# LHC Open Data



- the LHC collaborations make good chunks of their data publicly available
  - http://**opendata.cern**.ch/
- along with tools & software & examples
- for data visualisation and analysis
- from event reconstruction algorithms to machine learning challenges
- via virtual machines (with no need to install different software packages)
- few pointers
  - http://opendata.cern.ch/visualise/events/cms
  - http://www.i2u2.org/elab/cms/event-display/
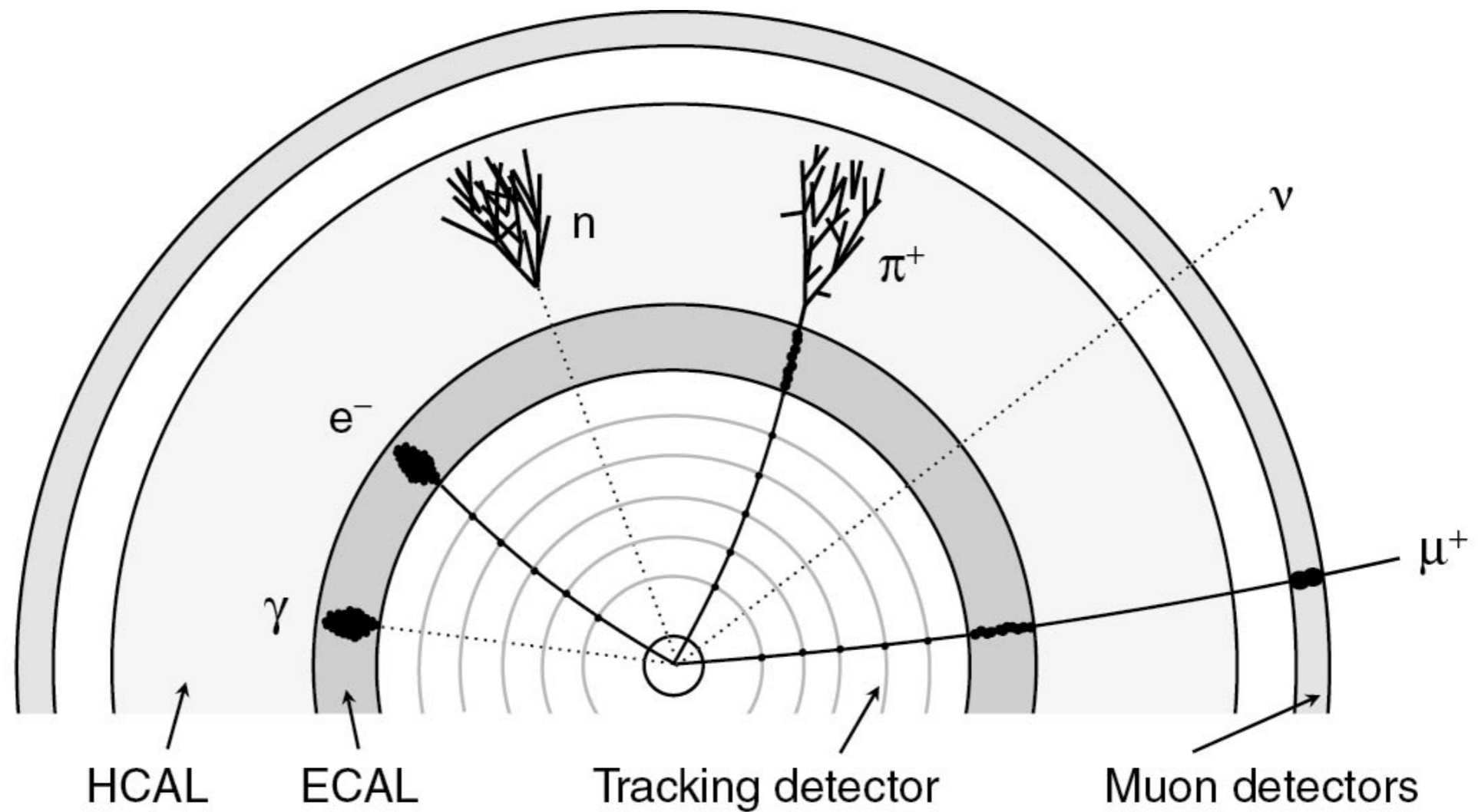- you're invited to **explore the LHC data** also on your own leisure

# goals

**perform a simple data analysis**

- visualise the data

- manipulate data ntuples

- produce, process, and display data histograms
    ‣ select different physics signals
    ‣ plot kinematic distributions, inspect detector/trigger effects

- extract physics parameters from data
    ‣ measure signal yields by performing a likelihood fit
    ‣ inspect statistical and systematic errors
    ‣ (extra) perform a differential measurement

