

# Physics at the Terascale: Computing Challenges and Solutions at the Large Hadron Collider

Andrew Melo  
Vanderbilt University

# About Me

- Nashville Native
  - Went to MLK for High School
- Sewanee C'07 - Computer Science and Physics
- Vanderbilt C'16 - PhD
- Currently a Post-Doc w/Vanderbilt



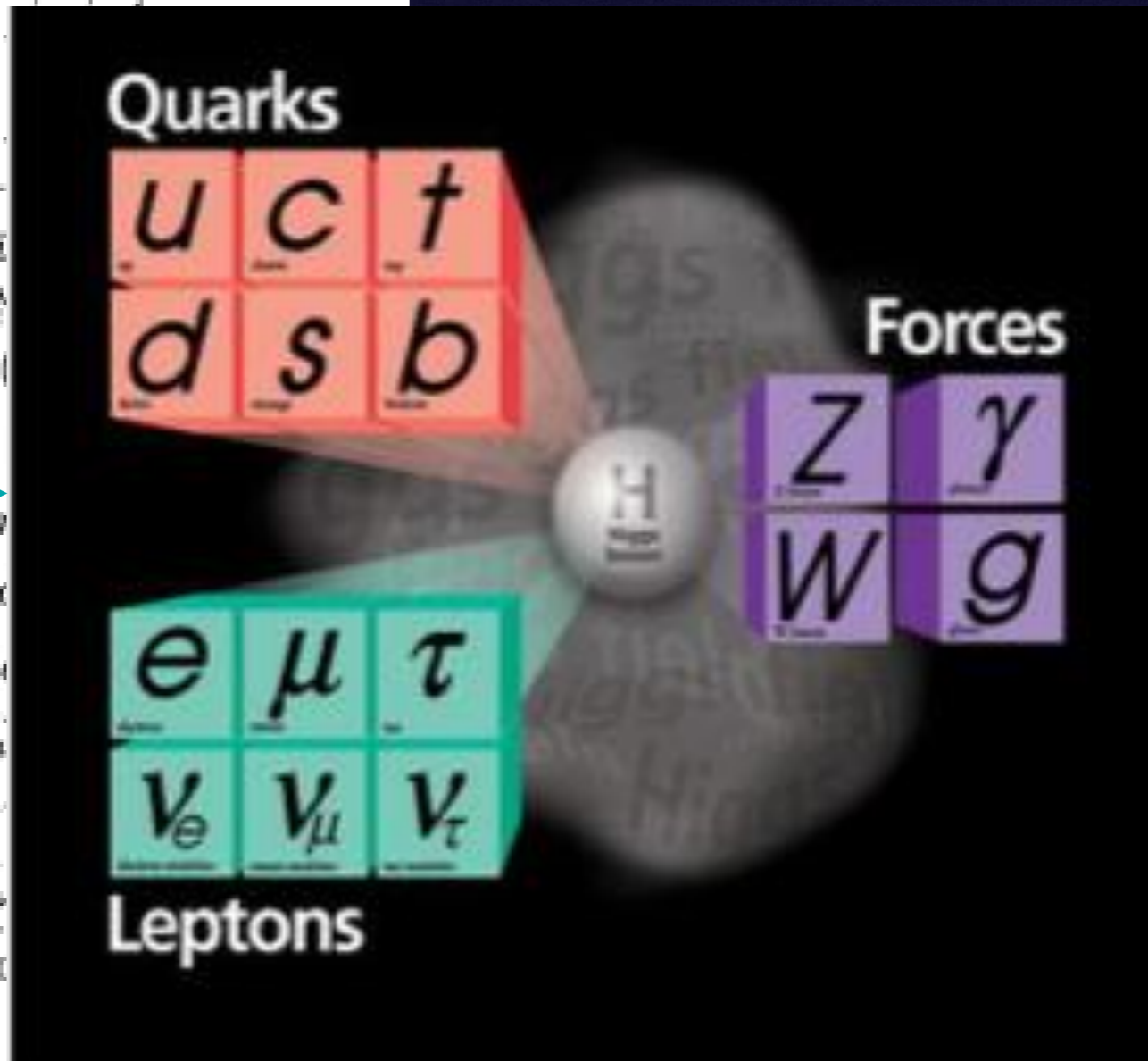
# What's the Point?

<<enter physics here>>





$$\begin{aligned}
& -\frac{1}{2}\partial_\nu g_\mu^a \partial_\nu g_\mu^a - g_s f^{abc} \partial_\mu g_\nu^a g_\mu^b g_\nu^c - \frac{1}{4}g_s^2 f^{abc} f^{ade} g_\mu^b g_\nu^c g_\mu^d g_\nu^e + \frac{1}{2}ig_s^2 (q_i^\sigma \gamma^\mu q_j^\sigma) g_\mu^a + \\
& \bar{G}^a \partial^2 G^a + g_s f^{abc} \partial_\mu \bar{G}^a G^b g_\mu^c - \partial_\nu W_\mu^+ \partial_\nu W_\mu^- - M^2 W_\mu^+ W_\mu^- - \frac{1}{2}\partial_\nu Z_\mu^0 \partial_\nu Z_\mu^0 - \frac{1}{2c_w^2} M^2 Z_\mu^0 Z_\mu^0 - \\
& \frac{1}{2}\partial_\mu A_\nu \partial_\nu A_\mu + \frac{1}{2}ig_s g_w A_\mu (\bar{e}^\lambda \gamma^\mu e^\lambda + \bar{\nu}_e^\lambda \gamma^\mu \nu_e^\lambda) + \frac{1}{2}ig_s g_w A_\mu (\bar{\mu}^\lambda \gamma^\mu \mu^\lambda + \bar{\nu}_\mu^\lambda \gamma^\mu \nu_\mu^\lambda) + \frac{1}{2}ig_s g_w A_\mu (\bar{\tau}^\lambda \gamma^\mu \tau^\lambda + \bar{\nu}_\tau^\lambda \gamma^\mu \nu_\tau^\lambda) + \\
& \frac{1}{2c_w^2} M \phi^0 \partial_\mu \partial_\mu \phi^0 - g M W_\mu^+ W_\mu^- H - \frac{1}{2}g \frac{M}{c_w^2} Z_\mu^0 Z_\mu^0 H \\
& W_\nu^+ W_\nu^- - A_\nu (W_\mu^+ \partial_\nu W_\mu^- - W_\mu^- \partial_\nu W_\mu^+) + A_\mu (W_\nu^+ \partial_\nu W_\nu^- - W_\nu^- \partial_\nu W_\nu^+) - \frac{1}{2}g^2 W_\mu^+ W_\mu^- W_\nu^+ W_\nu^- + \\
& \frac{1}{2}g^2 W_\mu^+ W_\nu^- W_\mu^- W_\nu^+ + g^2 c_w^2 (Z_\mu^0 W_\mu^+ Z_\nu^0 W_\nu^- - Z_\mu^0 Z_\nu^0 W_\mu^+ W_\nu^-) + g^2 s_w^2 (A_\mu W_\mu^+ A_\nu W_\nu^- - \\
& A_\mu A_\nu W_\mu^+ W_\nu^-) + g^2 s_w c_w [A_\mu Z_\nu^0 (W_\mu^+ W_\nu^- - W_\nu^+ W_\mu^-) - 2A_\mu Z_\mu^0 W_\nu^+ W_\nu^-] - g\alpha [H^3 + \\
& H\phi^0\phi^0 + 2H\phi^+\phi^-] - \frac{1}{2}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}g \frac{M}{c_w^2} Z_\mu^0 Z_\mu^0 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_\mu^0 (H \partial_\mu \phi^0 - \\
& \phi^0 \partial_\mu H) - ig \frac{s_w^2}{c_w} M Z_\mu^0 (W_\mu^+ \phi^- - W_\mu^- \phi^+) + ig s_w M A_\mu (W_\mu^+ \phi^- - W_\mu^- \phi^+) - ig \frac{1-2c_w^2}{2c_w} Z_\mu^0 (\phi^+ \partial_\mu \phi^- - \\
& \phi^- \partial_\mu \phi^+) + ig s_w A_\mu (\phi^+ \partial_\mu \phi^- - \phi^- \partial_\mu \phi^+) - \frac{1}{4}g^2 W_\mu^+ W_\mu^- [H^2 + (\phi^0)^2 + 2\phi^+\phi^-] - \\
& \frac{1}{4}g^2 \frac{1}{c_w^2} Z_\mu^0 Z_\mu^0 [H^2 + (\phi^0)^2 + 2(2s_w^2 - 1)^2 \phi^+\phi^-] - \frac{1}{2}g^2 \frac{s_w^2}{c_w} Z_\mu^0 \phi^0 (W_\mu^+ \phi^- + \\
& \frac{1}{2}ig^2 \frac{s_w^2}{c_w} Z_\mu^0 H (W_\mu^+ \phi^- - W_\mu^- \phi^+) + \frac{1}{2}g^2 s_w A_\mu \phi^0 (W_\mu^+ \phi^- + W_\mu^- \phi^+) + \frac{1}{2}ig^2 s_w \\
& W_\mu^- \phi^+) - g^2 \frac{s_w}{c_w} (2c_w^2 - 1) Z_\mu^0 A_\mu \phi^+ \phi^- - g^1 s_w^2 A_\mu A_\mu \phi^+ \phi^- - \bar{e}^\lambda (\gamma \partial + \\
& \partial^\lambda \gamma \partial \nu^\lambda - \bar{u}_j^\lambda (\gamma \partial + m_u^\lambda) u_j^\lambda - \bar{d}_j^\lambda (\gamma \partial + m_d^\lambda) d_j^\lambda + ig s_w A_\mu [-(\bar{e}^\lambda \gamma^\mu e^\lambda) + \frac{2}{3}(\bar{e}^\lambda \\
& \frac{1}{3}(\bar{d}_j^\lambda \gamma^\mu d_j^\lambda))] + \frac{ig}{4c_w} Z_\mu^0 [(\bar{\nu}^\lambda \gamma^\mu (1 + \gamma^5) \nu^\lambda) + (\bar{e}^\lambda \gamma^\mu (4s_w^2 - 1 - \gamma^5) e^\lambda) + (\bar{u}_j^\lambda \\
& 1 - \gamma^5) u_j^\lambda) + (\bar{d}_j^\lambda \gamma^\mu (1 - \frac{8}{3}s_w^2 - \gamma^5) d_j^\lambda)] + \frac{ig}{2\sqrt{2}} W_\mu^+ [(\bar{\nu}^\lambda \gamma^\mu (1 + \gamma^5) e^\lambda) + \\
& \gamma^5) C_{\lambda\kappa} d_j^\kappa] + \frac{ig}{2\sqrt{2}} W_\mu^- [(\bar{e}^\lambda \gamma^\mu (1 + \gamma^5) \nu^\lambda) + (\bar{d}_j^\kappa C_{\lambda\kappa}^\dagger \gamma^\mu (1 + \gamma^5) u_j^\lambda)] + \frac{ig}{2\sqrt{2}} \frac{m_u^\lambda}{M} \\
& \gamma^5) e^\lambda) + \phi^- (\bar{e}^\lambda (1 + \gamma^5) \nu^\lambda)] - \frac{g}{2} \frac{m_u^\lambda}{M} [H (\bar{e}^\lambda e^\lambda) + i\phi^0 (\bar{e}^\lambda \gamma^0 e^\lambda)] + \frac{ig}{2M\sqrt{2}} \phi^+ (-m \\
& \gamma^5) d_j^\kappa) + m_u^\lambda (\bar{u}_j^\lambda C_{\lambda\kappa} (1 + \gamma^5) d_j^\kappa) + \frac{ig}{2M\sqrt{2}} \phi^- [m_d^\lambda (\bar{d}_j^\lambda C_{\lambda\kappa}^\dagger (1 + \gamma^5) u_j^\lambda) - m_u^\lambda (C \\
& \gamma^5) u_j^\lambda] - \frac{g}{2} \frac{m_u^\lambda}{M} H (\bar{u}_j^\lambda u_j^\lambda) - \frac{g}{2} \frac{m_d^\lambda}{M} H (\bar{d}_j^\lambda d_j^\lambda) + \frac{ig}{2} \frac{m_u^\lambda}{M} \phi^0 (\bar{u}_j^\lambda \gamma^5 u_j^\lambda) - \frac{ig}{2} \frac{m_d^\lambda}{M} \phi^0 (\bar{d}_j^\lambda \\
& \bar{X}^+ (\partial^2 - M^2) X^+ + \bar{X}^- (\partial^2 - M^2) X^- + \bar{X}^0 (\partial^2 - \frac{M^2}{c_w^2}) X^0 + \bar{Y} \partial^2 Y + ig c_w W_\mu^+ \\
& \partial_\mu \bar{X}^+ X^0) + ig s_w W_\mu^+ (\partial_\mu \bar{Y} X^- - \partial_\mu \bar{X}^+ Y) + ig c_w W_\mu^- (\partial_\mu \bar{X}^- X^0 - \partial_\mu \\
& ig s_w W_\mu^- (\partial_\mu \bar{X}^- Y - \partial_\mu \bar{Y} X^+) + ig c_w Z_\mu^0 (\partial_\mu \bar{X}^+ X^+ - \partial_\mu \bar{X}^- X^-) + ig s_w A_\mu \\
& \partial_\mu \bar{X}^- X^-) - \frac{1}{2}g M [\bar{X}^+ X^+ H + \bar{X}^- X^- H + \frac{1}{c_w^2} \bar{X}^0 X^0 H] + \frac{1-2c_w^2}{2c_w} ig M [\bar{X} \\
& \bar{X}^- X^0 \phi^-] + \frac{1}{2c_w} ig M [\bar{X}^0 X^- \phi^+ - \bar{X}^0 X^+ \phi^-] + ig M s_w [\bar{X}^0 X^- \phi^+ - \bar{X}^0 \\
& \frac{1}{2}ig M [\bar{X}^+ X^+ \phi^0 - \bar{X}^- X^- \phi^0]
\end{aligned}$$

