# Overview (Exp)

| | Morning | Afternoon |
|---|---|---|
| Tuesday 16th | **Lecture 1**<br>**LHC Data Analysis** | **Exercise 1** |
| Wednesday 17th | **Lecture 2**<br>**LHC Statistics** | **Exercise 2** |
| Thursday 18th | | **Lecture 3**<br>**LHC Statistics** |
| Friday 19th | | **Exercise 3 (&4)** |

Slides for today inspired by
**Sourabh Dube (VSOP-28 2022)**

# Our understanding of matter



Physical Size

$\sim 10^{-10}$m    $\sim 10^{-14}$m    $< \sim 10^{-15}$m

Matter    Atom    Electron    Proton    Quarks    Forces

Nucleus    Neutron    Leptons

# Our understanding of matter



Physical Size

~$10^{-10}$m ~$10^{-14}$m < ~$10^{-15}$m

Matter  Atom  Electron  Proton  Nucleus  Neutron

Quarks  Forces  Leptons

de Broglie relationship

$$\lambda = \frac{hc}{E}$$

Energy (density)

Probing finer structure requires higher energy densities → Particle Collisions probe fine structure of Nature

# HISTORY OF THE UNIVERSE

Dark energy accelerated expansion

Cosmic Microwave Background radiation is visible

Structure formation

TODAY

Accelerators

RHIC & LHC heavy ions

LHC protons

High-energy cosmic rays

Size of visible universe

Inflation

Big Bang

POSSIBLE DARK MATTER RELICS

NUCLEONS FORM

NUCLEI FORM

$t = 10^{-36} s$
$E = 10^{16}$ GeV

$t = 10^{-10} s$
$E = 10^2$ GeV

$t = 10^{-4} s$
$E = 10^{-1}$ GeV

$t = 10^2 s$
$E = 10^{-4}$ GeV

$t = 3 \times 10^5$ y
$E = 3 \times 10^{-10}$ GeV

$t = 10^9$ y
$E = 10^{-12}$ GeV

$t = 13.8 \times 10^9$ y
$E = 2.3 \times 10^{-13}$ GeV

t = Time (seconds, years)
E = Energy of photons (units GeV = 1.6 × 10⁻¹⁰ joules)

## Key

| | | | | |
|---|---|---|---|---|
| q | quark | ν | neutrino | ion |
| g | gluon | W Z | bosons | star |
| e | electron | | | atom | galaxy |
| μ | muon | qq̄ | meson | photon |
| τ | tau | | baryon | black hole |

The concept for the above figure originated in a 1986 paper by Michael Turner.

# The LHC



**The Large Hadron Collider at CERN is a fundamental physics experiment!**

- 27 km in circumference
- 100m underground
- Accelerates protons to 99.9999991% x speed of light
- Proton circles 11,245 times per second!

At center-of-mass energies of **13.6 TeV**, proton collisions probe **physics around the time of the big-bang!**

# Proton Collisions

Unlike electron-positron colliders, proton collisions are messy but can probe **a huge range of energies simultaneously!**

$Q^2 = 10\ \text{GeV}^2$

g/10

u

d

c,c̄

s,s̄

ū

d̄

Fraction of proton momentum carried

# Open questions in Particle Physics

- Is the Higgs sector SM-like ?

- What is Dark Matter (DM)?

- Why is there more matter than anti-matter?

- What is the fundamental nature of neutrinos?

- What is (or is there) a quantum description of gravity?

- …

**Data analyses at the LHC**

Precision SM/Top/Higgs

Indirect searches for BSM

Spectroscopy & Flavour Physics

Direct searches for BSM/Exotic particles

Heavy Ions

**ATLAS and CMS** are the two **General Purpose Detectors** at the LHC

**LHCb** optimized for **flavour** physics and ALICE optimized for Heavy Ion collisions

Each is designed to detect the products that are produced in the proton-proton collisions

Extremely large-scale machines are required to reconstruct the microscopic events



Large Hadron Collider beauty (LHCb)

A Toroidal LHC ApparatuS (ATLAS)

Compact Muon Solenoid (CMS)

A Large Ion Collider Experiment (ALICE)

Data

Reconstruction & Particle ID

Calibrations

Event Selections & Distributions

RESULTS

Data

Reconstruction & Particle ID

Calibrations

Event Selections & Distributions

RESULTS

# Co-ordinate system

Co-ordinate system chosen around design of detector & collision system

$$\eta = -\ln\left[\tan\left(\frac{\theta}{2}\right)\right]$$

p

$\eta = 0$

$\eta < 0$

$\eta > 0$

$\vec{p}$

$\eta$

$\phi$

$\theta$

$\phi$

IP

$x$

$\eta = -\infty$

p

$\eta = \infty$

$z$

N

Jura

ATLAS

center of
the LHC

$y$

$\phi$

$x$

Typically deal with *transverse* projections as in this plane the incoming momentum is zero!