# Detector & Event Reconstruction & Visualisation

**calorimeters**:
measure particle's
energy by absorbing it

**trackers**:
detect trajectory
of charged particles
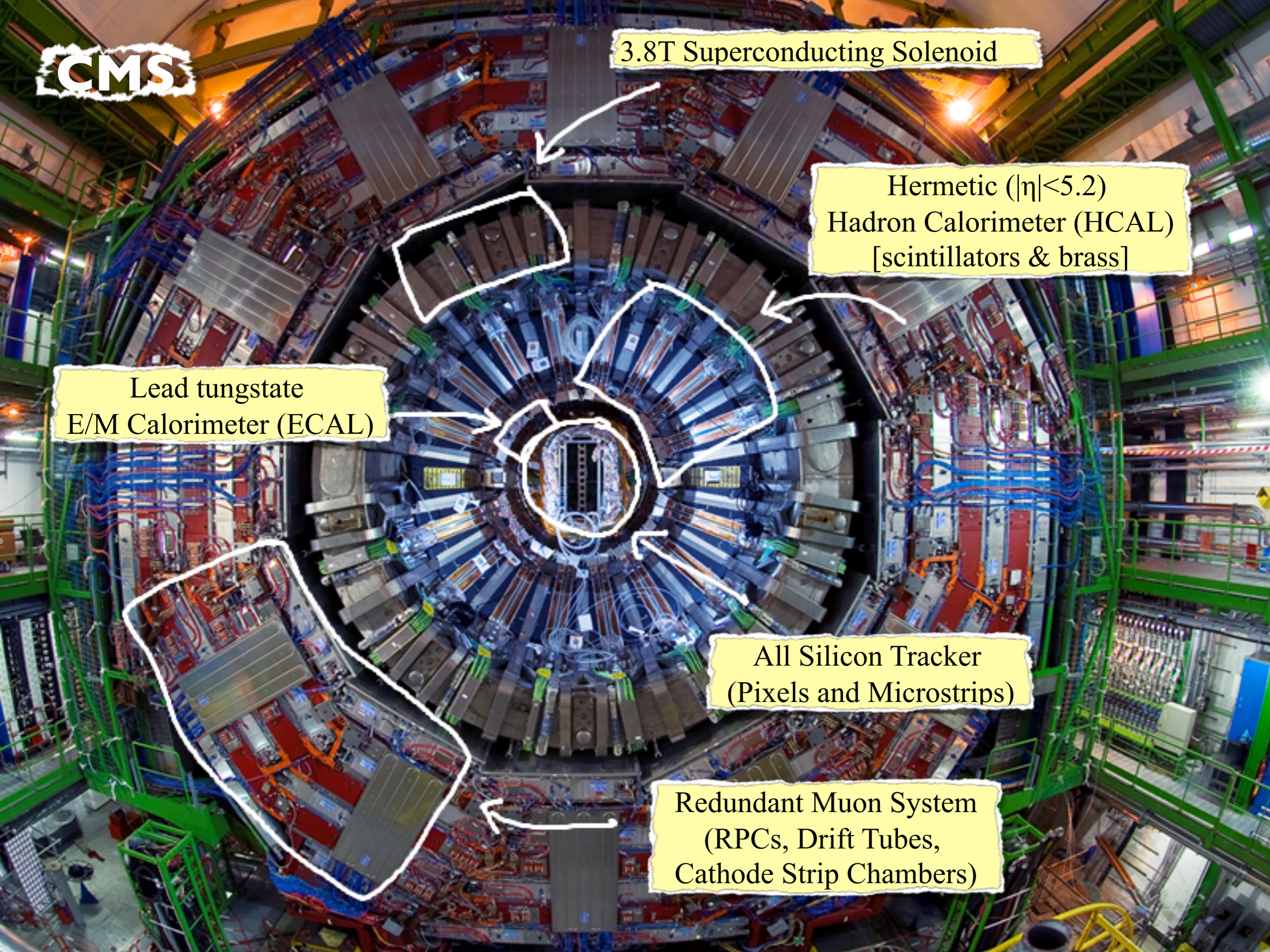
**muons**:
detected in outer
detector layers

CMS

3.8T Superconducting Solenoid

Hermetic ($|\eta|<5.2$)
Hadron Calorimeter (HCAL)
[scintillators & brass]

Lead tungstate
E/M Calorimeter (ECAL)

All Silicon Tracker
(Pixels and Microstrips)

Redundant Muon System
(RPCs, Drift Tubes,
Cathode Strip Chambers)

CMS

3.8T Superconducting Solenoid

Hermetic (|η|<5.2)
Hadron Calorimeter (HCAL)
[scintillators & brass]

Lead tungstate
E/M Calorimeter (ECAL)

All Silicon Tracker
(Pixels and Microstrips)