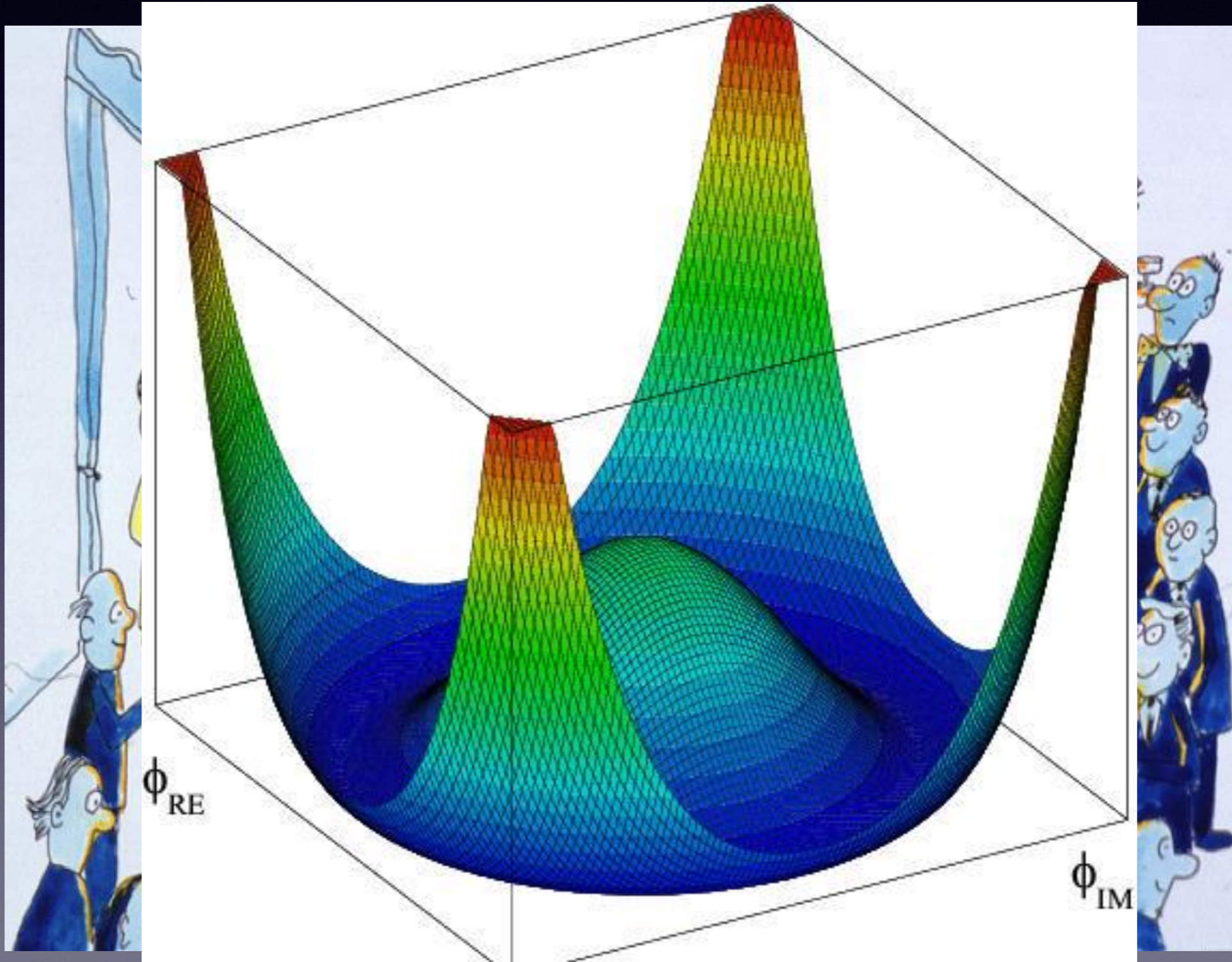


# The Higgs Boson

NOT “The God Particle”