**The CMS Detector**

**The CMS Detector**

Different elements of the detector designed to identify and reconstruct different stable particles that are produced

Silicon Tracker
|η| < 2.4

Electromagnetic Calorimeter

Hadron Calorimeter

Superconducting Solenoid

Iron return yoke interspersed with muon chambers

4T

2T

0 m        1 m        2 m        3 m        4 m        5 m        6 m        7 m

Key:

—— Muon        —— Electron        —— Charged hadron (e.g. pion)

- - - Neutral hadron (e.g. neutron)        - - - - Photon

Nicholas Wardle

16

# Forming Tracks



B

r

p          p

Charged particles travelling through silicon track layers (pixels/strips) will create electrons / hole pairs
→ Electrons drift where charge can be read-out
→ Localized "**hits**" in the tracker layer

r

# Forming Tracks

Tracking algorithm combines hits along path → track is formed!

- Radius of curvature → momentum
- Charge ID from direction of bending
- Angles of trajectory wrt beamline
- Impact parameters (offset wrt interaction point)

$$F = q\mathbf{v} \times \mathbf{B}$$

*B*

r

p    p



Track 27.
pt = 0.51
eta = -0.078
phi = -2.477

CMS Experiment at LHC, CERN
Data recorded: Sun Oct  2 03:37:08 2011 PDT
Run/Event: 177719 / 636545915
Lumi section: 647

# Calorimetery

PbW0$_4$ CMS, X$_0$=0.89 cm



*e* →

Calorimeter layers are designed to absorb particle energy: E.g electron bremsstrahlung in ECAL / pair production produces showers which evolve through calorimeter material



Pair produce
e$^+$e$^-$ pair

Radiated
photon

Incoming
electron

Calorimeter

X$_0$     2X$_0$     3X$_0$

$x$

Electromagnetic Calorimeter (ECAL) designed to stop electrons/photons

Hadronic Calorimeter (HCAL) designed to stop hadrons



HCAL

ECAL

Different materials have different radiation lengths (X$_0$)

$$-\frac{dE}{dx} = \frac{E}{X_0}$$

# Calorimetery

Remember that different components of our detector will respond differently to different particles

Electron (e)

Photon (γ)

Electrons and photons can be identified by deposits of energy in the ECAL without **NO** deposits in the HCAL

— Track

ECAL

HCAL

Muon Chamber

# "Super clustering"

Electrons bend in the presence of a magnetic field B
→ Radiation from acceleration of charged particle
→ Photons must be included in reconstruction of electrons to maintain a good energy measurement

In each collision, the detector components measure energy deposits forming hits / tracks



CMS Experiment at the LHC, CERN
Data recorded: 2012-Aug-09 22:43:53.319400 GMT
Run / Event / LS: 200600 / 200821634 / 125

Tracks

Muon chamber hits

beamline

ECAL/HCAL deposits
(→ superclusters)

# Calorimeters

Remember that different components of our detector will respond differently to different particles

Electron (e)

Charged Hadron ($\pi^+$, p)

Track

ECAL

HCAL

Muon Chamber

Photon ($\gamma$)

Neutral Hadron ($\rho$, n)

Muon ($\mu$)

# Jet Clustering

Coloured particles (quarks & gluons) produced in proton collisions do not reach the detector components

→ Part of the production energy/momentum is used to produce additional quark/antiquark pairs – which then form hadrons. It is the hadrons that exist/escape from the collision and can be detected

q/g trajectory

Hadron

p

p

calorimeter

How can we determine energy & momentum of the original coloured particle?

# Jet Clustering

Coloured particles (quarks & gluons) produced in proton collisions do not reach the detector components

 → Part of the production energy/momentum is used to produce additional quark/antiquark pairs – which then form hadrons. It is the hadrons that exist/escape from the collision and can be detected

- - -▸ q/g trajectory

——▸ Hadron

p      p

Clustering collects particles* with original quark/gluon into single four-vector using energy-momentum conservation!

