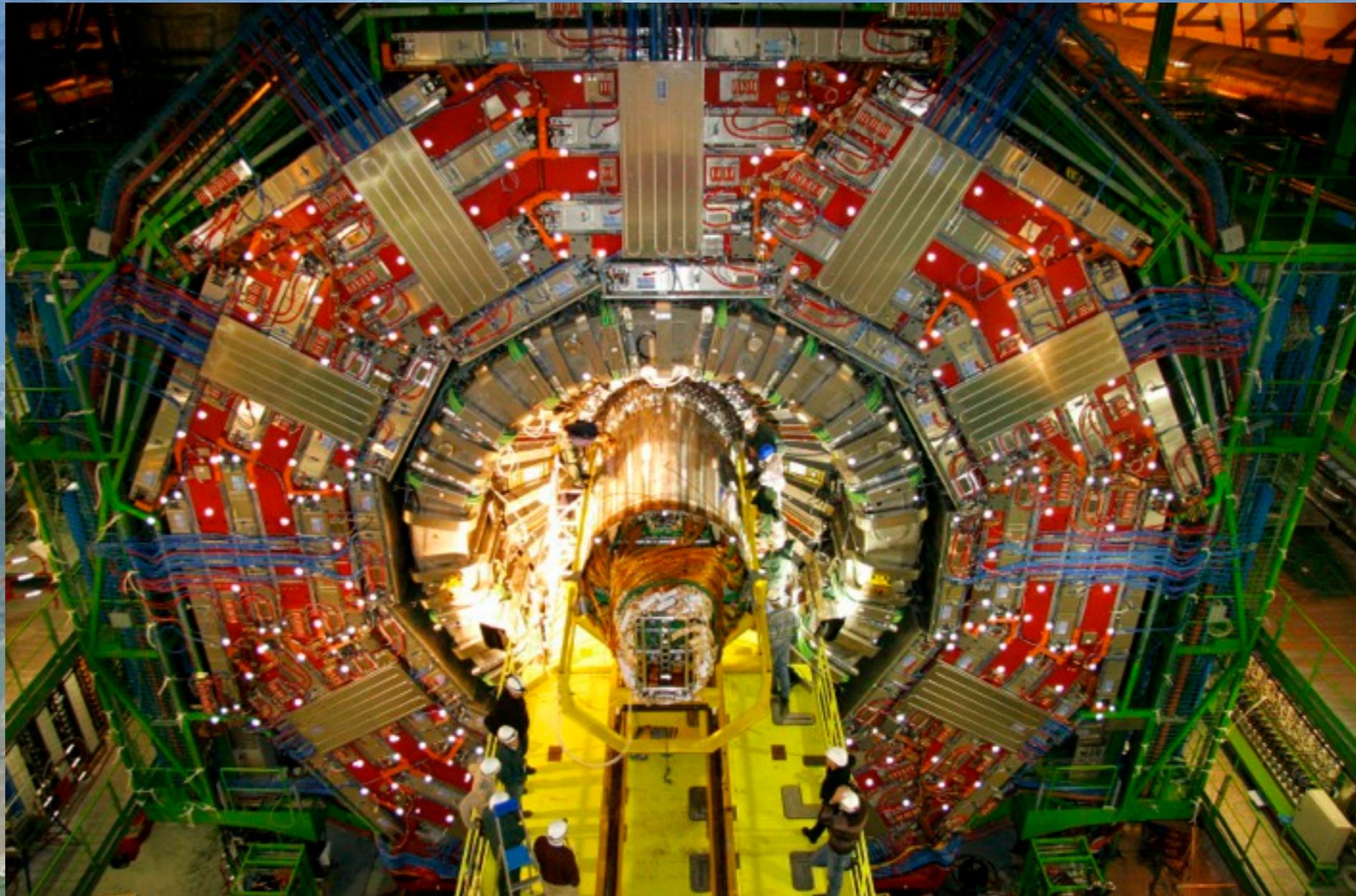


ALICE





The CMS Experiment at the LHC (ND involvement)



The CMS Experiment at the LHC (ND involvement)



OCEANS OF DATA



OCEANS OF DATA

✻ One Higgs boson produced every 3 billion collisions

OCEANS OF DATA

- ✻ One Higgs boson produced every 3 billion collisions
- ✻ Peak rate of 9 Higgs bosons/minute

OCEANS OF DATA

- ✱ One Higgs boson produced every 3 billion collisions
- ✱ Peak rate of 9 Higgs bosons/minute
- ✱ Total number of collisions produced to find Higgs: 690 trillion

OCEANS OF DATA

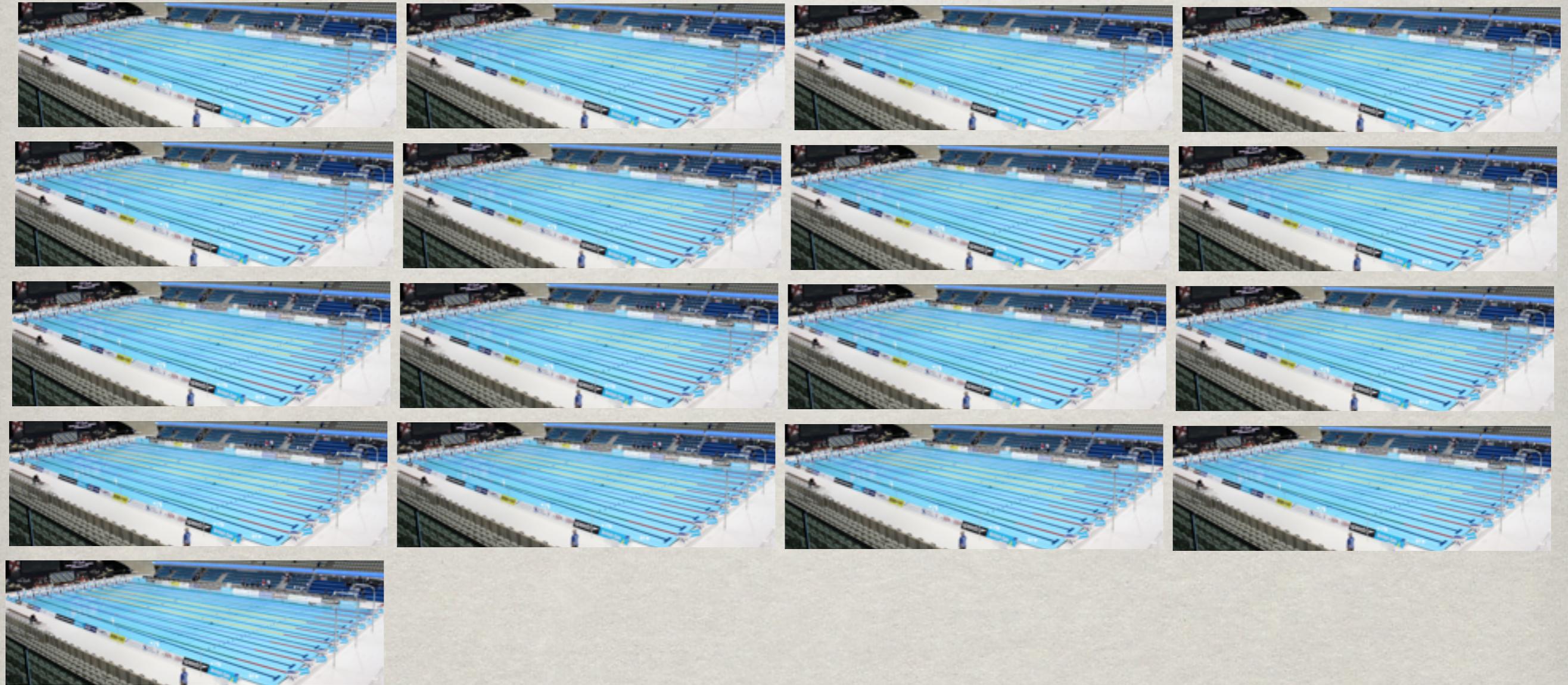
- ✱ One Higgs boson produced every 3 billion collisions
- ✱ Peak rate of 9 Higgs bosons/minute
- ✱ Total number of collisions produced to find Higgs: 690 trillion
- ✱ If each collisions were one grain of sand... 🟡