# What Is the Higgs Field?

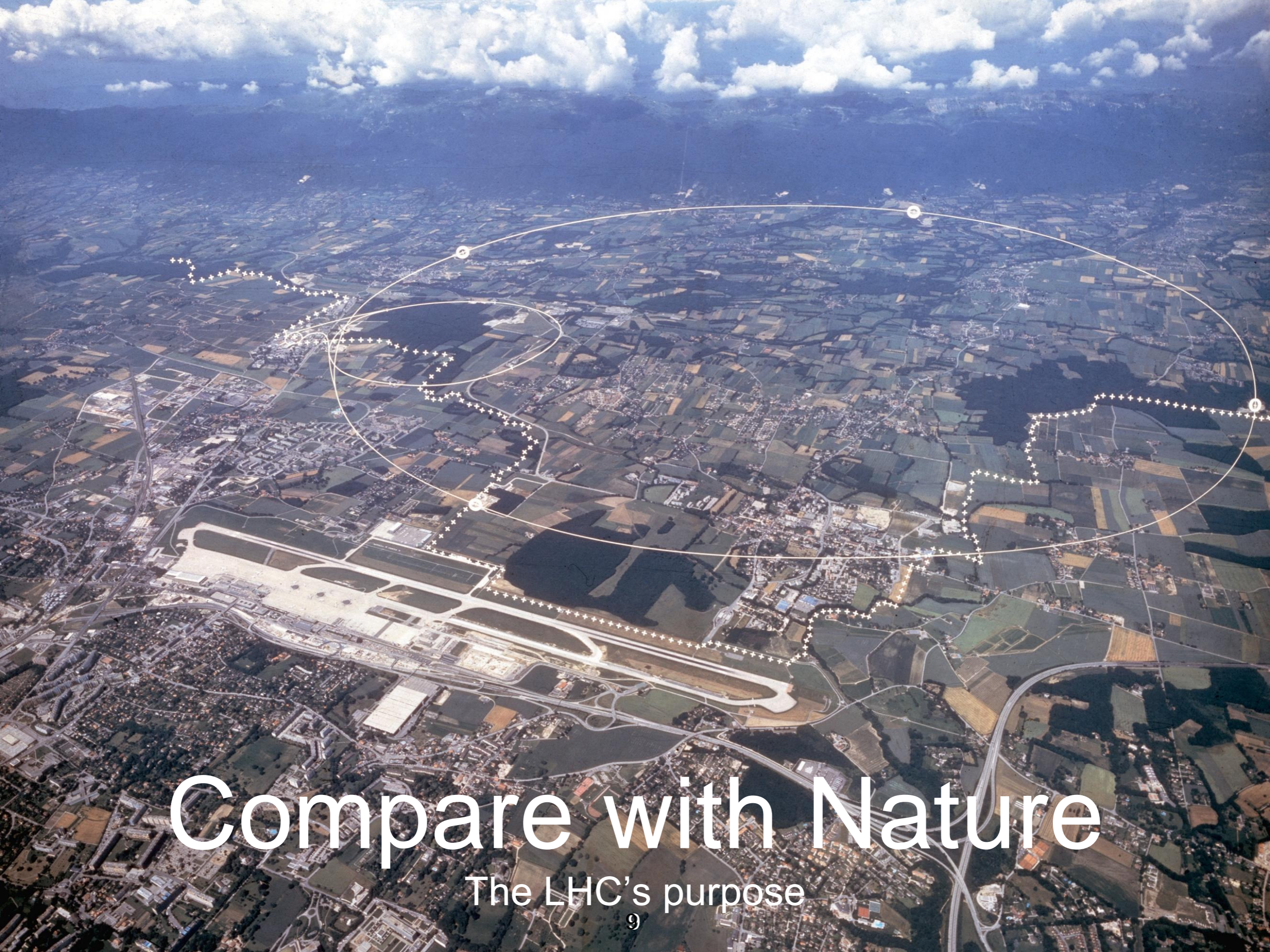




# Predict effects

- To compare with data, we need to simulate observables the detector would record
- QCD and QED don't have closed solutions
- To get enough statistics, we need  $O(1-100M)$  events
- Each event takes ~1min to simulate
  - It gets worse...(later)
- There are ~100 different backgrounds that need to be simulated
  - Plus upgrades for different detector alignments...





# Compare with Nature

The LHC's purpose



# How Is New Physics Found?

- Simulate the signal and all relevant backgrounds
- Record data from a detector
- Foreach event in (signal, background, data):
  - If event passes analysis-specific selections:
    - Save event somewhere to the side
  - If sum(signal+backgrounds) matches data:
    - Something is probably wrong, do a ton of cross-checks
  - If it still matches:
    - PressConference()

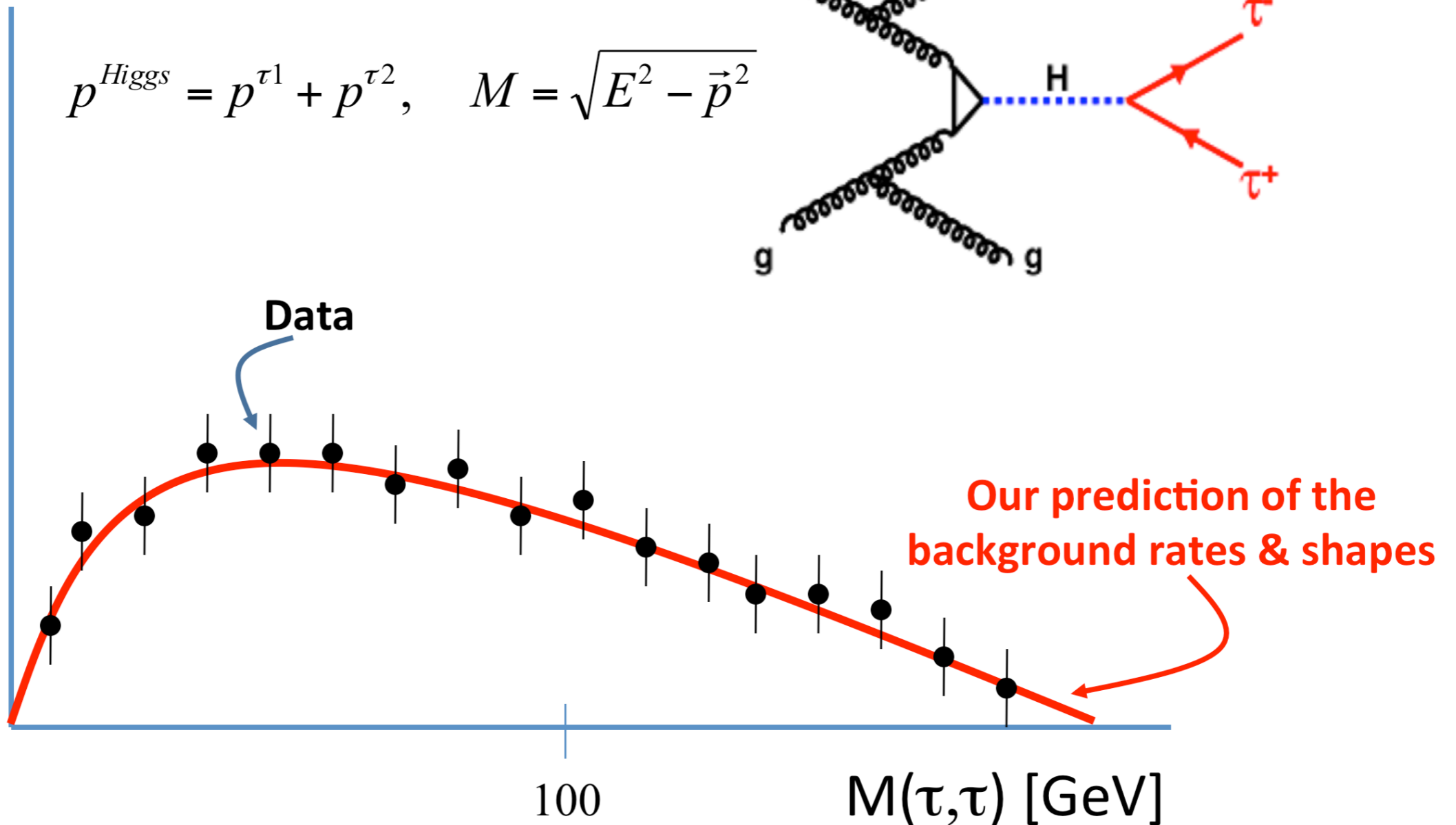
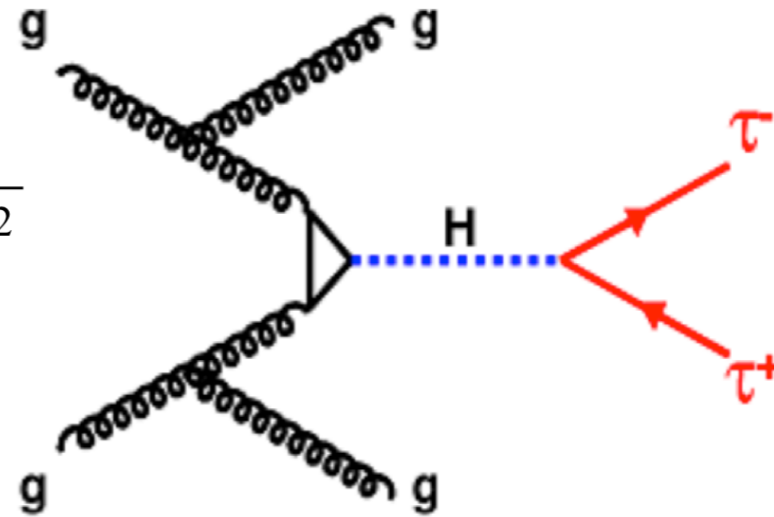




# Hunting for “Bumps”

Suppose there was NO Higgs ...  
What would our distributions look like?