*or tracks/energy deposits …

Initial particles

Combine the 2 particles with smallest $d_{ij}$

Continue iteratively combining particles (at each step combine the protojets with smallest $d_{ij}$ )

$$d_{ij} = \frac{1}{R^2}(\Delta\eta_{ij}^2 + \Delta\phi_{ij}^2) \cdot \min\left(\frac{1}{p_{T,i}^2}, \frac{1}{p_{T,j}^2}\right)$$

$\phi$

$\eta$

$d_{ij} > \frac{1}{p_{T,i}^2}$
→ stop clustering

Found 4 jets

4 jets, each with N constituents

# b/c-jet

Identifying which particle initiated each jet requires lots of combined information about the constituents of the jet and the vertices it contains



We use sophisticated **machine learning methods** to perform this task

# Tau-leptons

τ leptons have very short lifetime → they decay into **leptons** or **hadronically**

**τ decay mode probabilities**

17.4%   17.8%

$\pi^{\pm}$   $\pi^{\pm}$   $\pi^{\pm}$   $\pi^{\mp}$   $\pi^{\pm}$

HCAL

$\mu\nu\nu$   $e\nu\nu$

$\gamma$   $\gamma$

2.7%

$3\pi^{\pm}\pi^{0}\nu$

leptonic

ECAL

9.0%   $3\pi^{\pm}\nu$   3p   $\pi^{\pm}\nu$   10.8%

$\nu_{\tau}$   $\nu_{\tau}$   $\nu_{\tau}$

$\pi^{0}$   $\rho^{0}$   hadronic

$\rho^{\pm}$   $a_{1}^{\pm}$   others   1p

tracker   7.5%

$\tau^{\pm}$   $\tau^{\pm}$   $\tau^{\pm}$

$\pi^{\pm}2\pi^{0}\nu$   $\pi^{\pm}\pi^{0}\nu$

$\tau^{\pm} \to \pi^{\pm}\nu_{\tau}$   $\tau^{\pm} \to \pi^{\pm}\pi^{0}\nu_{\tau}$   $\tau^{\pm} \to \pi^{\pm}\pi^{\mp}\pi^{\pm}\nu_{\tau}$

9.3%   25.5%

Most modern τ-ID strategies use **machine learning** to identify the decay mode and reconstruct the τ four-momentum

# Physics objects are formed by clustering certain tracks & energy deposits

CMS Experiment at the LHC, CERN
Data recorded: 2012-Aug-09 22:43:53.319400 GMT
Run / Event / LS: 200600 / 200821634 / 125

**Electrons (e) & Photons (γ)**

**Jets & taus (τ)**

beamline

**Muons (μ)**

Heavy particles (H, W, Z, t) must be reconstructed from decay products

Nicholas Wardle

# Missing momentum



Neutrinos do not interact with any component of the detector

We infer the presence of neutrinos through an imbalance of momentum in the transverse plane
→ missing transverse momentum

$$\vec{p}_{T,\mathrm{miss}} = -\sum_i \vec{p}_{T,i}$$

**Data**

**Reconstruction & Particle ID**

**Calibrations**

**Event Selections & Distributions**

**RESULTS**

# Standard Candles

Z, J/ψ and Υ decays in data provide standard candles to calibrate energy/momentum measurements

Large source of clean events with a well described mass peak

# Relative corrections

More complicated objects (eg jets) require several stages of correction → Use previously calibrated objects to calibrate jet momentum!



$$R_{\text{jet},p_T} = \frac{p_{T,jet}}{p_{T,z/\gamma}}$$

Data

Reconstruction
&
Particle ID

Calibrations

Event Selections
&
Distributions

z

μ/e

μ/e

jet

RESULTS

$\gamma_2$

$\phi$

$\gamma_1$

$m_{\gamma\gamma} = 2E_1 E2(1 - \cos\phi)$

# Collecting events



Each event that we select this way builds a picture of the underlying physics → if we're lucky, we might find something new

To extract the Physics, we use **distributions of observables** across many events

# Collecting events



Each event that we select this way builds a picture of the underlying physics → if we're lucky, we might find something new

To extract the Physics, we use **distributions of observables** across many events