OCEANS OF DATA

✻ Would fill 17 Olympic-sized swimming pools

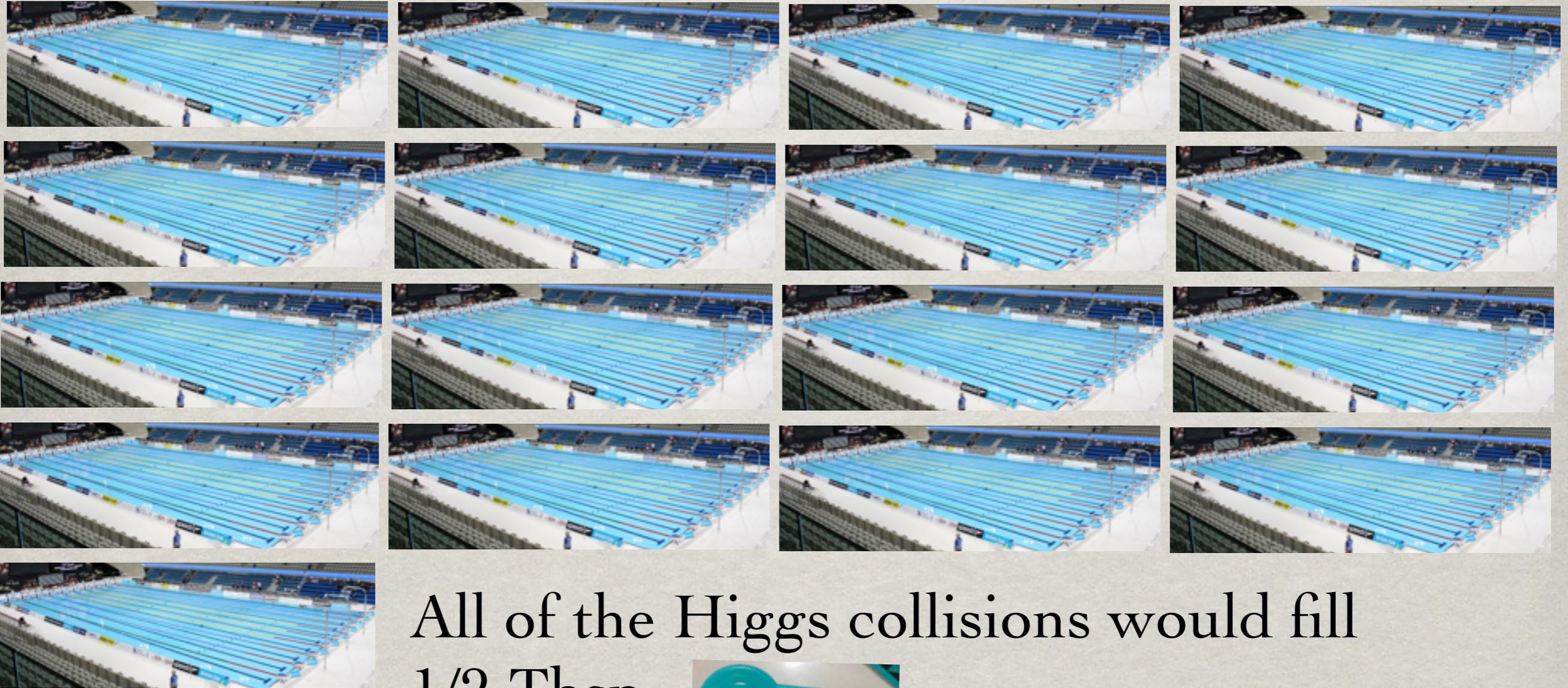
OCEANS OF DATA

☼ Would fill 17 Olympic-sized swimming pools



OCEANS OF DATA

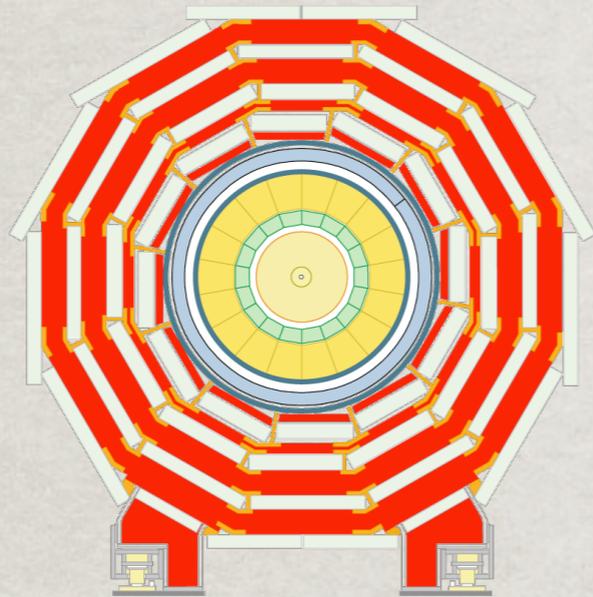
☼ Would fill 17 Olympic-sized swimming pools



All of the Higgs collisions would fill
1/2 Tbsp



HOW FAST DO WE NEED TO GO?



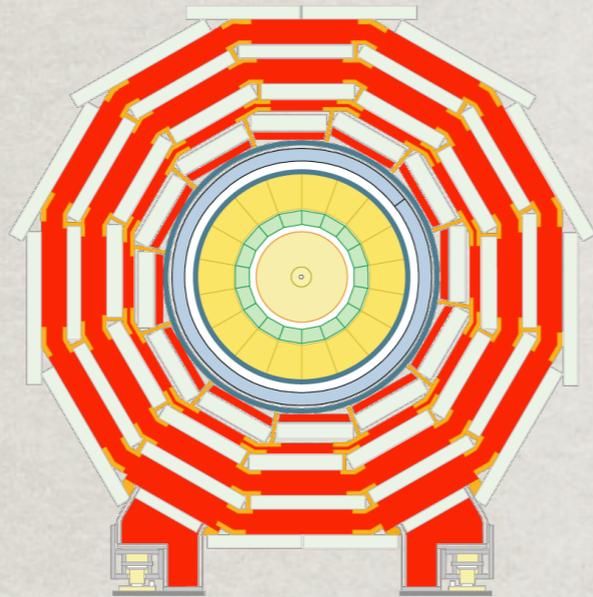
Basic facts:

- ➔ Data from detector: 200 kB/ collision
- ➔ Processing time for analysis: 5 sec (basic)

	Proton Collisions in Detector	Level 1 Trigger	High Level Trigger
Data Rate	16 MHz	60 kHz	300 Hz
Data Collected	50 EB	200 PB	1-2 PB
Processing time	45 Million CPU years!	170 Thousand CPU years	860 CPU years

For 1 year's
worth of data

HOW FAST DO WE NEED TO GO?



Basic facts:

- ➔ Data from detector: 200 kB/ collision
- ➔ Processing time for analysis: 5 sec (basic)

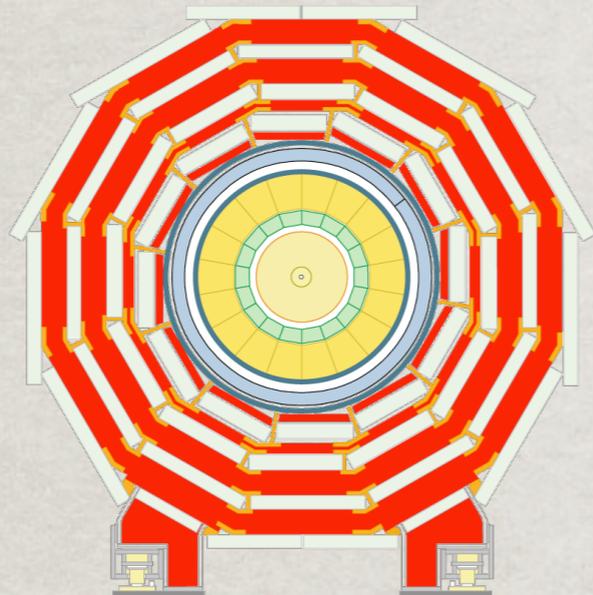
	Proton Collisions in Detector	Level 1 Trigger	High Level Trigger
Data Rate	16 MHz	60 kHz	300 Hz
Data Collected	50 EB	200 PB	1-2 PB
Processing time	45 Million CPU years!	170 Thousand CPU years	860 CPU years

For 1 year's worth of data

1 MILLION TB!

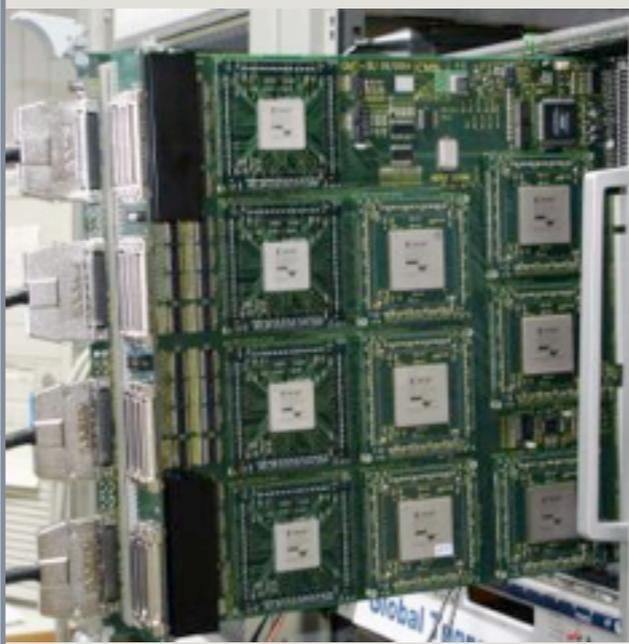


HOW FAST DO WE NEED TO GO?



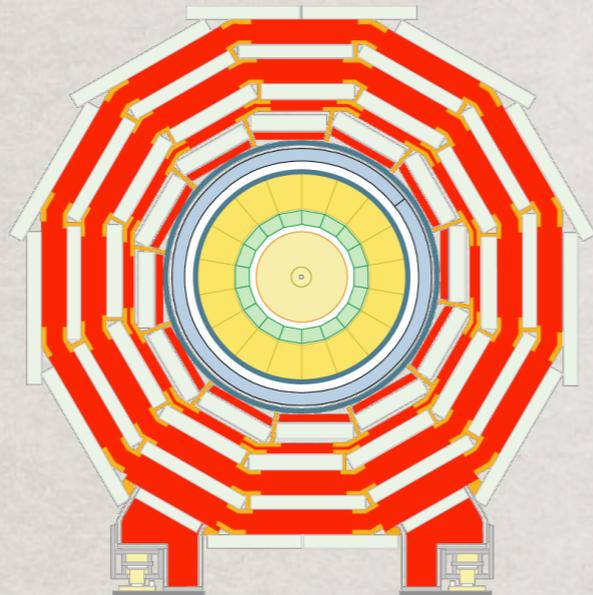
Basic facts:

- ➔ Data from detector: 200 kB/ collision
- ➔ Processing time for analysis: 5 sec (basic)

	Proton Collisions in Detector	Level 1 Trigger
Data Rate	16 MHz	
Data Collected	50 EB	
Processing time	45 Million CPU years!	

For 1 year's
worth of data

HOW FAST DO WE NEED TO GO?



Basic facts:

- ➔ Data from detector: 200 kB/ collision
- ➔ Processing time for analysis: 5 sec (basic)

FPGA Chips do
very simple
analysis
~ μ s to analyze
data

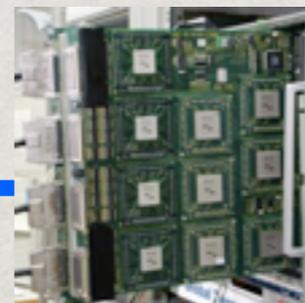
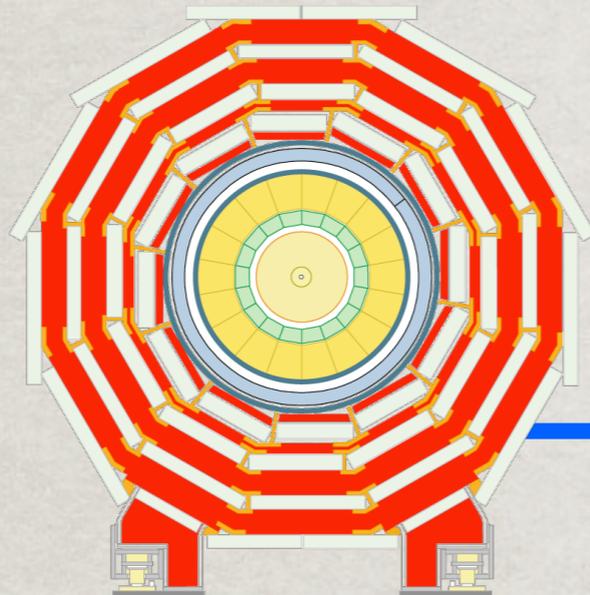
	Proton Collisions in Detector	Level 1 Trigger
Data Rate	16 MHz	
Data Collected	50 EB	
Processing time	45 Million CPU years!	

For 1 year's
worth of data

HOW FAST DO WE NEED TO GO?