$$p^{Higgs} = p^{\tau^1} + p^{\tau^2}, \quad M = \sqrt{E^2 - \vec{p}^2}$$

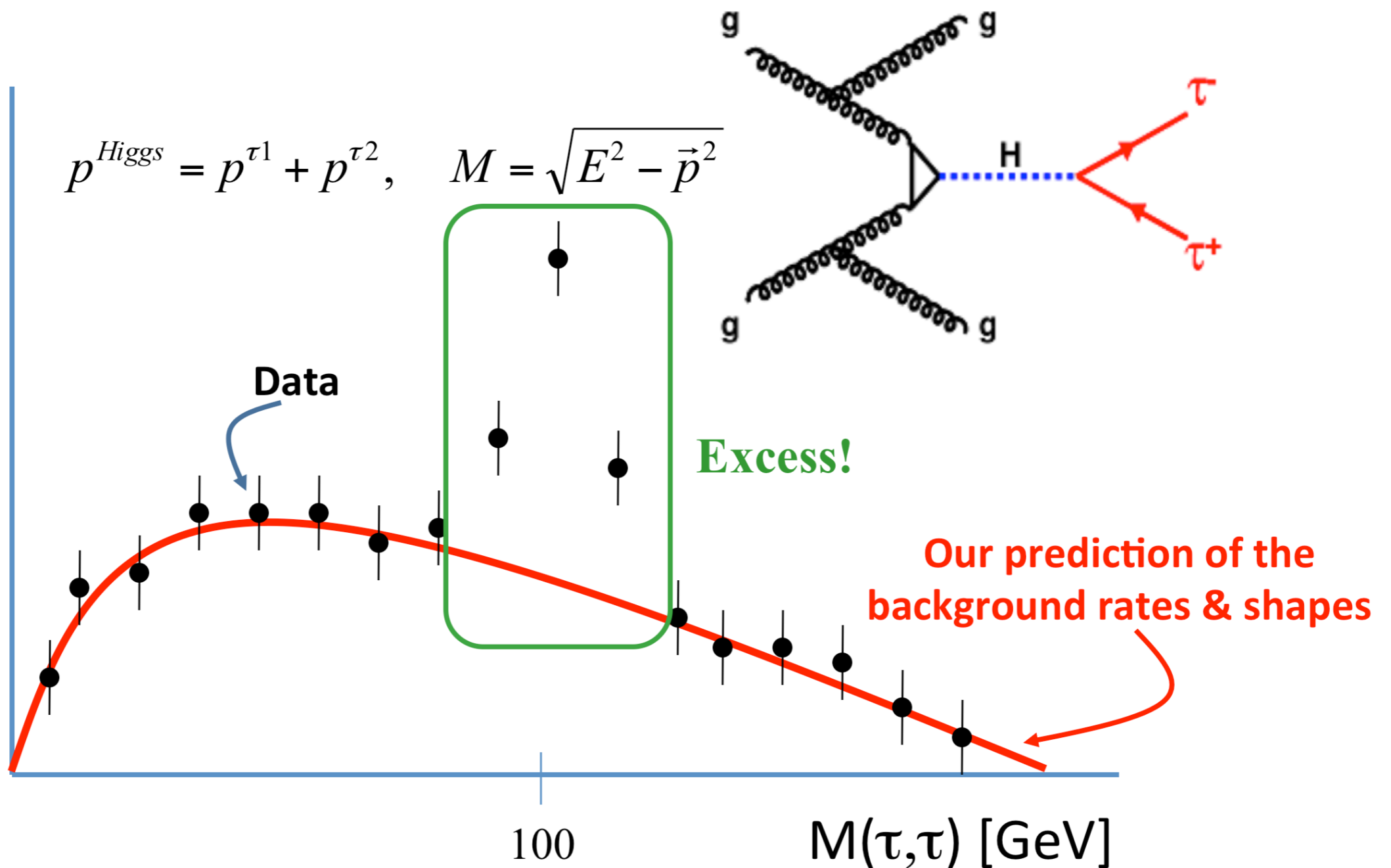






# Hunting for “Bumps”

Suppose there WAS a Higgs ...  
What would our distributions look like?

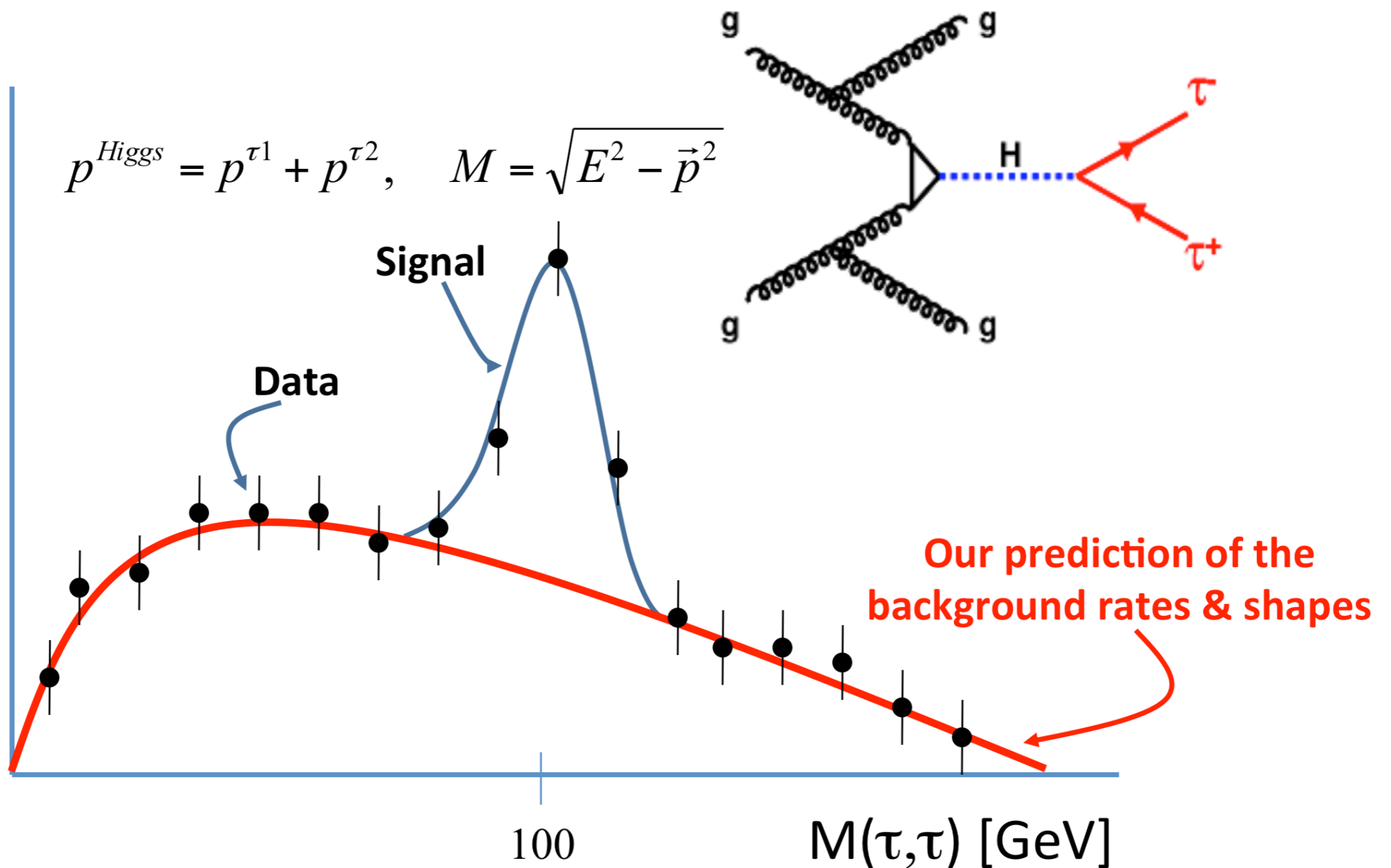






# Hunting for “Bumps”

Suppose there WAS a Higgs ...  
What would our distributions look like?





# The Problem

# How Many H Are Around Us?

0

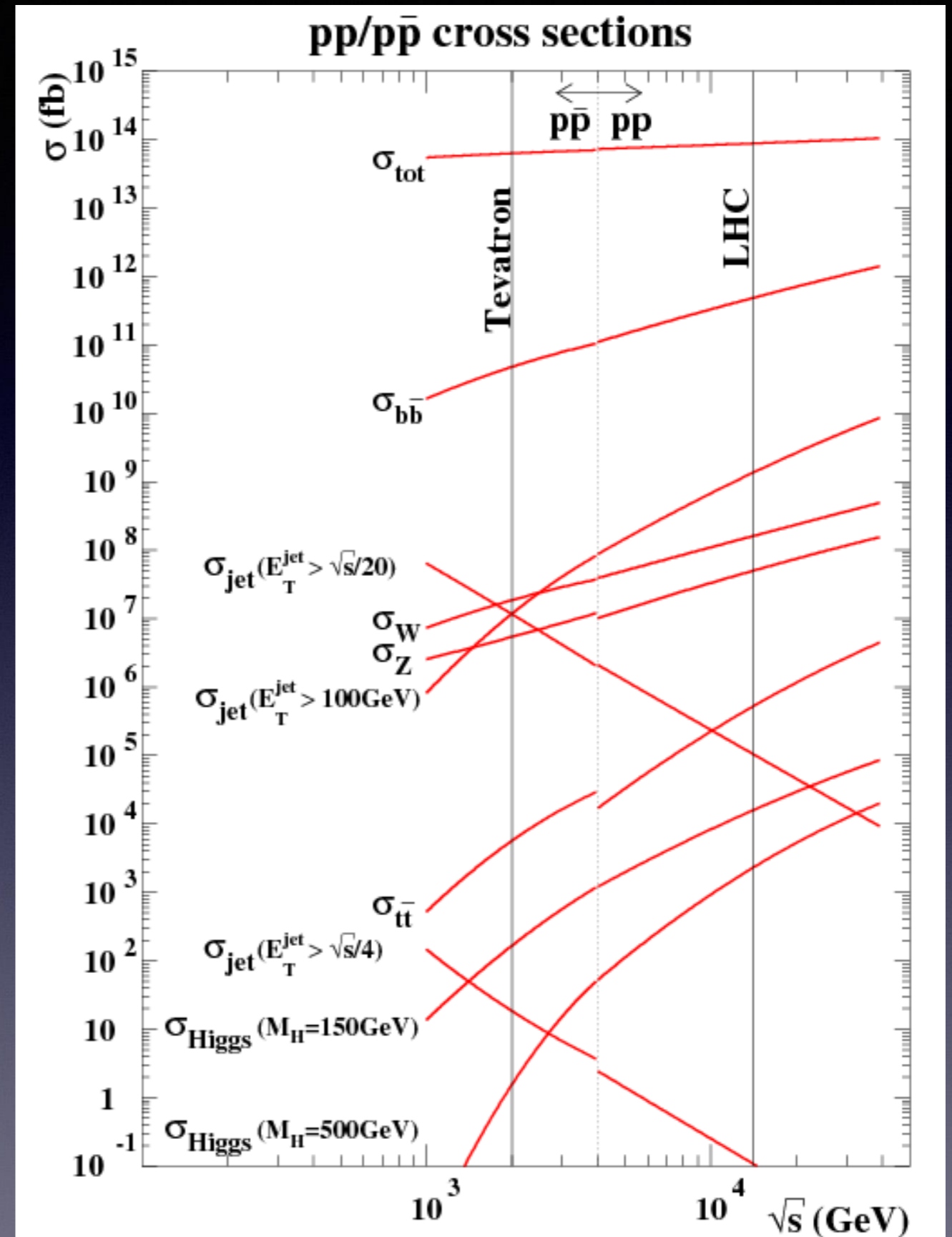
The Higgs weighs  $125 \text{ GeV}/c^2$  - Need  $125 \text{ GeV}$  of energy to produce one