# Collecting events

Collisions (i.e. bunch crossings) happen at **40 MHz**

Integrated luminosity (units of 1/area) $= \int c \cdot \dfrac{N_1 N_2}{4\pi\delta_{x,y}} \cdot f n_b dt$

$$N = L\sigma$$

Number of events

**CMS**

$L_{inst} \sim 10^{34}\ cm^{-2}\ s^{-1}$

Total integrated luminosity (fb$^{-1}$)

| | |
|---|---|
| — | 2010, 7 TeV, 45.0 pb$^{-1}$ |
| — | 2011, 7 TeV, 6.1 fb$^{-1}$ |
| — | 2012, 8 TeV, 23.3 fb$^{-1}$ |
| — | 2015, 13 TeV, 4.3 fb$^{-1}$ |
| — | 2016, 13 TeV, 41.6 fb$^{-1}$ |
| — | 2017, 13 TeV, 49.8 fb$^{-1}$ |
| — | 2018, 13 TeV, 67.9 fb$^{-1}$ |
| — | 2022, 13.6 TeV, 41.5 fb$^{-1}$ |
| — | 2023, 13.6 TeV, 31.9 fb$^{-1}$ |
| — | 2024, 13.6 TeV, 40.0 fb$^{-1}$ |

× 50

Date (UTC)

# Collecting events

Collisions (i.e. bunch crossings) happen at **40 MHz**

Integrated luminosity (units of 1/area) $= \int c \cdot \dfrac{N_1 N_2}{4\pi \delta_{x,y}} \cdot f n_b dt$

$$N = L\sigma$$

Number of events

Cross-section (units of area)

Example: For every **1,000,000,000** inelastic proton-proton collisions, **only expect one of them to produce a Higgs boson**!

We typically have to select events based on these observables to dig out the signal from the background (noise!)



proton - (anti)proton cross sections

1 barn = $10^{-28}$ m$^2$

$\sigma_{tot}$

Tevatron     LHC

$\sigma_b$

$\sigma_{jet}(E_T^{jet} > \sqrt{s}/20)$

$\sigma_W$
$\sigma_Z$

$\sigma_{jet}(E_T^{jet} > 100 \text{ GeV})$

$\sigma_t$

$\sigma_{jet}(E_T^{jet} > \sqrt{s}/4)$

$\sigma_{Higgs}(M_H = 150 \text{ GeV})$

$\sigma_{Higgs}(M_H = 500 \text{ GeV})$

$\sigma$ (nb)

events/sec for L = $10^{33}$ cm$^{-2}$ s$^{-1}$

$\sqrt{s}$  (TeV)

# Event Selection

By knowing ahead of time the kind of events we are interested in, we impose selections on the events to reduce the background as much as possible while maintaining the signal

**Signal (H)**

**Background (tt)**

# Event Selection

We often use **Machine Learning techniques** to combine as much information as possible for this task



**How can we choose these selections/train our Machine Learning models?**

**Data**

**Reconstruction & Particle ID**

**Calibrations**

**Event Selections & Distributions**

**Simulation**

μ/e

z

μ/e

jet

**RESULTS**

# Simulation

Generate large number of simulated events for each process contributing to our analysis (**signals** and backgrounds)



Simulated events must be *weighted* to the get the correct predicted yield for a given dataset

$$L_{\text{eff}} = \frac{N_{\text{gen}}}{\sigma} \implies \text{weight} = \frac{L}{L_{\text{eff}}}$$

# Data-driven background



jet

proton

$\nu$

$\nu$

z

jet

proton

Missing transverse momentum

Nicholas Wardle

# Data-driven background



μ

Z

μ

Estimate the normalization of the
Z→neutrinos background using data!

$$N_{Z(\to\nu\nu)} \approx N_{Z(\to\mu\mu)} \frac{B(Z \to \nu\nu)}{B(Z \to \mu\mu)} A(\mu)\epsilon(\mu)$$

Data

Simulation

1001011010001100100100001
010...1010100010101
...010
100...100
101010101011...10...
01000100010010101101010

Reconstruction
&
Particle ID

Calibrations

z

μ/e

μ/e

jet

Event Selections
&
Distributions

RESULTS

**Huge computing power required to acquire and analyse LHC data**

**Online Selection of events "Trigger" determines which events to keep in around 4 micro-seconds!**

**Huge collaborations of people required for Data Analysis at the LHC**

# Now it's your turn!

This afternoon, we are going to have a go at doing a data analysis with some real CMS proton-proton collision data!

All of the instructions for getting setup and the exercises are available here:
https://nucleosynthesis.github.io/LHCDataStatisticsICISE2024/

You will also see links to the lecture slides (**password VSOPLHC2024**)

If you haven't already done so, **please go through the "Getting started" section** before this afternoon's session!

This afternoon, we will be working through **Exercise 1**

# (Extra Slide) Data Tiers

RAW data output from experiments is too large for direct analysis

  e.g 2018 data from CMS O(10) PB at RAW data level

→ Reduce content through processing at different data tiers to make analysis manageable

→ Less information but content is closer to final analysis objects (hits → particles)

# (Extra Slide) A Real CMS analysis selection flow