Basic facts:

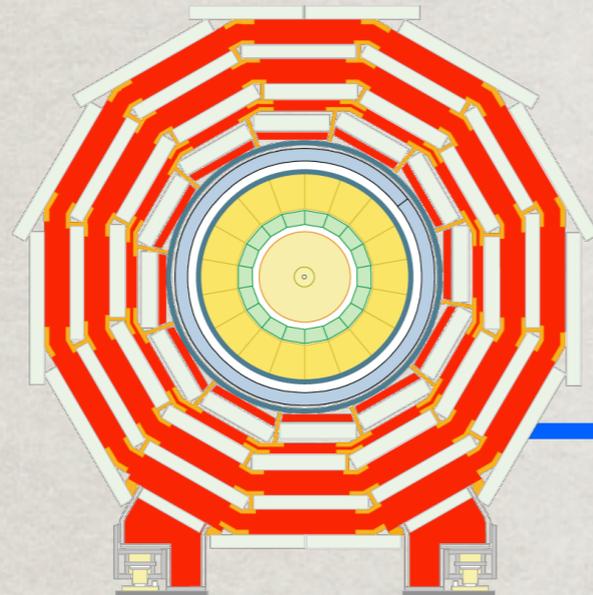
- ➔ Data from detector: 200 kB/ collision
- ➔ Processing time for analysis: 5 sec (basic)



	Proton Collisions in Detector	Level 1 Trigger
Data Rate	16 MHz	60 kHz
Data Collected	50 EB	200 PB
Processing time	45 Million CPU years!	170 Thousand CPU years

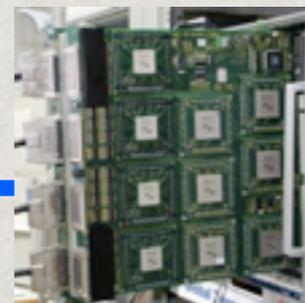
For 1 year's
worth of data

HOW FAST DO WE NEED TO GO?



Basic facts:

- ➔ Data from detector: 200 kB/ collision
- ➔ Processing time for analysis: 5 sec (basic)

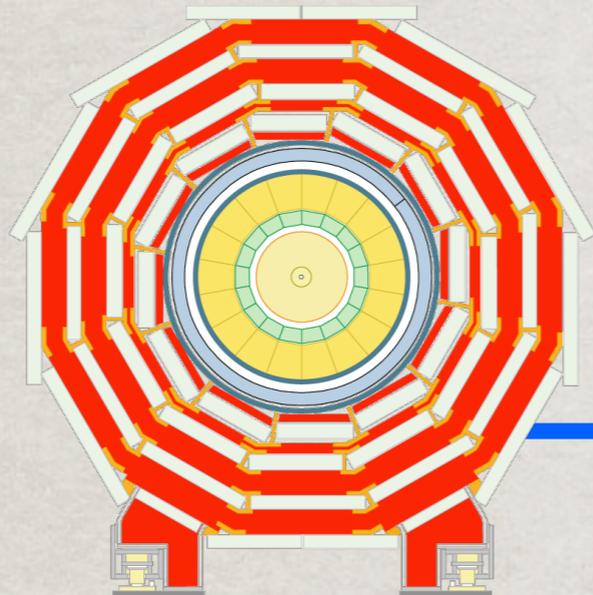


	Proton Collisions in Detector	Level 1 Trigger	High Level Trigger
Data Rate	16 MHz	60 kHz	Computer farm with ~ 3000 CPUs High speed network/ switch
Data Collected	50 EB	200 PB	
Processing time	45 Million CPU years!	170 Thousand CPU years	



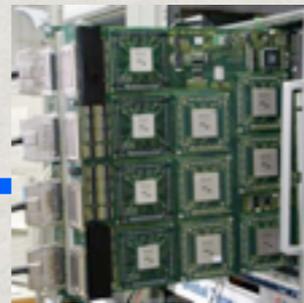
For 1 year's worth of data

HOW FAST DO WE NEED TO GO?



Basic facts:

- ➔ Data from detector: 200 kB/ collision
- ➔ Processing time for analysis: 5 sec (basic)



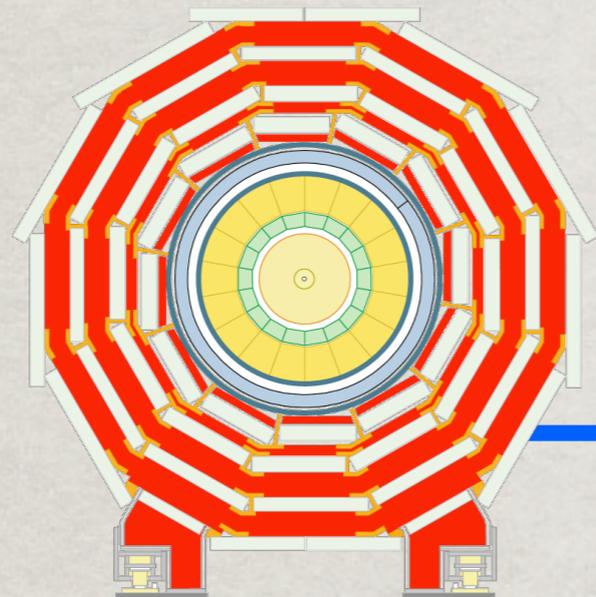
Simplified analysis code
 ~ ms to analyze data

	Proton Collisions in Detector	Level 1 Trigger	High Level Trigger
Data Rate	16 MHz	60 kHz	Computer farm with ~ 3000 CPUs High speed network/ switch
Data Collected	50 EB	200 PB	
Processing time	45 Million CPU years!	170 Thousand CPU years	



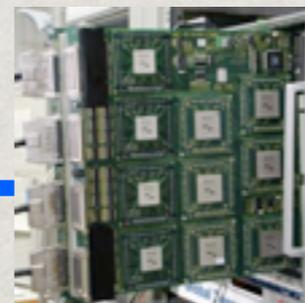
For 1 year's worth of data

HOW FAST DO WE NEED TO GO?



Basic facts:

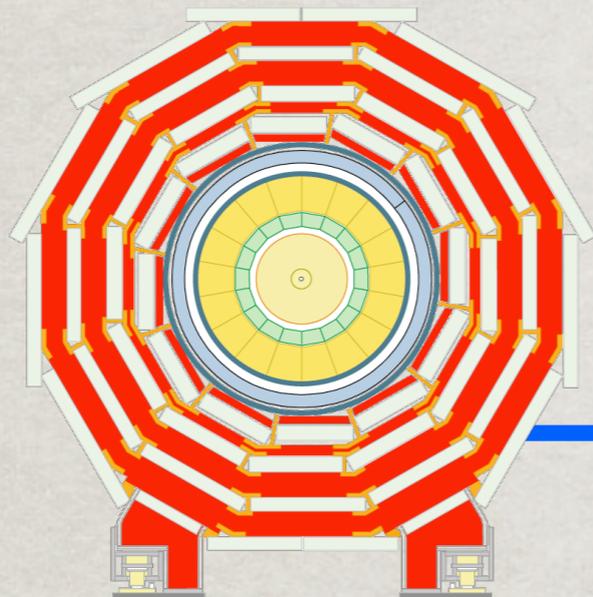
- ➔ Data from detector: 200 kB/ collision
- ➔ Processing time for analysis: 5 sec (basic)



	Proton Collisions in Detector	Level 1 Trigger	High Level Trigger
Data Rate	16 MHz	60 kHz	300 Hz
Data Collected	50 EB	200 PB	1-2 PB
Processing time	45 Million CPU years!	170 Thousand CPU years	860 CPU years

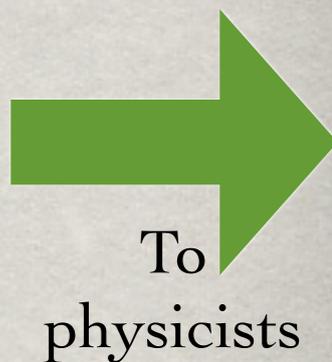
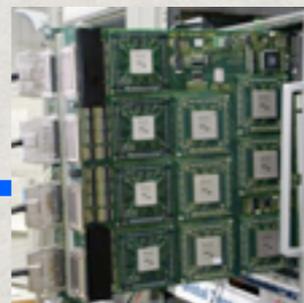
For 1 year's worth of data

HOW FAST DO WE NEED TO GO?



Basic facts:

- ➔ Data from detector: 200 kB/ collision
- ➔ Processing time for analysis: 5 sec (basic)



	Proton Collisions in Detector	Level 1 Trigger	High Level Trigger
Data Rate	16 MHz	60 kHz	300 Hz
Data Collected	50 EB	200 PB	1-2 PB
Processing time	45 Million CPU years!	170 Thousand CPU years	860 CPU years