Equivalent to 10.7 million degrees K



# “We’re gonna need a bigger machine”

- Higher energies "unlock" new processes
- At increasing energies, heavier particles are increasingly preferred

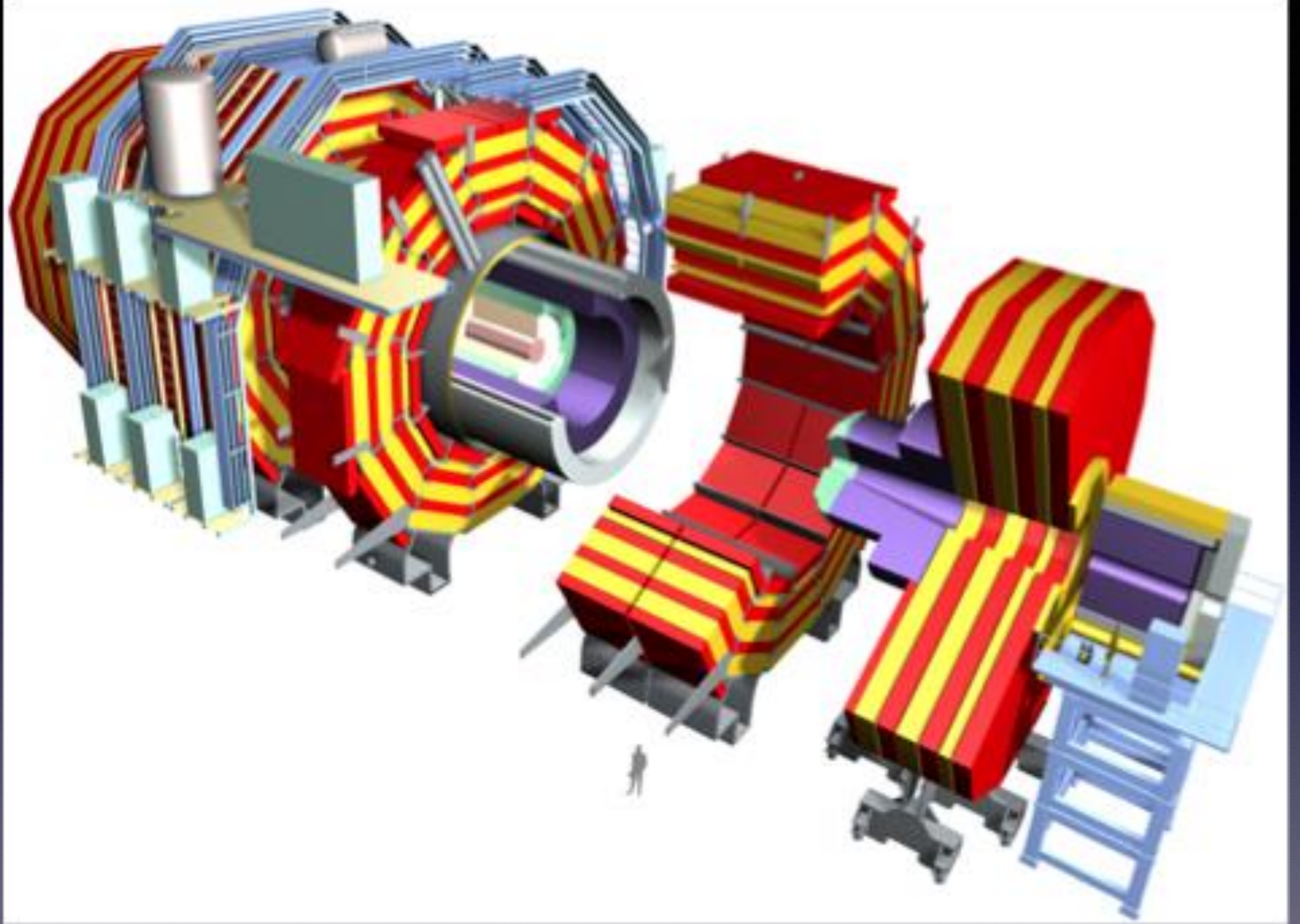


# The LHC Provides Higher Energies

- One proton-proton collision =  $10^{-2}(\text{Barn}^{-1})$
- Higgs production cross section @ 8TeV =  $10^{-12}\text{Barn}$
- Chance of producing a Higgs in a given collision:  $10^{-14}$
- We need a LOT of collision events to produce even one higgs!







How Much Data Is Produced?

# CMS Trigger System

Bunch Crossing Rate - 40MHz

Level1 Trigger - 100KHz

High Level Trigger (HLT) - 1KHz

To: Data Acquisition (DAQ) and Stable Storage



# CMS Offline System

- Once the raw data is streamed to a buffer, it needs to be reshuffled to a permanent storage format and injected to be made available for subsequent steps to analyze
- The "offline" system happens asynchronously
  - The "online" system must be active while the detector is running
- I work here!

# The Grid

- Was doing "The Cloud" before clouds were hip
- Allows the experiment to transfer data between and run jobs at ~80 sites in ~30 countries
- Federation of authentication and authorization (authn/authz) across administrative domains
- A consequence of the reality of funding



## 7 Tier-1 sites

(CPU, disk & tape)

## 52 Tier-2 sites

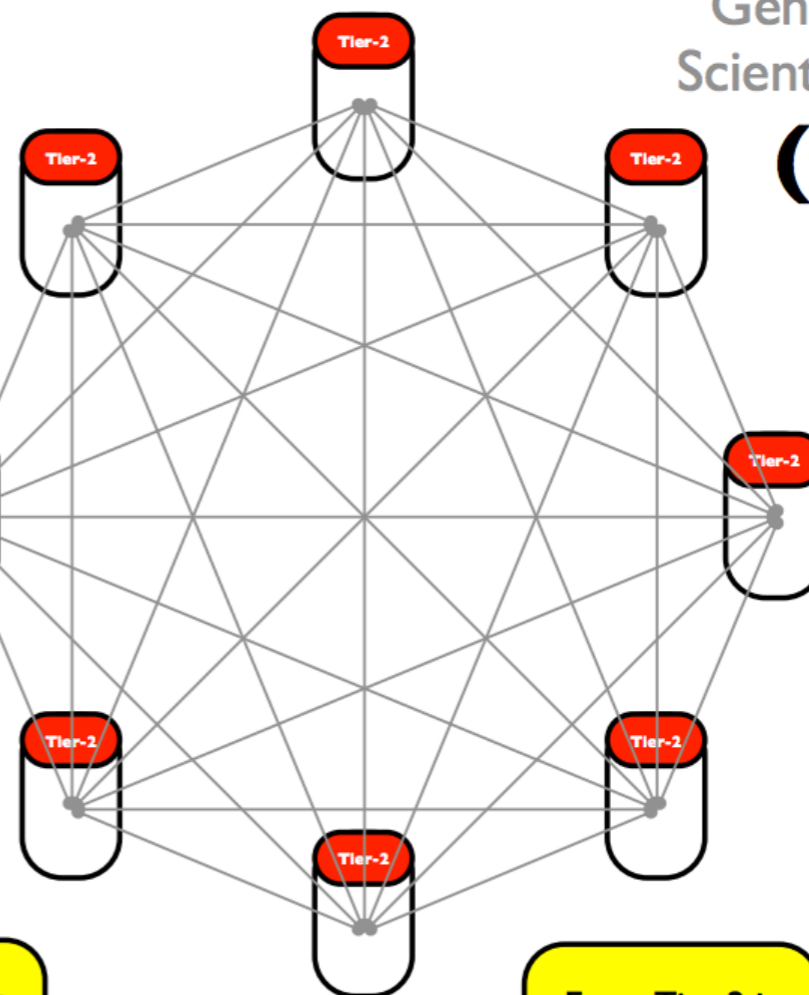
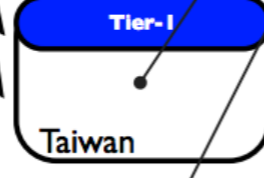
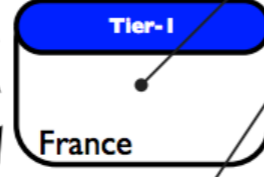
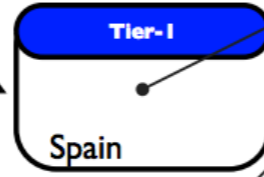
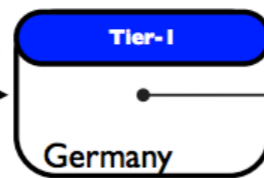
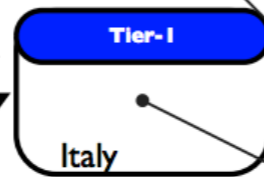
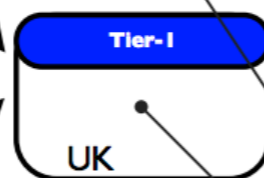
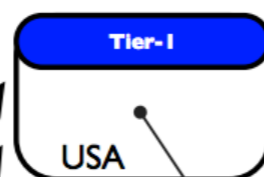
(CPU, disk)