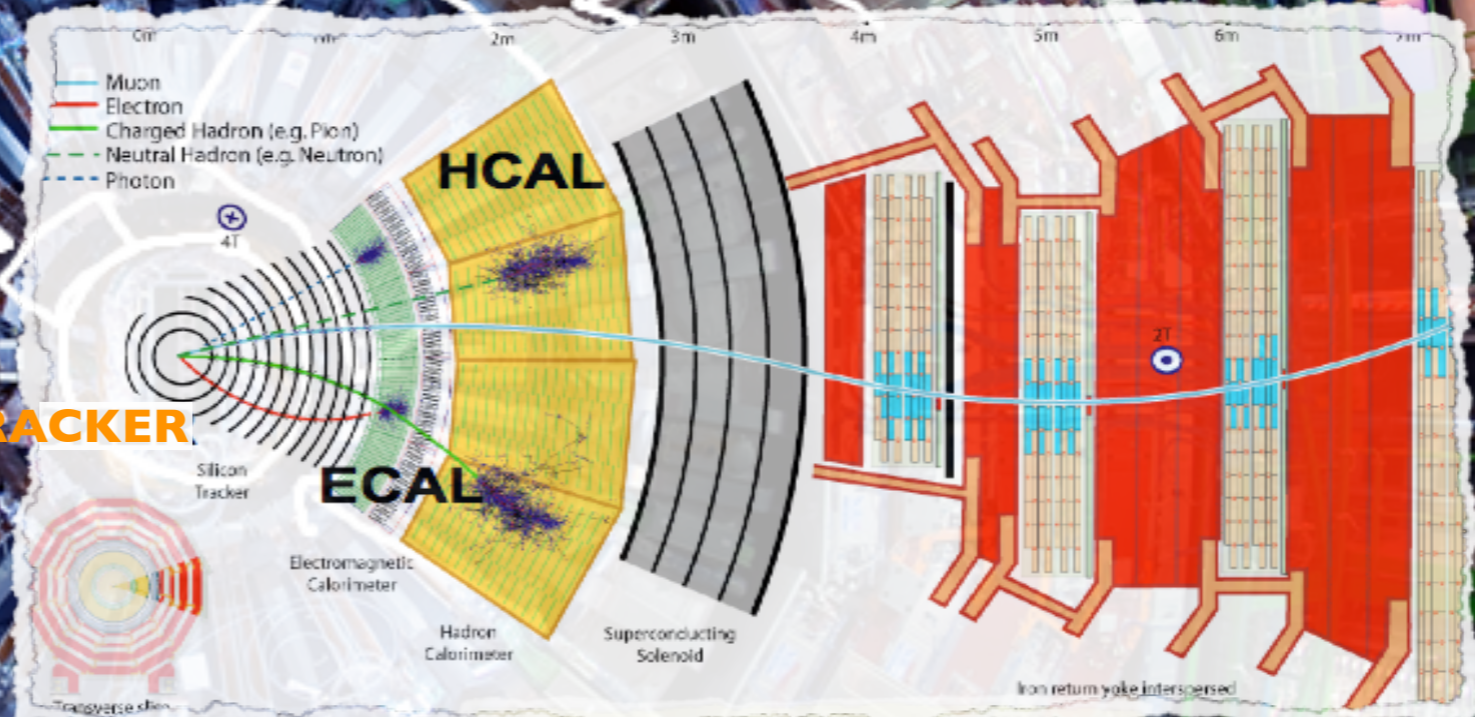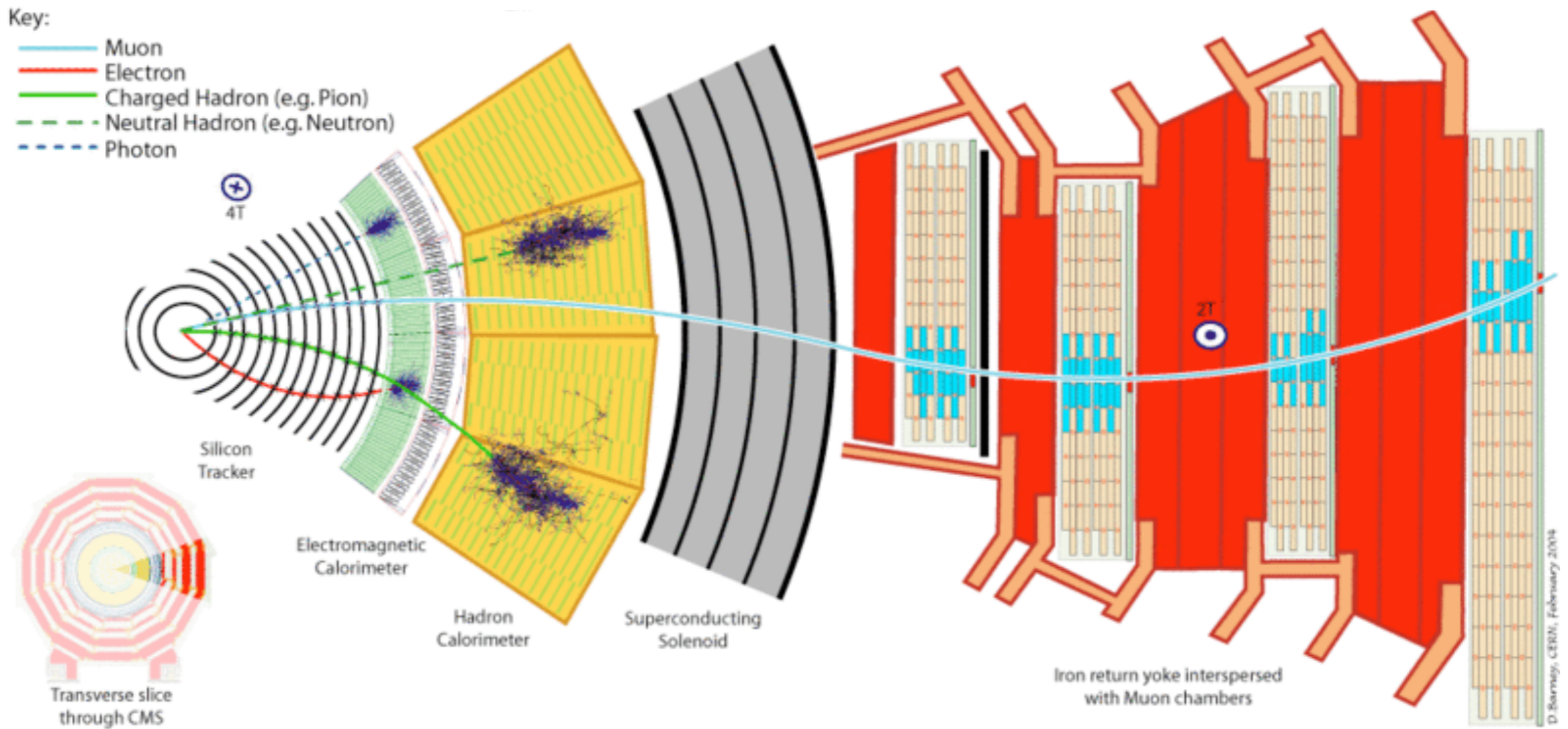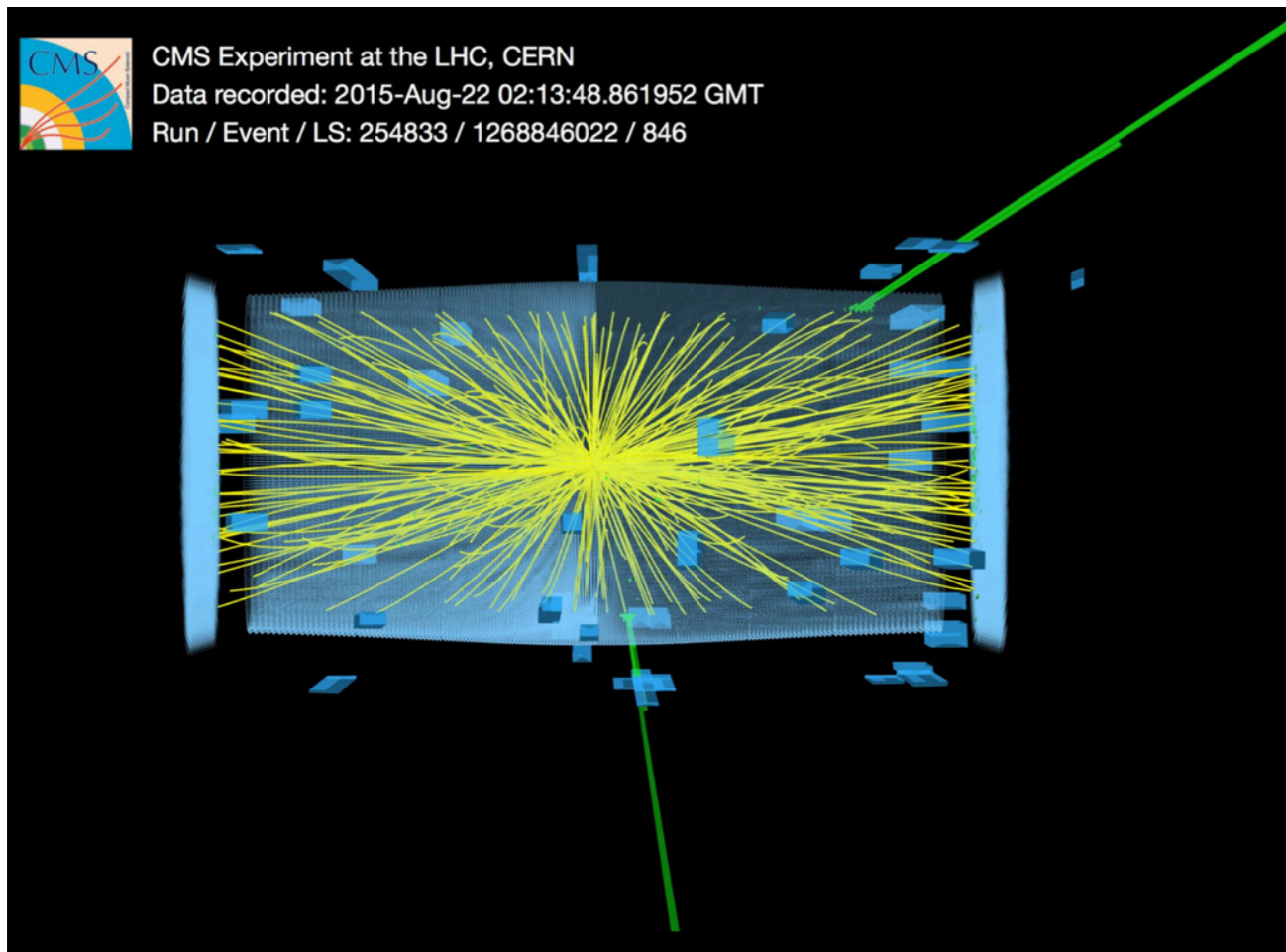