For 1 year's
worth of data

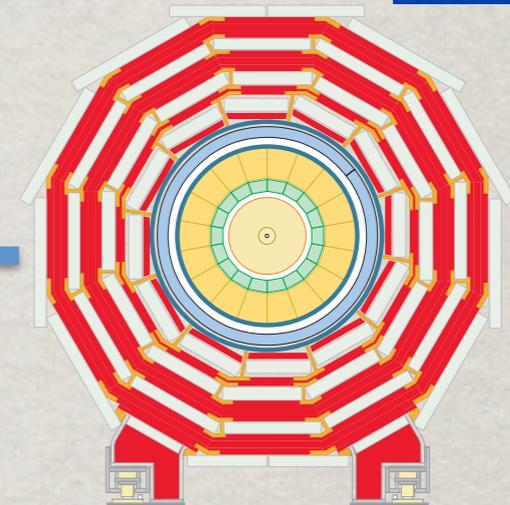
PROCESSING CMS DATA

~ 5-10 PB/year

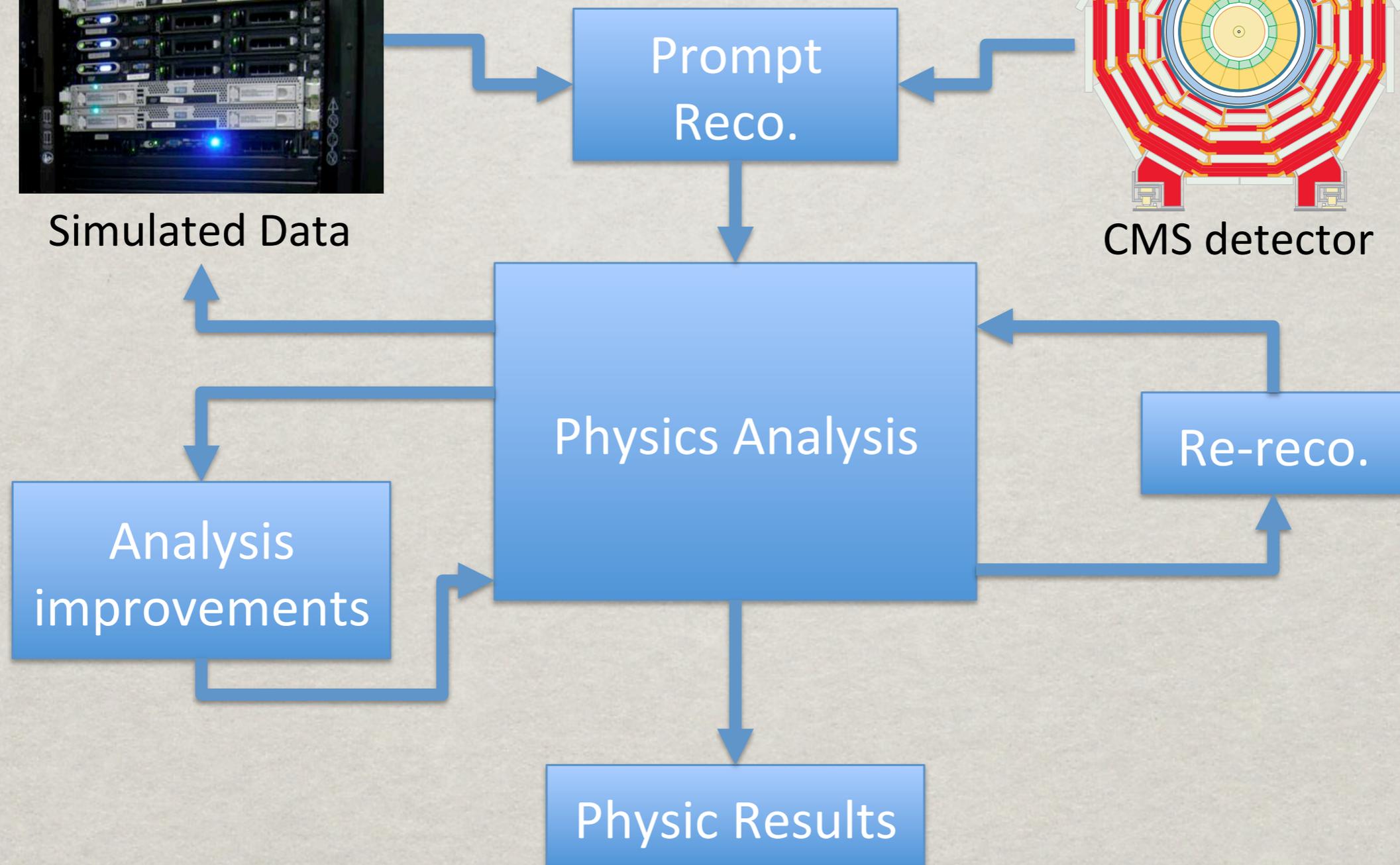
~ 2 PB/year



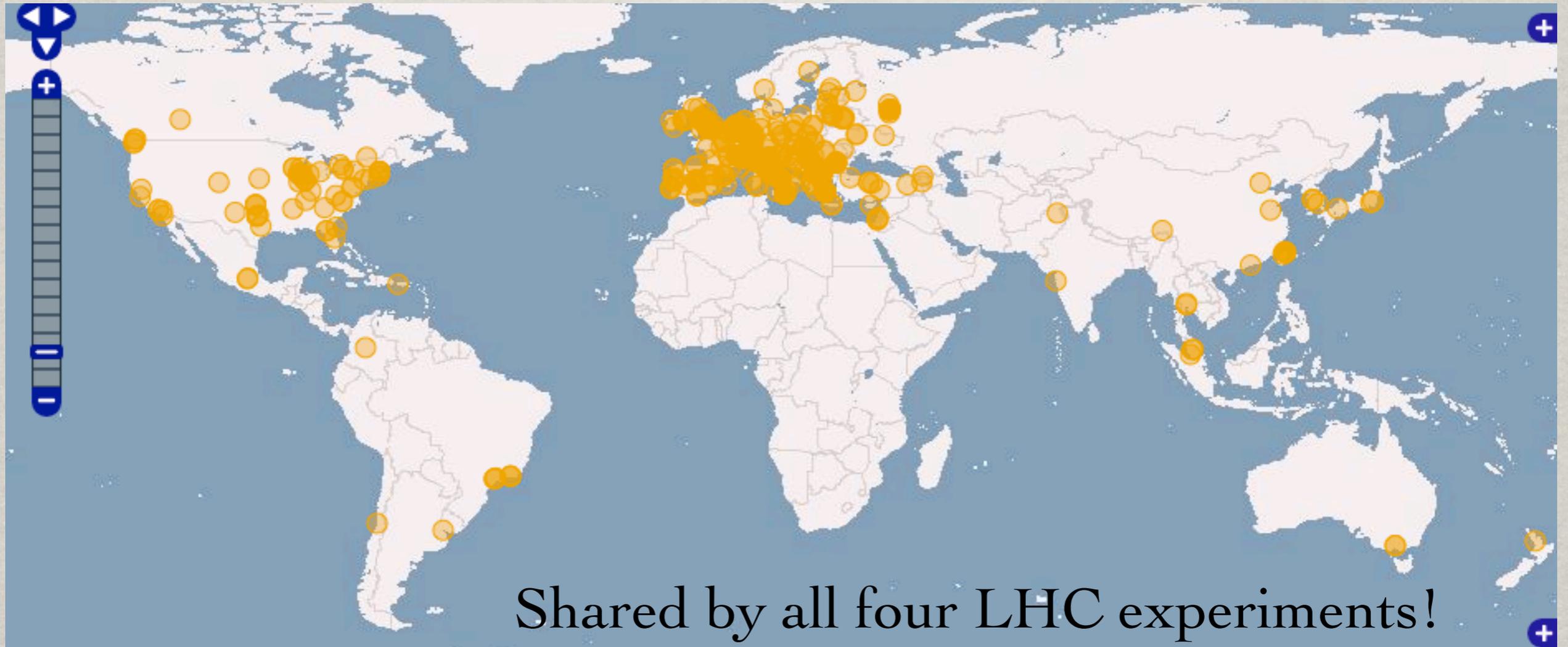
Simulated Data



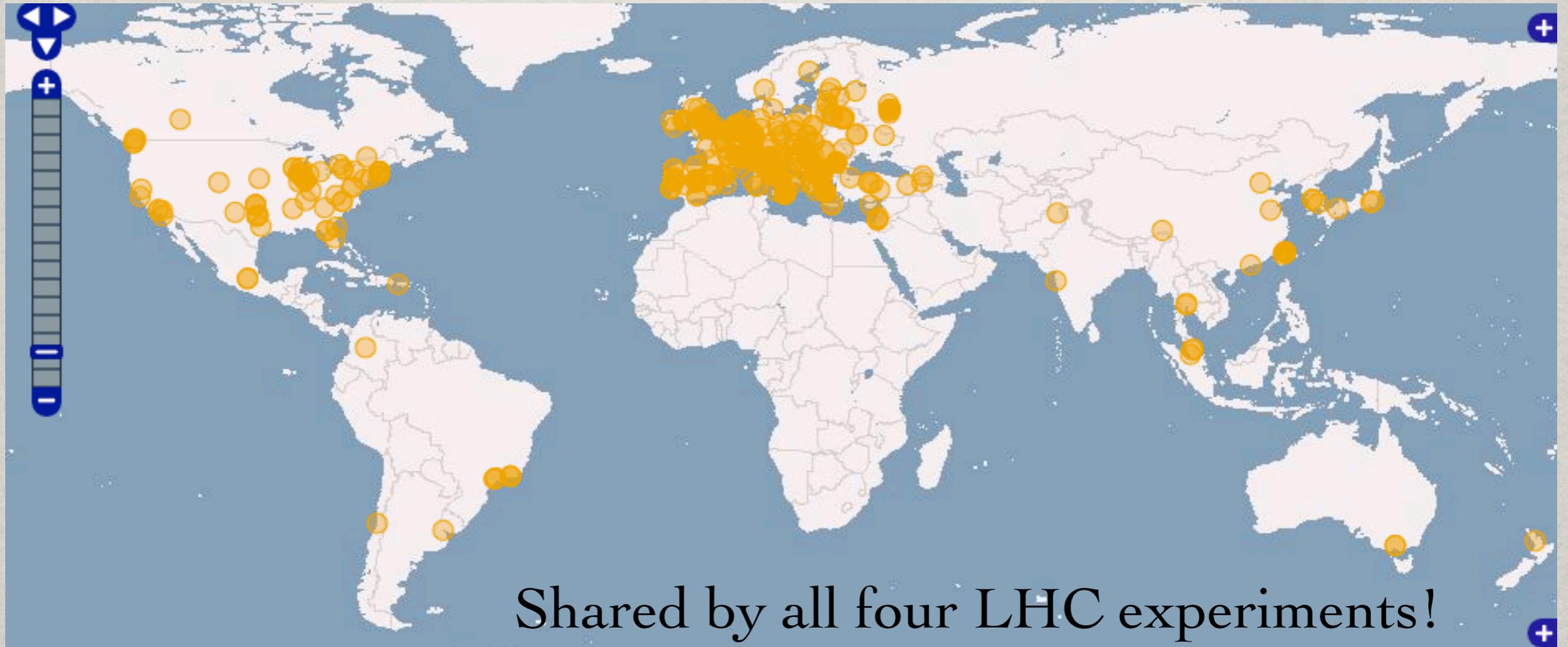
CMS detector



WORLDWIDE LHC COMPUTING GRID



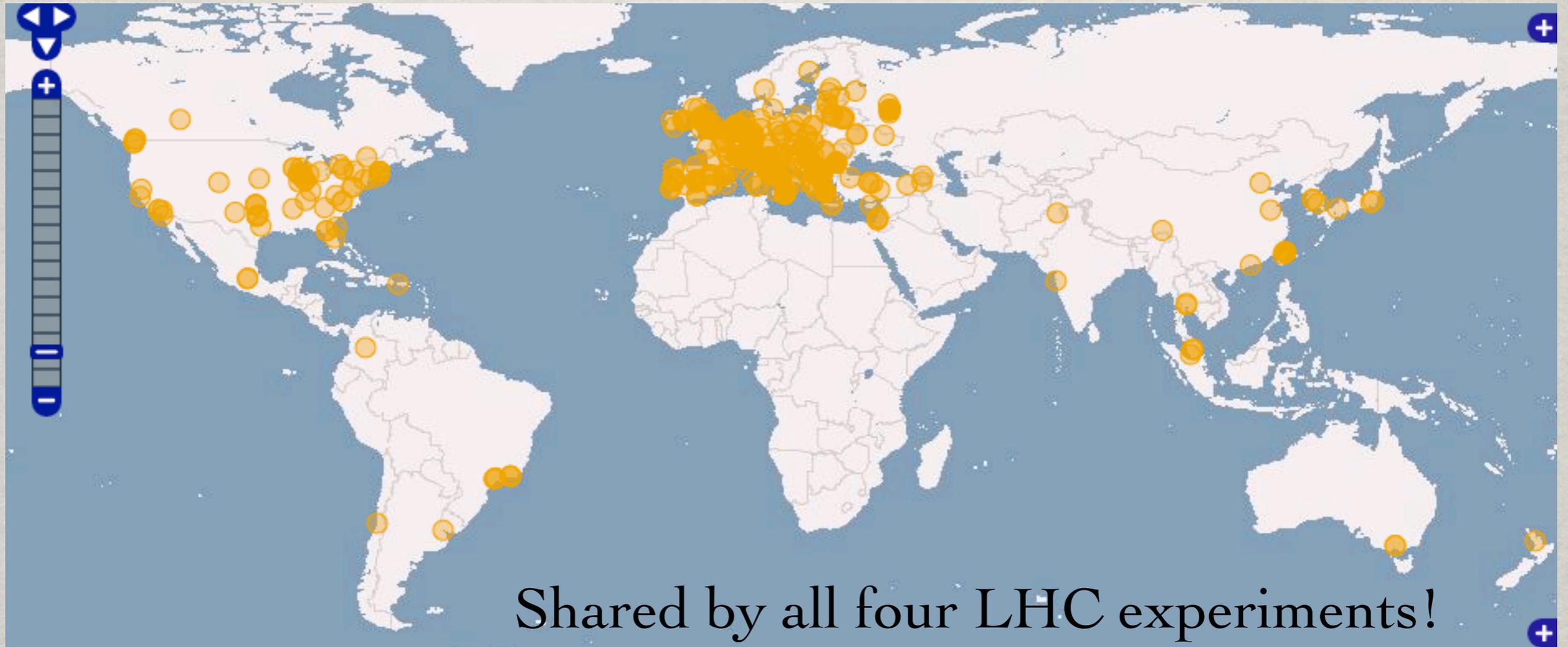
WORLDWIDE LHC COMPUTING GRID



Shared by all four LHC experiments!