General purpose  
Scientific Networks

**(GPN)**

Dedicated LHC  
Optical Private  
Network

**(LHCOPN)**



Every Tier-1 is  
connected to the Tier-0  
and other Tier-1 sites

Every Tier-2 is  
connected to  
every Tier-1 site

Every Tier-2 is  
connected to  
every Tier-2 site

**Full-Mesh Network Topology**

## Archiving

**Tape based:** grouping by physics content of related files on separate sets of tapes (tape families) have to be defined before files start arriving at a Tier-I site

## Serving

### T0 → T1

#### ▶ **RAW data**

- ▶ RAW data is recorded at Tier-0
- ▶ Stored on tape at CERN as “cold” backup copy
- ▶ RAW data is distributed across the Tier-1s via the LHCOPN network links
- ▶ Archived on tape and current year kept on disk for quick access

### T2 → T1

#### ▶ **Simulation output**

- ▶ Output of Monte Carlo generation and simulation is produced mainly on Tier-2s (CPU intensive workflows)
- ▶ Archived on tape distributed across the Tier-1s via the LHCOPN and GPN network links

### T1 → T1

#### ▶ **Analysis Object Data (AOD) production**

- ▶ Output of reconstructions of RAW data and simulation is produced on Tier-1s (strong I/O capabilities required)
- ▶ Archived on tape at Tier-1 site that has archival copy of input and ran reconstruction workflow

### T1 → T2

#### ▶ **Analysis Object Data (AOD) access**

- ▶ Access to reconstructed data and simulated events on Tier-2 sites
- ▶ Samples are distributed to the Tier-2s via the GPN network links
- ▶ Tier-2 disk is logically separated into managed portions for common samples and unmanaged areas for user files (ntuples)physics samples



## Production workflows

### T0

#### ▶ **Data recording**

- ▶ Collisions are recorded by the detector, selected by the trigger and sent from the detector pit to the CERN computing center in binary format
- ▶ 10% of the selected collisions are express repacked into the CMS ROOT format and reconstructed
- ▶ Latency of express reconstruction is 1 hour to allow for prompt alignment and calibration workflows and data quality monitoring
- ▶ 100% of the selected collisions are repacked into the CMS ROOT format and promptly reconstructed
- ▶ Prompt reconstruction starts 48 hours after events have been recorded to incorporate updated calibrations

### T1/2

#### ▶ **Simulation**

- ▶ Collisions are generated using theory software packages and the detector response is simulated using GEANT4
- ▶ CPU intensive workflow
- ▶ Needs little or no input

### T1

#### ▶ **Data and simulated event reconstruction**

- ▶ Collisions are reconstructed with different software versions and/or calibrations
- ▶ RAW data and simulated events archived on Tier-1s need to be pre-staged to disk for efficient access
- ▶ Reconstruction of simulated events needs to access additional datasets as input to simulate additional interactions in the detector (PileUp)

## Analysis

### T2

#### ▶ **Data and simulated event analysis**

- ▶ Collisions are accessed by physicists using officially released and self-written code
- ▶ Input is distributed beforehand through transfer system and analysis jobs are sent to data location
- ▶ Output is stored on Tier-2s outside the transfer system and catalogues but can be elevated
- ▶ Multi-user access instead of single user production mode

# Storing PBytes of Data

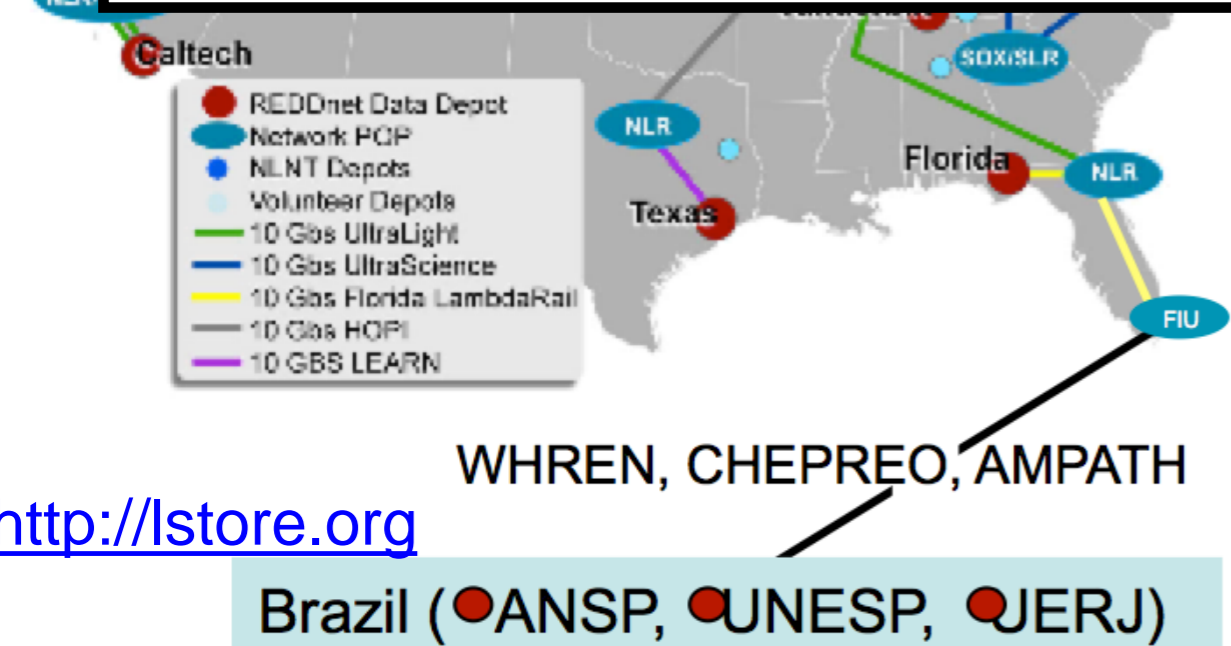
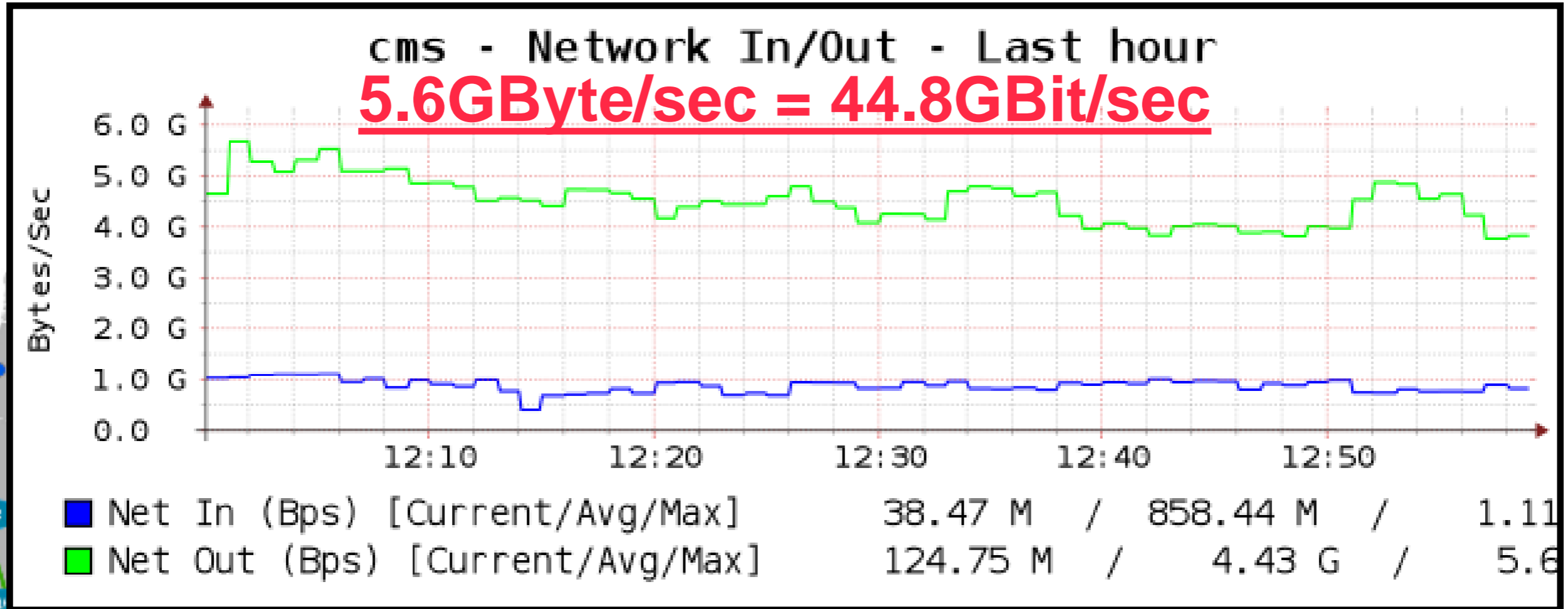