Redundant Muon System
(RPCs, Drift Tubes,
Cathode Strip Chambers)
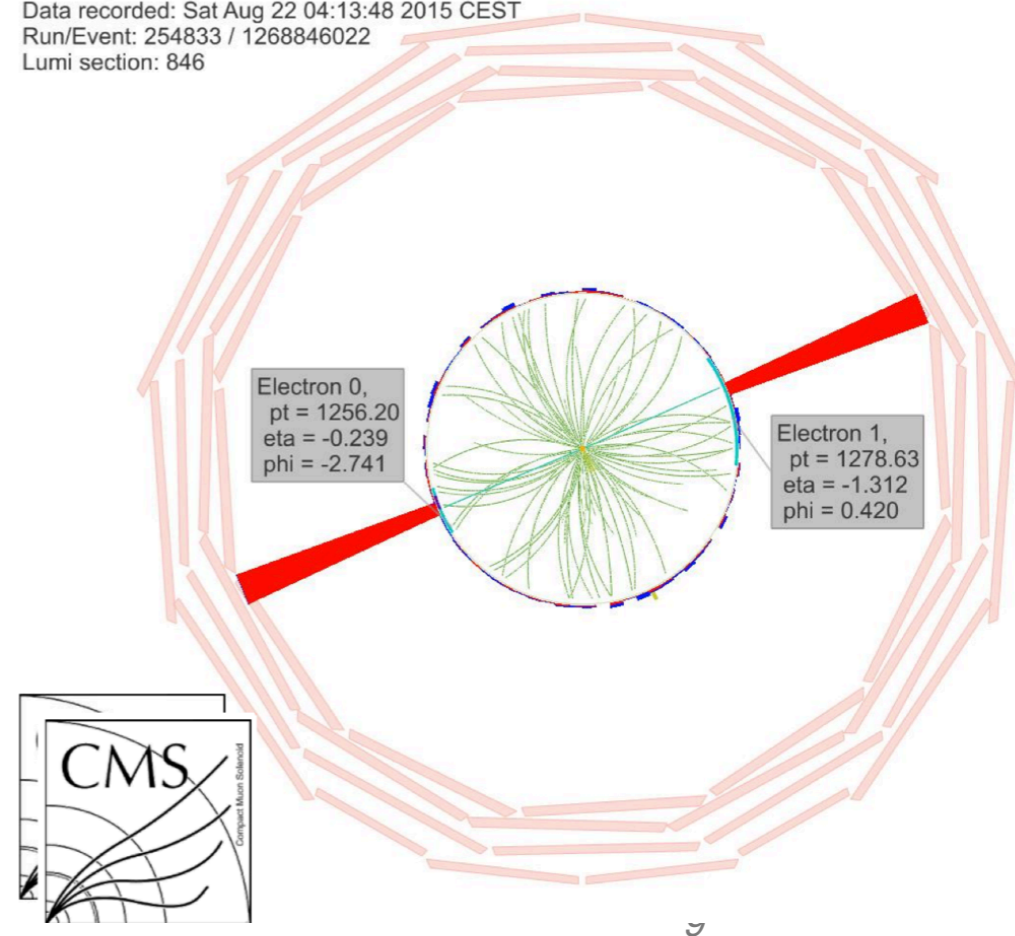
Muon
Electron
Charged Hadron (e.g. Pion)
Neutral Hadron (e.g. Neutron)
Photon