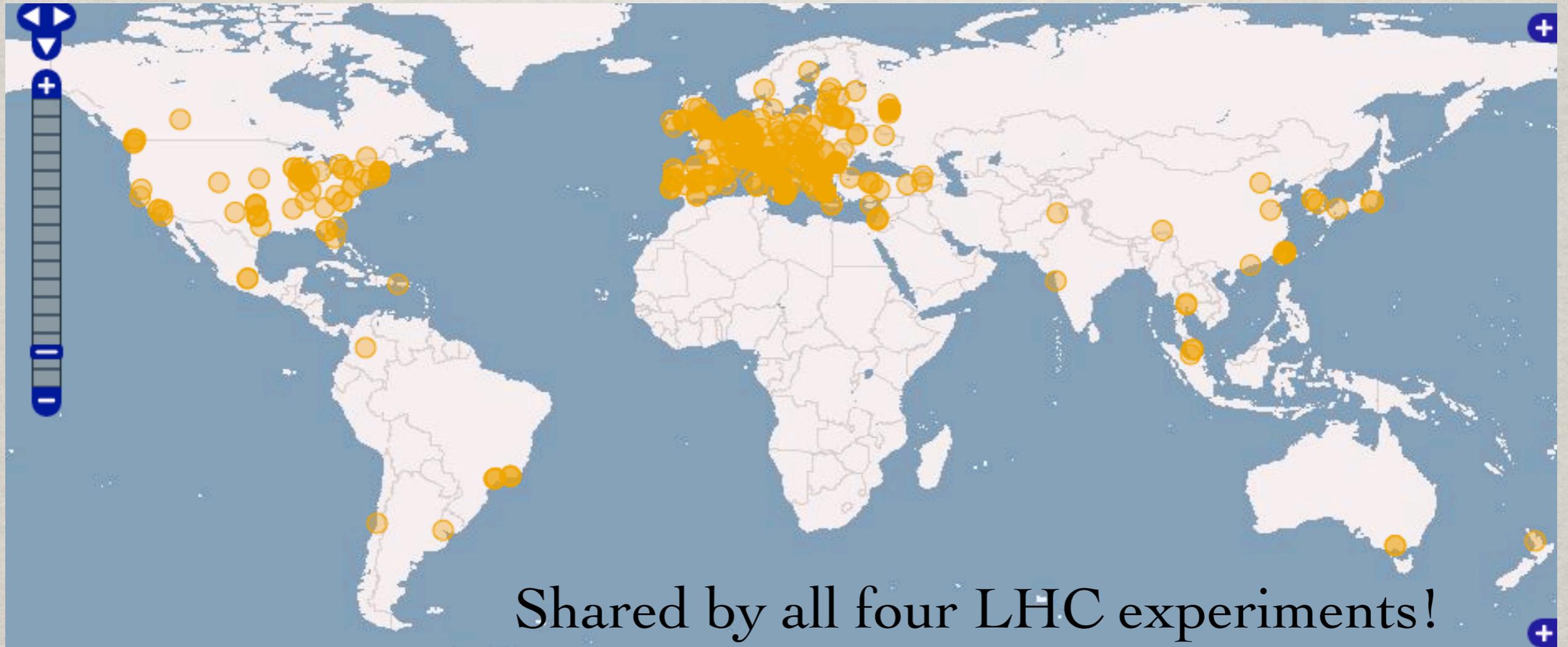
- ✻ Over 160 sites around world (including OSG sites in US)

WORLDWIDE LHC COMPUTING GRID



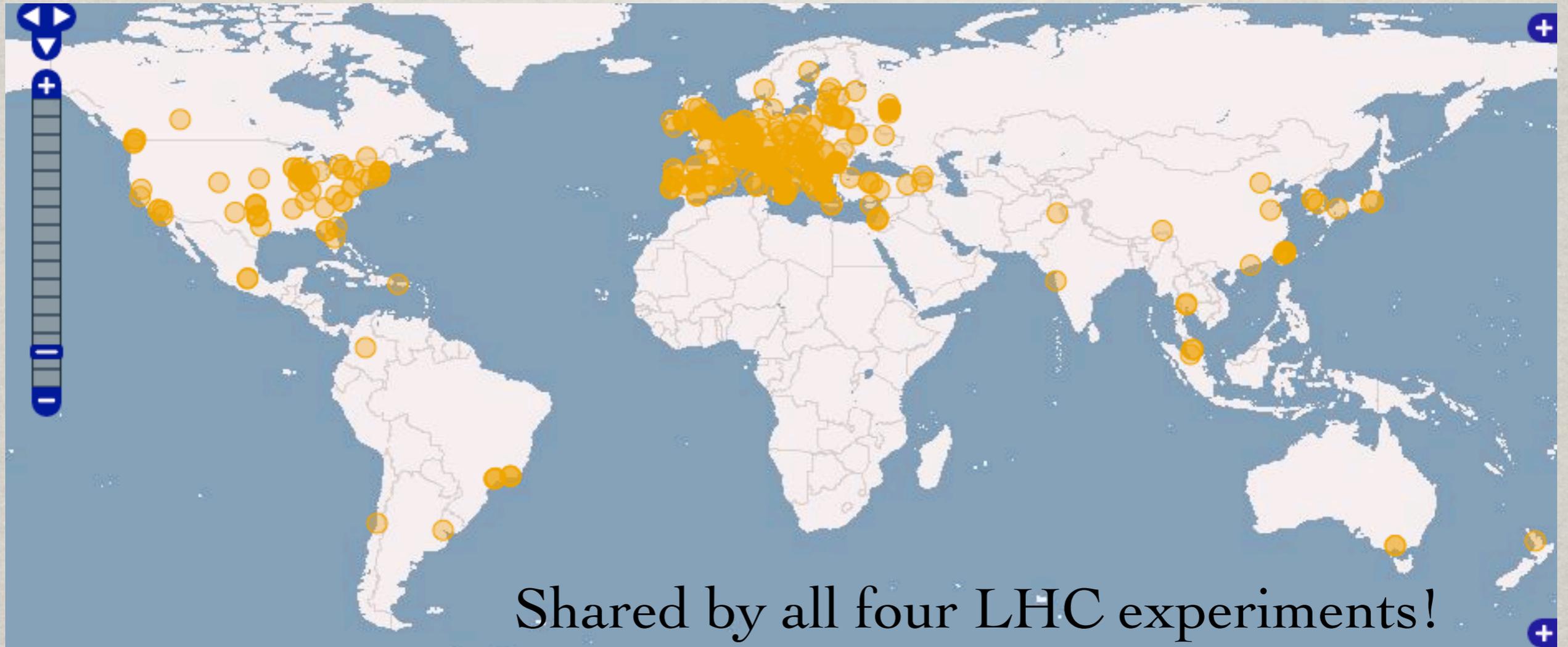
- ✻ Over 160 sites around world (including OSG sites in US)
- ✻ > 200k CPU cores available

WORLDWIDE LHC COMPUTING GRID



- ✿ Over 160 sites around world (including OSG sites in US)
- ✿ > 200k CPU cores available
- ✿ Has gone as high as ~ 1 million jobs submitted in a single day

WORLDWIDE LHC COMPUTING GRID



Shared by all four LHC experiments!

- ✿ Over 160 sites around world (including OSG sites in US)
- ✿ > 200k CPU cores available
- ✿ Has gone as high as ~ 1 million jobs submitted in a single day
- ✿ > 300 PB of total storage available

ORGANIZATION



ORGANIZATION

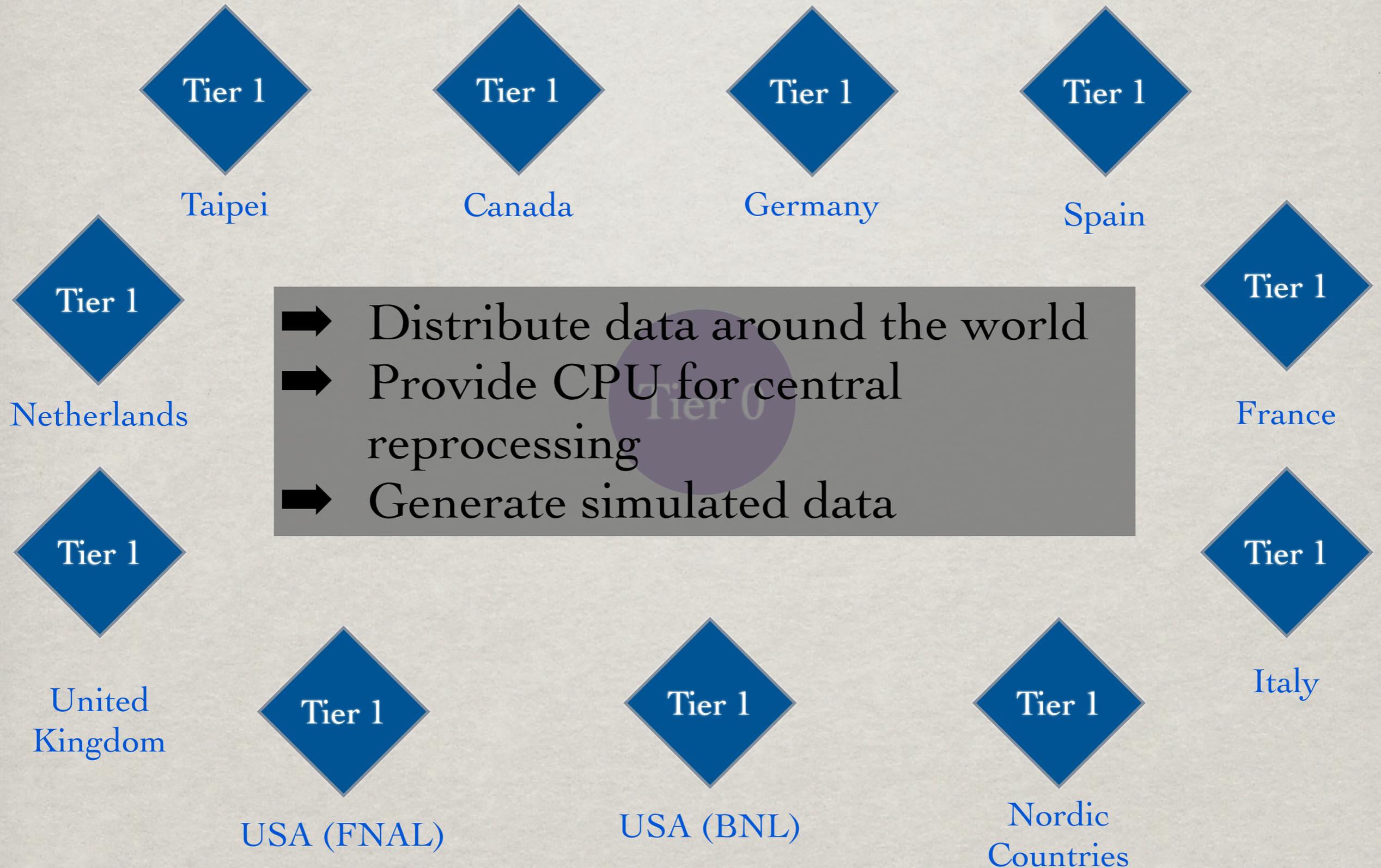


- ➔ All LHC data passes through T0 for initial processing
- ➔ Provides less than 20% of total CPU resources for LHC experiments
- ➔ Basic data processing common to all analyses