- Rule #1 - Has to be cheap
- Rule #2 - Has to be fast
- Rule #3 - Has to be reliable





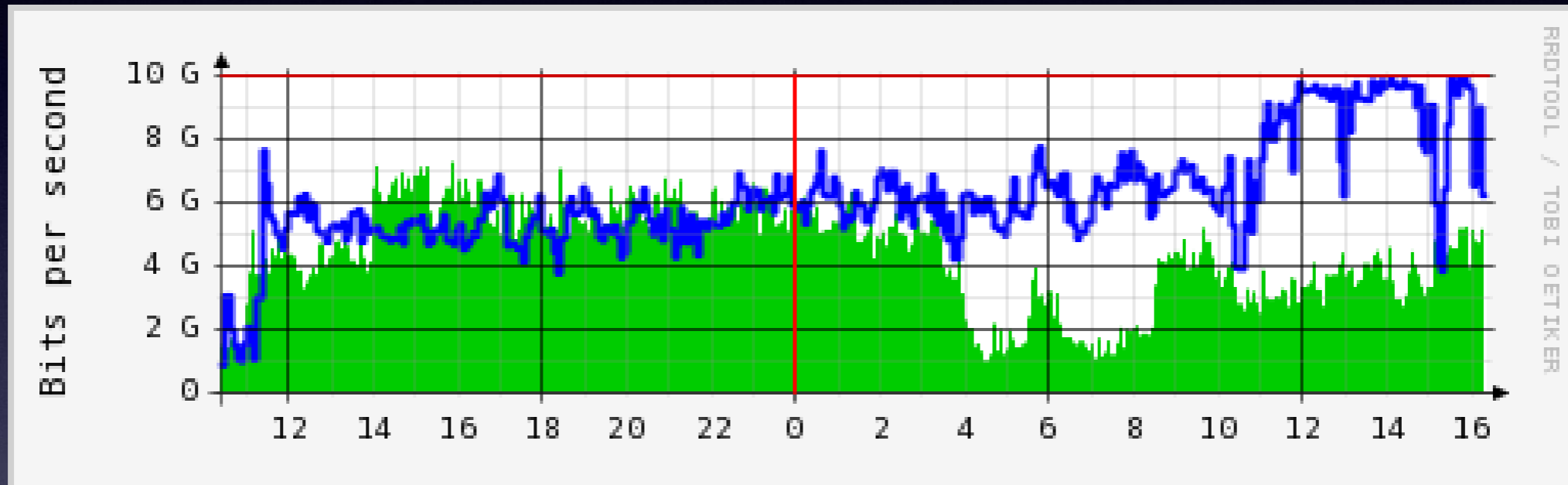
# LStore: Vanderbilt's Solution



- Satellite Imagery (AmericaView)
- High Energy Physics: LHC
- Terascale Supernova Initiative
- Structural Biology
- *Vanderbilt Campus Depot*

# Transferring over the WAN

- Managed to fully saturate a 10GBit/s link!



- Vanderbilt ITS was VERY unhappy

> =====

> A traffic policer will be configured on the research 6509 to protect enterprise Internet access. It will effectively limit traffic from the research networks to <9Gbps. The policer will not be applied to the interfaces until an outage is actually occurring, so there is not an expected service impact from this change.

>

> The impetus for this change was the recent Google outage caused by link saturation.

<http://github.com/PerilousApricot/gridftp-lfs>



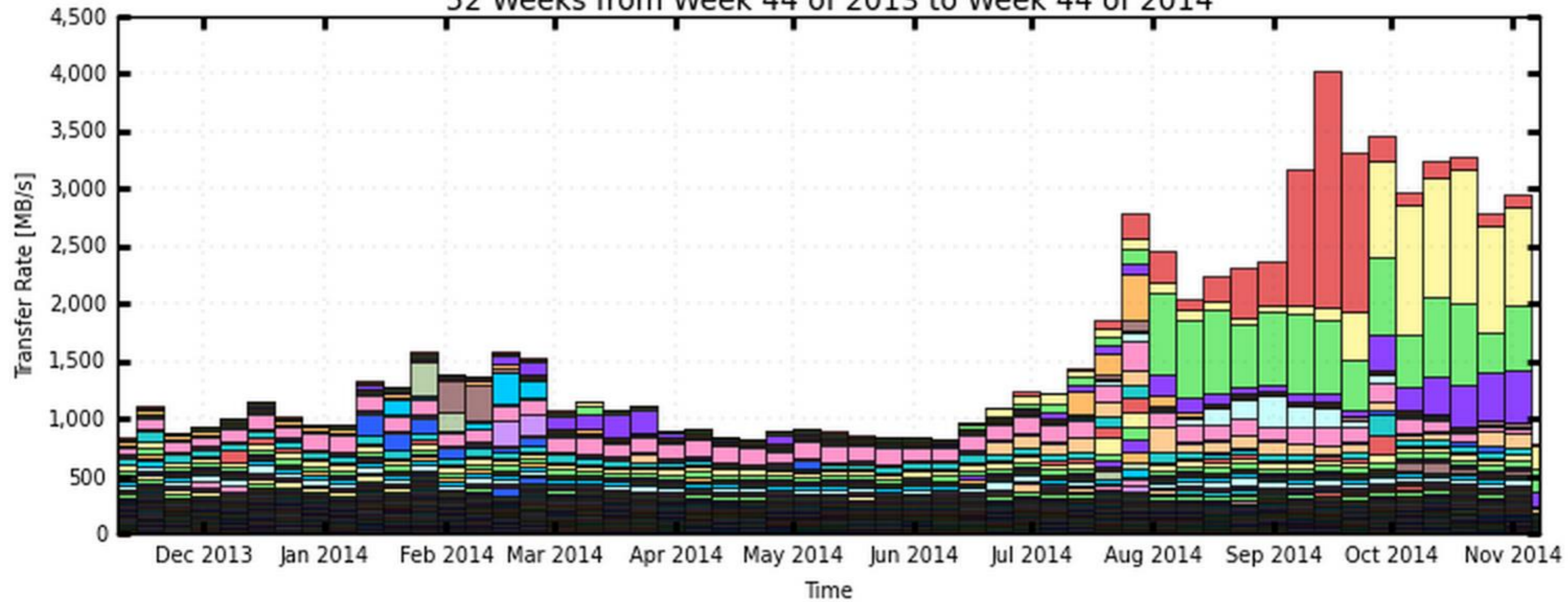
# Global Bulk Transfers

- PhEDEx is responsible for delivering data globally
  - Clever acronym, say it aloud
- Global agents make queues for files pending transfer to sites
- Local agents at each site handle moving the files they need

# ~200PB/Year

## CMS PhEDEx - Transfer Rate

52 Weeks from Week 44 of 2013 to Week 44 of 2014





# Workflow Management

- A typical analysis may use  $O(\text{PB})$  of data
- What tools are needed to leverage our computing resources to enable these workflows?
- Remember: resources are at sites with varying levels of support and performance
  - Assume the worst and expect to retry
- Choose to optimize throughput over latency
  - We have more jobs than CPUs
  - Everything will have to wait anyway

# Optimize for Users Needs

- Production system
  - Requests are well defined and long-term: “re-reconstruct all of the data from 2012”
  - Should be extensively automated:  $n\text{Tasks} \ggg n\text{Humans}$
- Analysis system
  - Requests are short-term and ill-defined: users are REAL good at breaking things
  - Can fall back to the user more often



# WMAgent and CRAB

- Each is built off the same framework (I spent 3 years here)
- WMAgent - for production
  - Complete lifecycle for data from detector to scientist-usable form
  - A central request manager generates WorkQueueElements for distributed and isolated WMAgents to consume
    - Long queues = resiliency to failure!
- CRAB - for analysis (I worked on the current rewrite)
  - Lets a user say “Let me find all of the events with two electrons”