HCAL
ECAL
TRACKER

Silicon
Tracker

Electromagnetic
Calorimeter

Hadron
Calorimeter

Superconducting
Solenoid

Iron return yoke interspersed

Transverse slice

Key:
— Muon
— Electron
— Charged Hadron (e.g. Pion)
-- Neutral Hadron (e.g. Neutron)
···· Photon

4T

Silicon Tracker

Electromagnetic Calorimeter

Hadron Calorimeter

Superconducting Solenoid

2T

Iron return yoke interspersed with Muon chambers

Transverse slice through CMS

D. Barney, CERN, February 2004

# a di-electron event



Event Display of a Candidate Electron-Positron Pair with an Invariant Mass of 2.9 TeV
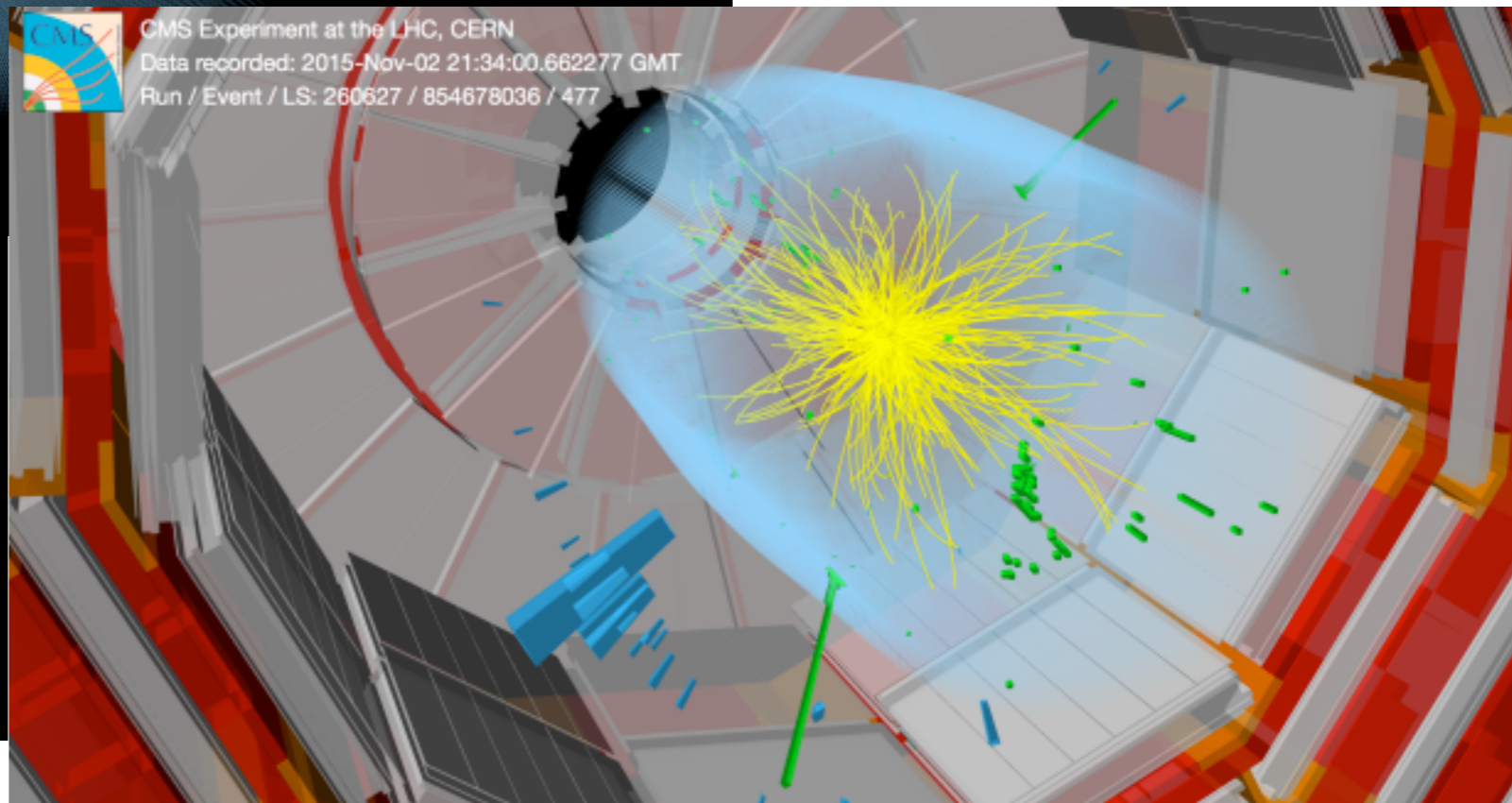
# di-photons