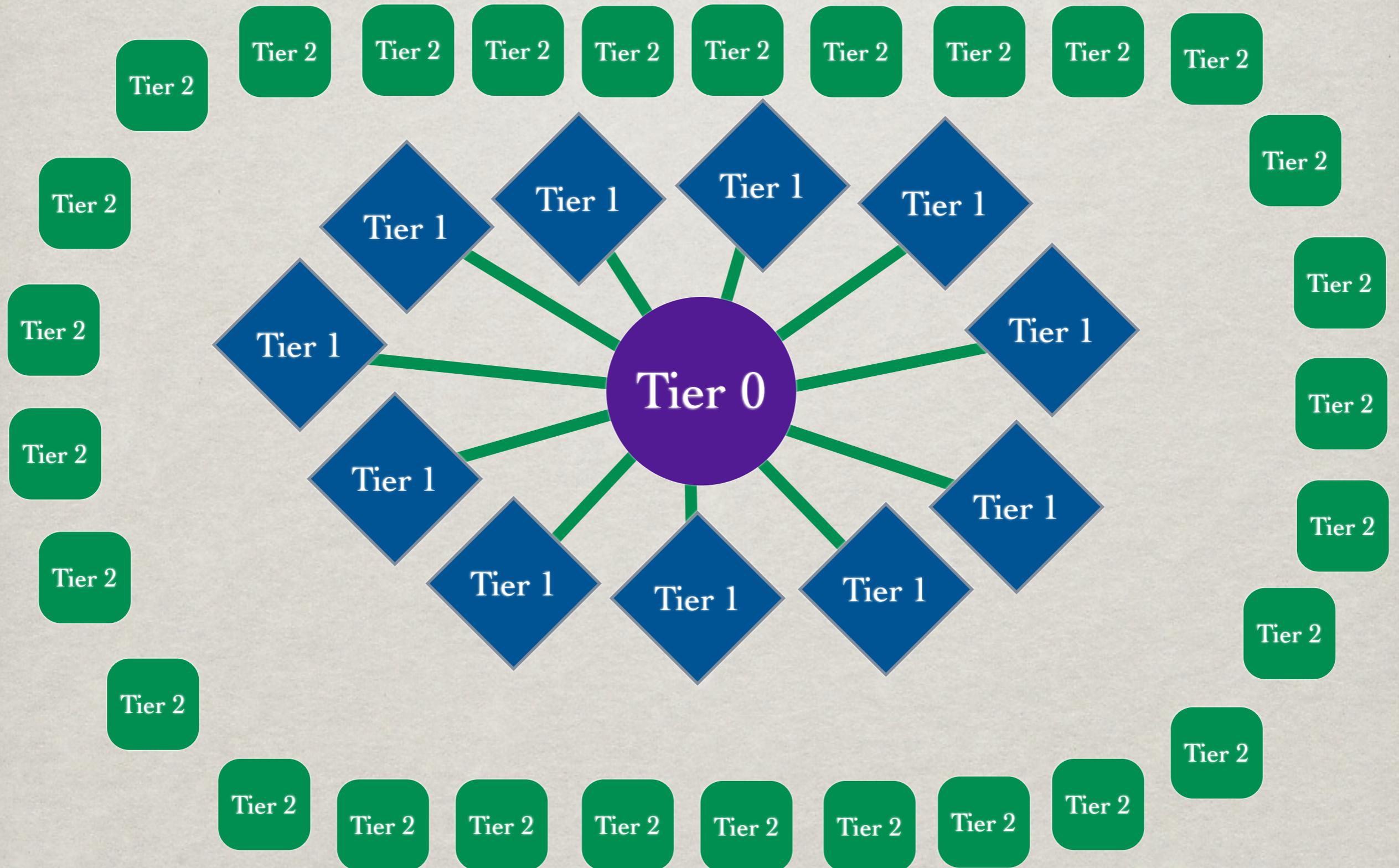
ORGANIZATION



ORGANIZATION



ORGANIZATION



ORGANIZATION



➔ Over 140 T2 sites throughout the world

➔ Average site has ~800 CPU's and 300 TB storage

▶ Some have much more: > 1 PB storage!

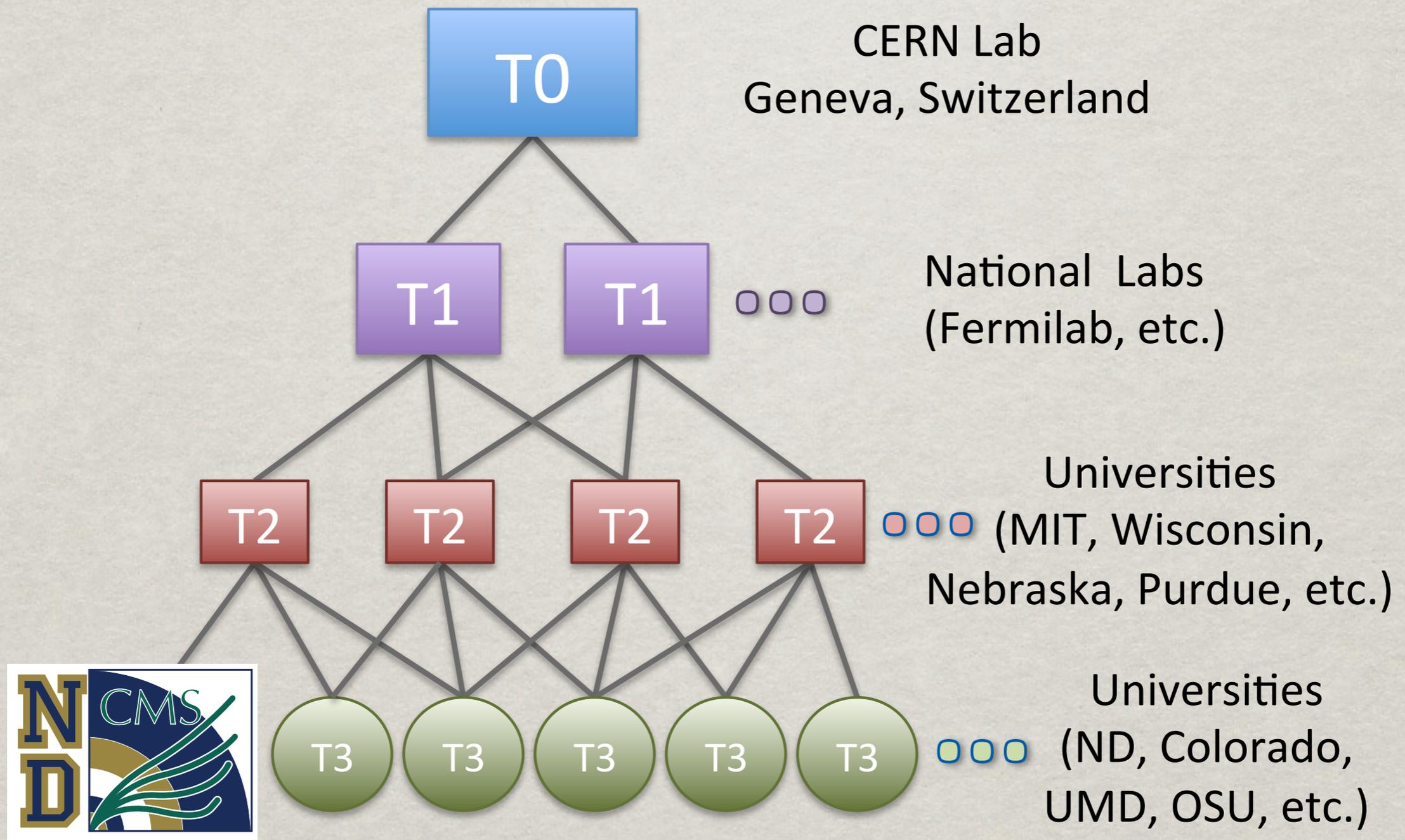
➔ Provide CPU and storage for analysis of data, plus some simulation