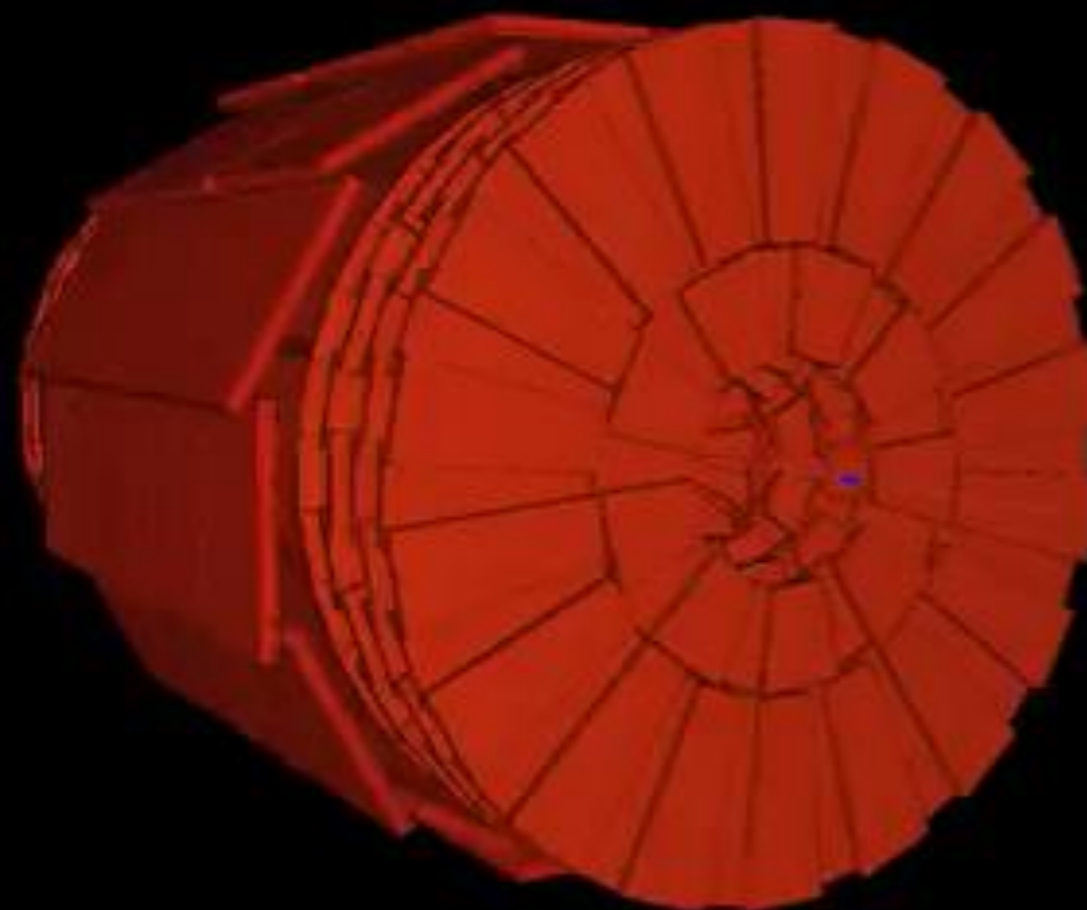
# CMS Software (CMSSW)

- Software framework and executables which handle reading/writing/analyzing all CMS data
- ~1.5M lines of C++
  - Previous team member on C++ standardization committee
  - Very active in GCC development
- C++ modules linked by simple python configuration language
- Makes the easy stuff easy and the hard stuff possible



# Future Plans

- LHCs upgrade completed last year, which means not only a higher luminosity, but much more complex events
- GPU-ization of algorithms
  - Tracking/Simulation
- Better networking
  - Most sites moving to 100GBt WAN links
- Simplified on-disk formats
  - Faster to read
  - Smaller on disk
- Multicore processing



# The End Result

One candidate Higgs → 4mu event

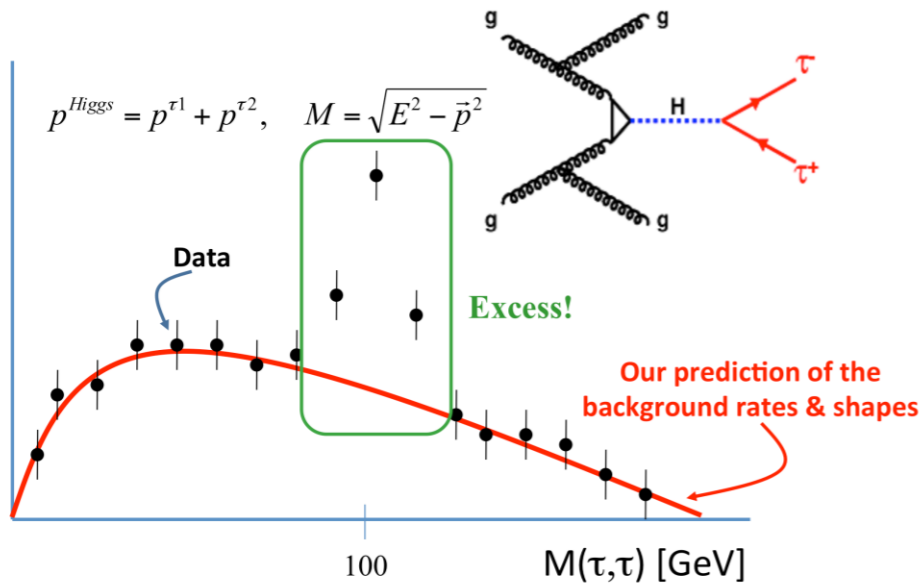


# The End Result



## Hunting for "Bumps"

Suppose there WAS a Higgs ...  
What would our distributions look like?

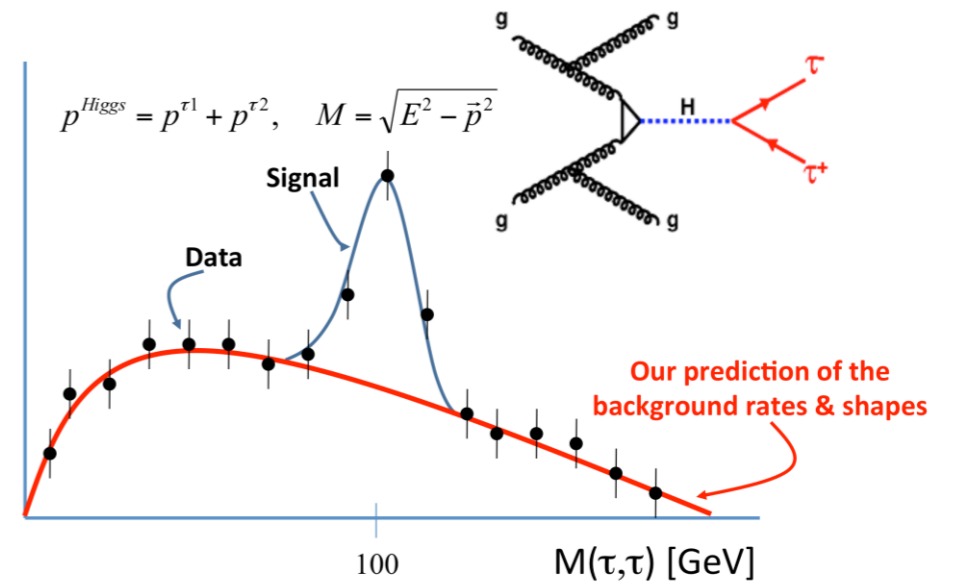


27



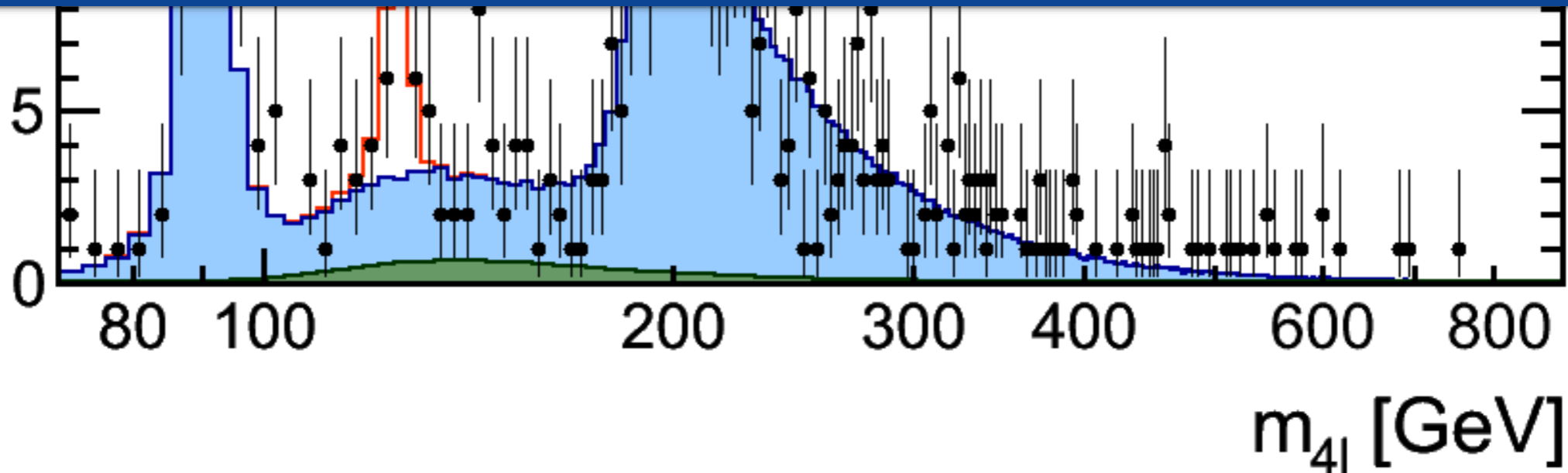
## Hunting for "Bumps"

Suppose there WAS a Higgs ...  
What would our distributions look like?



28

# What we hoped for!



# See more!

The image shows a browser window displaying the CERN Open Data Portal. The browser's address bar shows the URL `opendata.cern.ch`. The website header includes the logo "opendata CERN" and navigation links for "ABOUT", "SEARCH", "EDUCATION", and "RESEARCH". A prominent red stamp with the word "DEMO" is overlaid on the top navigation area. The main content area features a central diagram of particle tracks with various Greek letters ( $\mu$ ,  $\gamma$ ,  $\tau$ ,  $e$ ,  $q$ ) marking specific points. Two large white boxes are positioned on either side of the diagram. The left box is titled "Education" and contains the text "Visualise events, check reconstructed data, find new tools or build your own!" with a "Start learning" button below it. The right box is titled "Research" and contains the text "Get the genuine working environments, virtual machines and datasets to support your research" with a "Start analysing" button below it. A vertical sidebar of three circles is visible on the far right edge of the page.

opendata  
CERN

**DEMO**

ABOUT SEARCH EDUCATION RESEARCH

## Education

Visualise events, check reconstructed data, find new tools or build your own!

Start learning

## Research

Get the genuine working environments, virtual machines and datasets to support your research

Start analysing