EM-Cal

tracker

Hadr-Cal

$X \rightarrow \gamma\gamma$

$m_{\gamma\gamma} \sim 750$ GeV
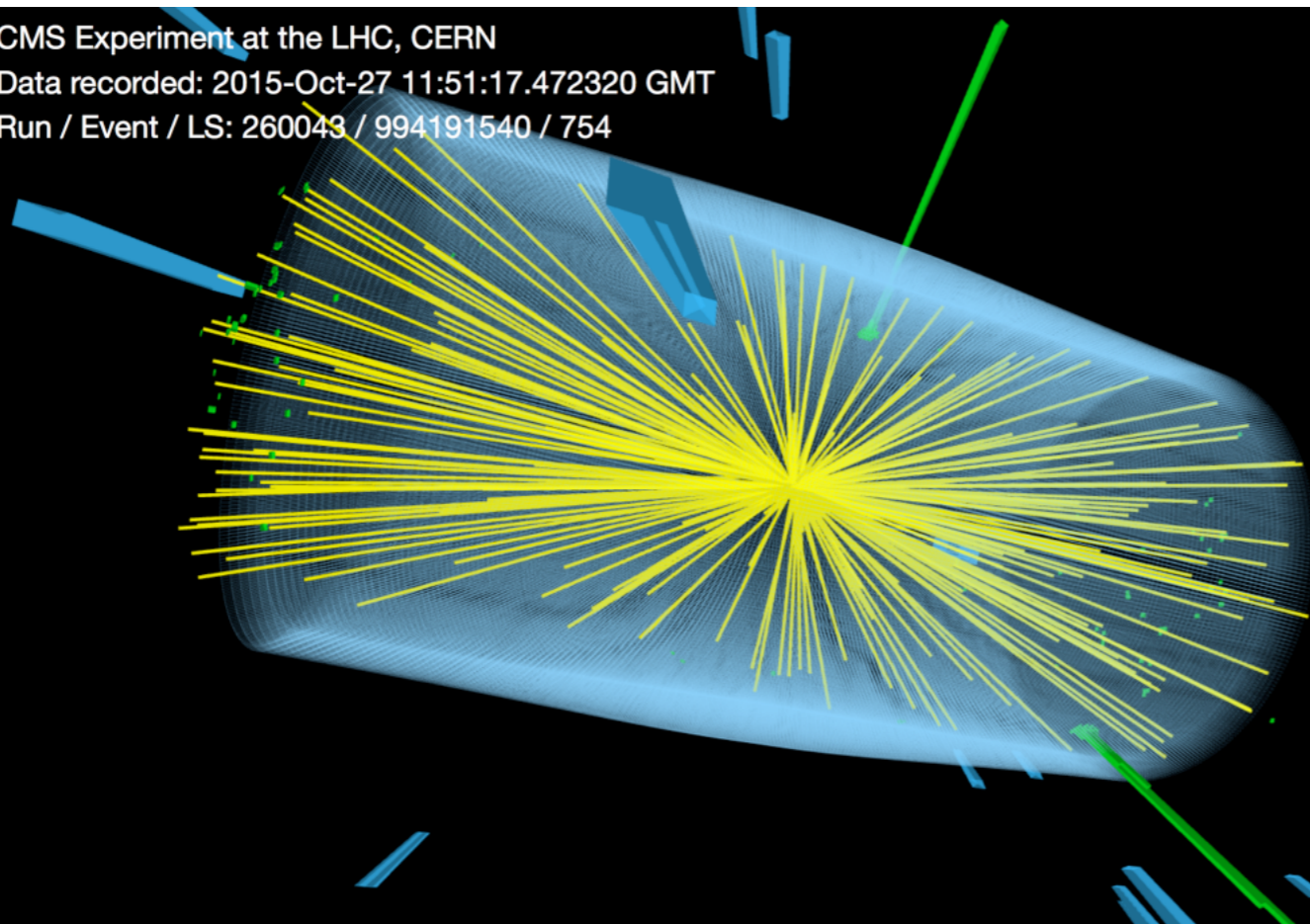
CMS-PHO-EVENTS-2015-007

CMS Experiment at the LHC, CERN
Data recorded: 2015-Nov-02 21:34:00.662277 GMT
Run / Event / LS: 260627 / 854678036 / 477