➔ This is where “average user” runs analysis

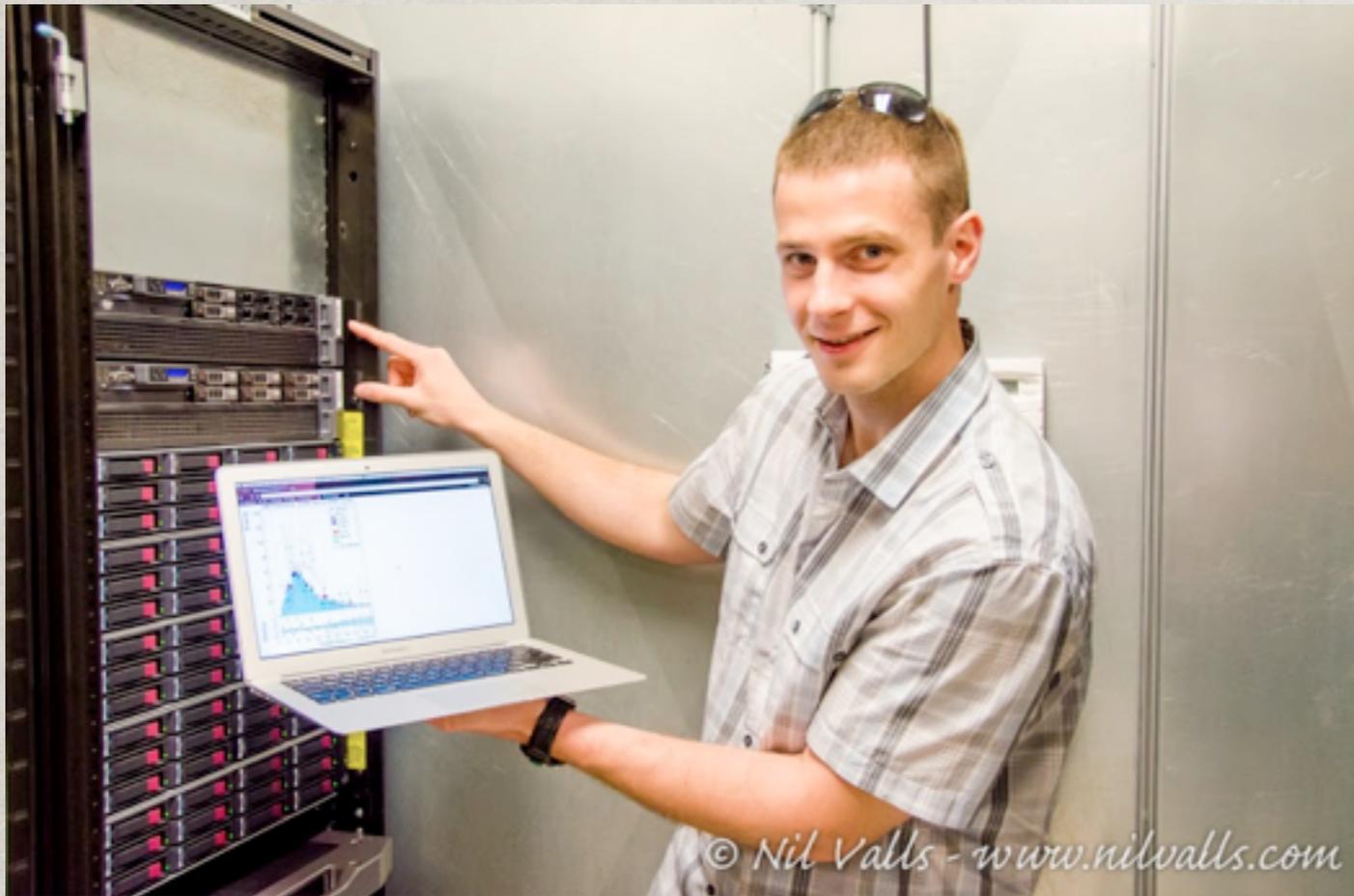
➔ Connected via regional internet links



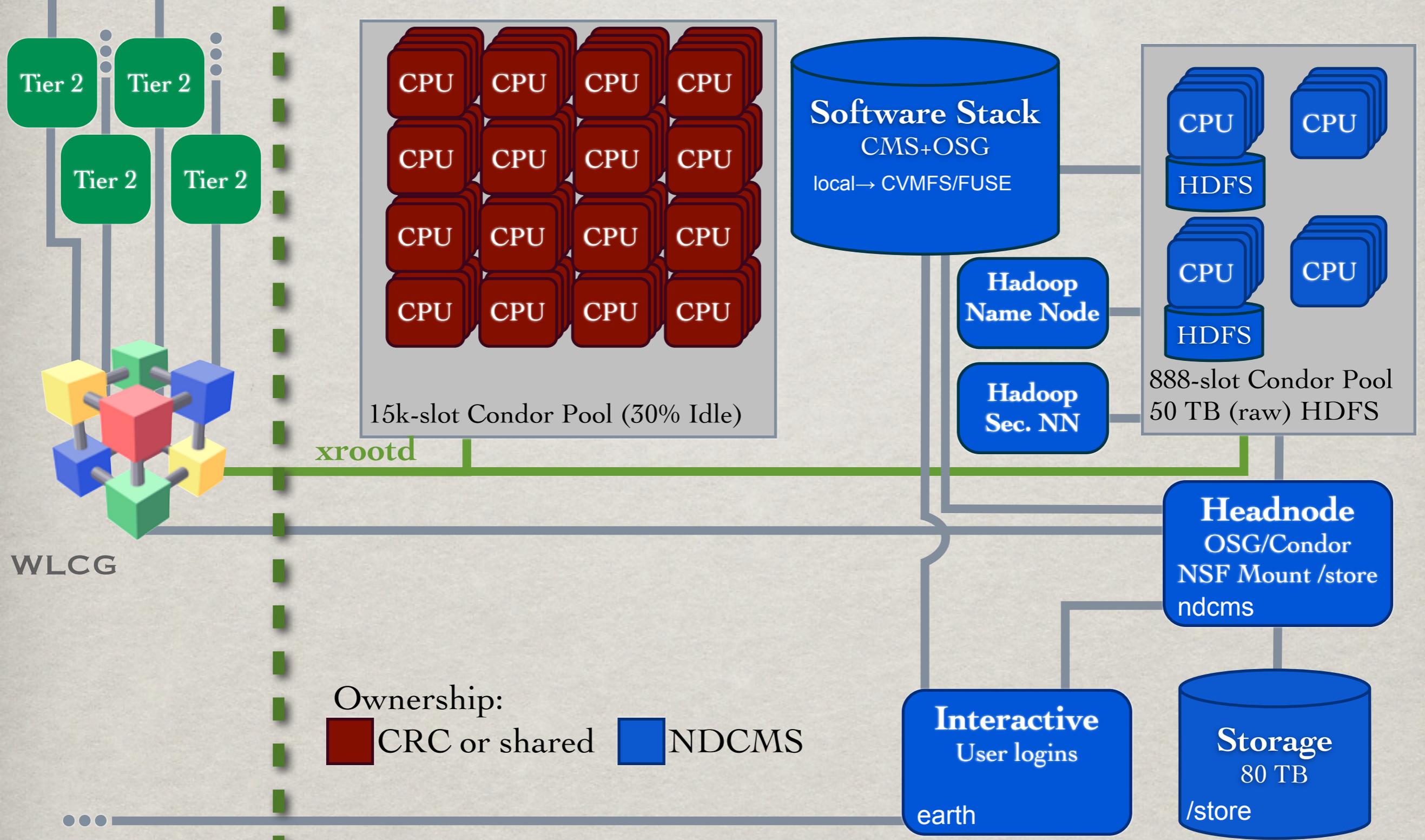
WHERE ND FITS IN



ND T3



BASIC IDEA IN PICTURES



WHAT DO WE USE IT FOR?

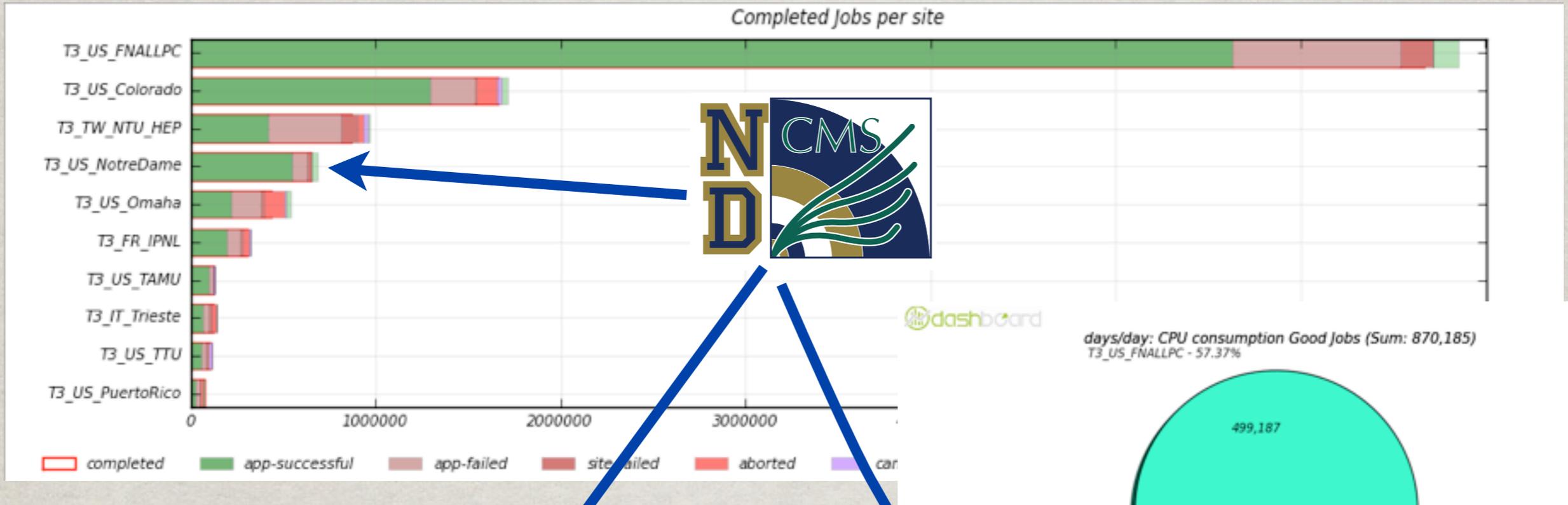
- ✱ Short answer: physics analysis
- ✱ More details
 - ✱ Start with large general purpose dataset (simulations + real data): ~400TB
 - ✱ Reduce by keeping only essential information for an analysis: ~20 TB
 - ✱ Select most interesting subset of reduced data: ~100 GB
 - ✱ Make plots
- ✱ First step done by reading data stored at T2 over network (xrootd)
 - ✱ Goes against trend of “moving code to data” (I think)
 - ✱ Solves two problems
 - ✱ Lack of storage (and need to manage large amounts of data needed only briefly)
 - ✱ Transfer time (overlaps with processing time)
- ✱ Rest done with local data access

[Talk by Brian Bockelman \(last year\)](#)



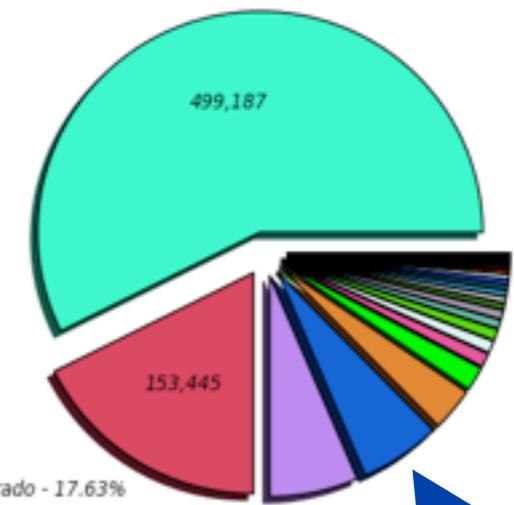
USAGE STATISTICS

Completed Jobs per site



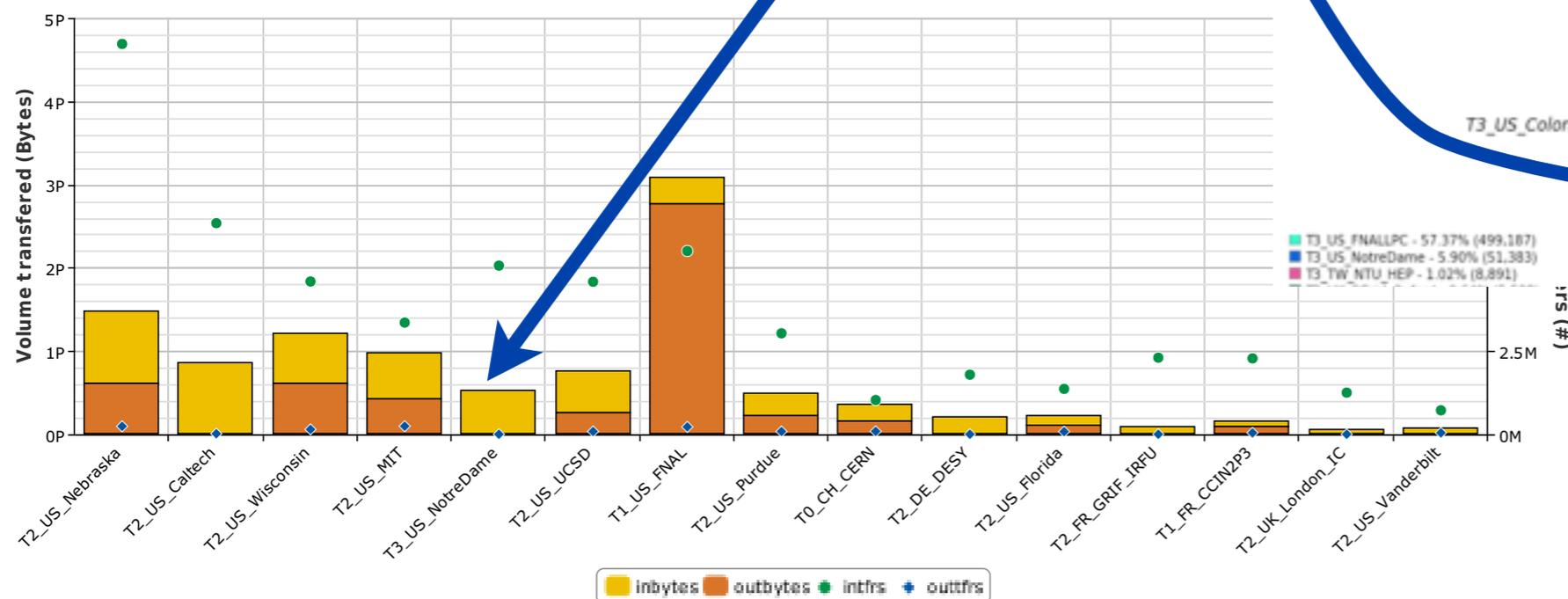
dashboard

days/day: CPU consumption Good Jobs (Sum: 870,185)
T3_US_FNALLPC - 57.37%



TRAFFIC STATISTICS PER SITE

2012-10-10 00:00 - 2013-10-10 00:00 UTC

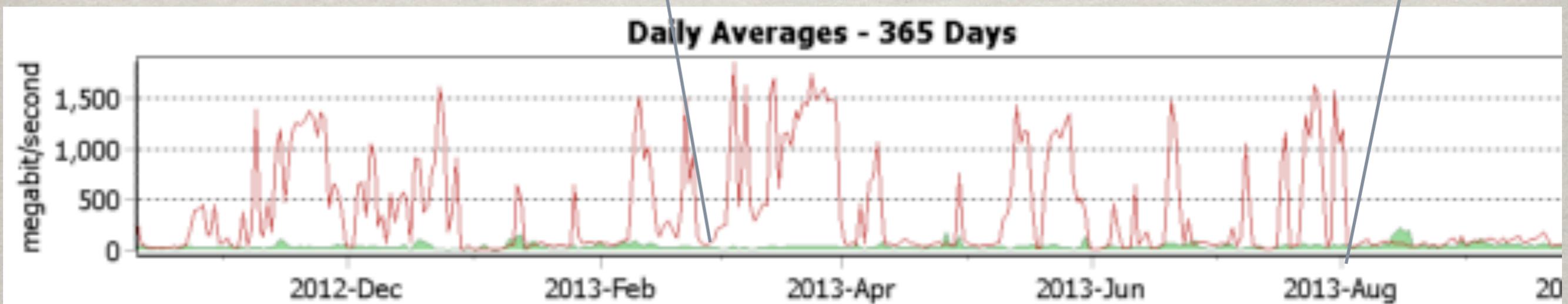
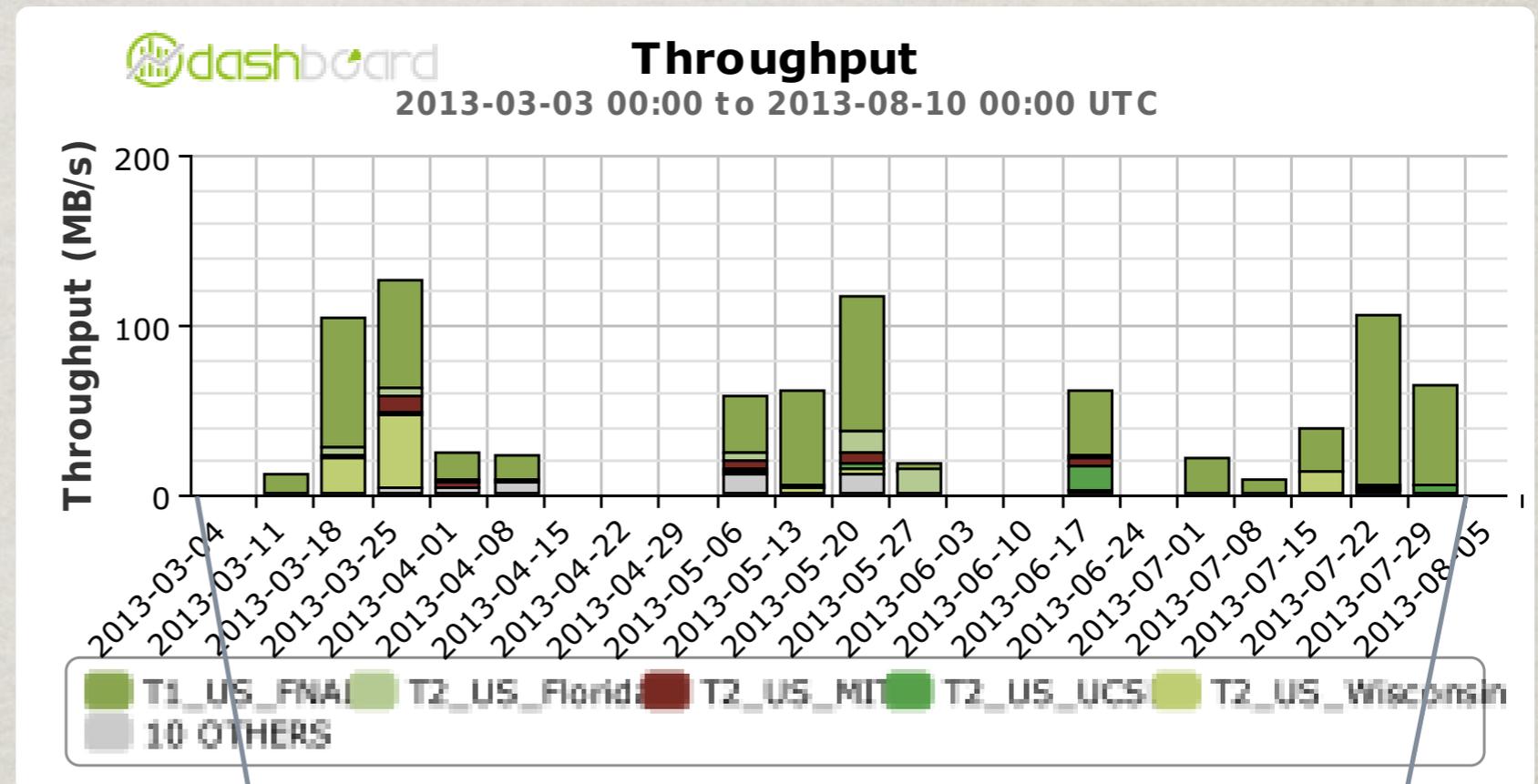


T3_US_FNALLPC - 57.37% (499,187)
T3_US_Colorado - 17.63% (153,445)
T3_US_NotreDame - 5.90% (51,383)
T3_TW_NTU_HEP - 1.02% (8,891)
T3_US_Omaha - 6.44% (56,053)
T3_FR_IPNL - 3.19% (27,759)
T3_IT_Trieste - 1.75% (15,214)
T3_US_TAMU - 0.88% (7,641)
T3_US_TTU - 0.82% (7,095)



BANDWIDTH USAGE

Haven't saturated campus bandwidth yet!



SCALING FROM HERE

- ✱ Main bottleneck (right now): CPU
 - ✱ Most intensive part of the processing takes ≥ 3 weeks using full T3 resources
 - ✱ Would like it to be faster
- ✱ Plan: top into opportunistic resources on campus
 - ✱ CRC pool
 - ✱ Potentially 15k cores available
 - ✱ Only $\sim 30\%$ idle at any moment: still 5k cores!
 - ✱ Comparable resources to T2

CHALLENGES

☼ Problem: CRC condor machines don't have CMS/OSG software

☼ Solution: Use CVMFS + Parrot

[Good talk last year by Dan Bradley](#)

☼ Status: Successful small scale tests: working to scale to 1000's of cores

CRC REU:
Dillon Skeeahan
(+Paul Brenner)

☼ Problem: CRC machines need access to data at T2 sites

☼ Solution: xrootd (just like T3 workers)

☼ Status: After reconfiguration of CRC network to allow outside access on nodes, successful small scale tests

CHALLENGES

- ✱ Problem: Data stored at ND T3 not accessible to CRC nodes
 - ✱ Solution: Still TBD
 - ✱ Run xrootd server for ND storage?
 - ✱ Chirp?
- ✱ Problem: Resource management
 - ✱ Ideally jobs would overflow automatically from dedicated T3 resources to opportunistic when necessary
 - ✱ Right now, need to manually decide where to run
 - ✱ Solution: TBD
 - ✱ Natural sort of approach: Condor glidein
 - ✱ Could WorkQueue be an interesting alternative?

CHALLENGES

☼ Problem: Preemption

☼ Serious problem:

- ☼ Jobs don't own slot. Will be evicted if owner wants machine
- ☼ Current CMS workflow control doesn't handle preemption well: no checkpointing or automatic restarting

☼ Solution: ???

- ☼ In principle: address no restart issue with tuning to condor submit file generated by CMS workflow tool
- ☼ Doesn't handle issue of wasted resources when evicted jobs progress lost
- ☼ Need solution to handle: checkpointing? Run in VM (stop/restart on eviction)?
- ☼ Can imagine building solution using Chirp and/or Parrot to store checkpoint information across network to allow for graceful restarts

CONCLUSIONS

- ✿ ND T3 has transformed our group's ability to do data analysis
- ✿ Entering new territory in terms of scaling T3 resources (not counting FNAL)
- ✿ CCL tools critical to successes so far
 - ✿ Everything rides on CVMFS + Parrot
 - ✿ See opportunities for other tools to play important role(s)
- ✿ Several opportunities to innovate (network data access, workflow management, preemption, etc.)
 - ✿ Innovations can be fed back to larger T1/2/3 community