CMS Experiment at the LHC, CERN
Data recorded: 2015-Oct-27 11:51:17.472320 GMT
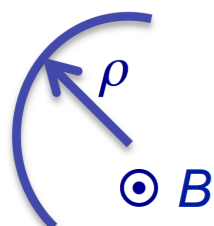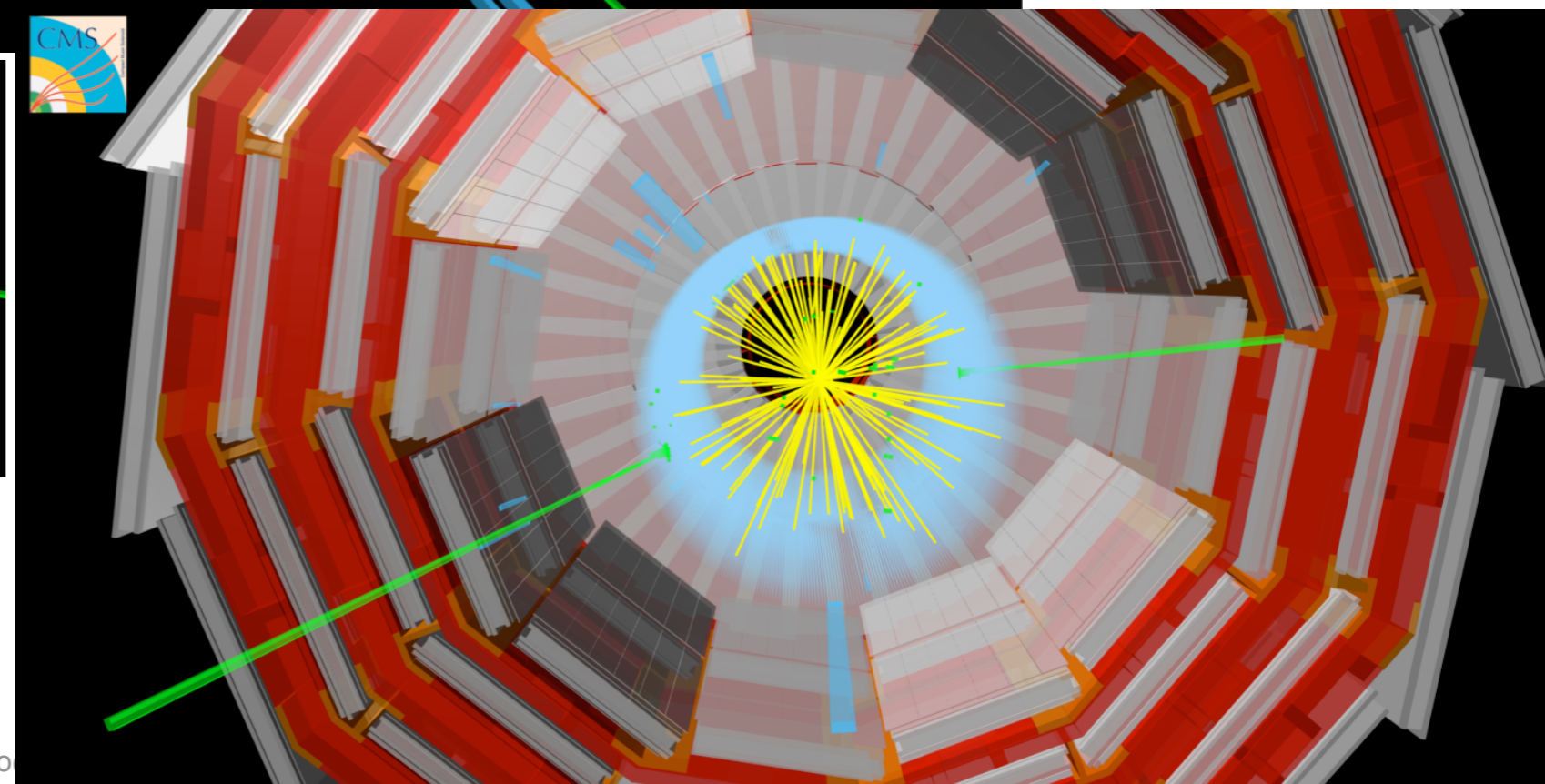Run / Event / LS: 260043 / 994191540 / 754

?

m$_{??}$ ~ 800 GeV

CMS

CMS Experiment at the LHC, CERN
Data recorded: 2015-Sep-11 22:46:54.589056 GMT
Run / Event / LS: 256353 / 437637379 / 244

$\rho$

$\odot$ $B$

$$\rho = \frac{p}{ZeB}$$

# di-jets

## X → j j



Jet 1,
  pt = 2.88 TeV
  eta = -0.364
  phi = 1.915

Jet 0,
  pt = 3.04 TeV
  eta = 0.059
  phi = -1.235

CMS Experiment at LHC, CERN
Data recorded: Mon Oct 12 2015 EEST
Run/Event: 258749 / 549864773
Lumi section: 355
Dijet Mass: 6.14 TeV

Dijet Mass = 6.14 TeV

Jet 1,
  pt = 2.88 TeV
  eta = -0.364
  phi = 1.915

Jet 0,
  pt = 3.04 TeV
  eta = 0.059
  phi = -1.235

# a di-muon event

X ➞ μμ

# a $\mu^+\mu^-e^+e^-$ event

Run 251244 Event 204117665
$\sqrt{s}$ = 13 TeV

$e_1$
$p_T$ = 63.3 GeV
$\eta$ = 1.2

$\mu_1$
$p_T$ = 58.7 GeV
$\eta$ = 1.8

$pp \rightarrow ZZ \rightarrow 2e2\mu$

$m_{\mu\mu}$= 91.1 GeV
$m_{ee}$ = 88.2 GeV
$m_{4\ell}$ = 208.9 GeV
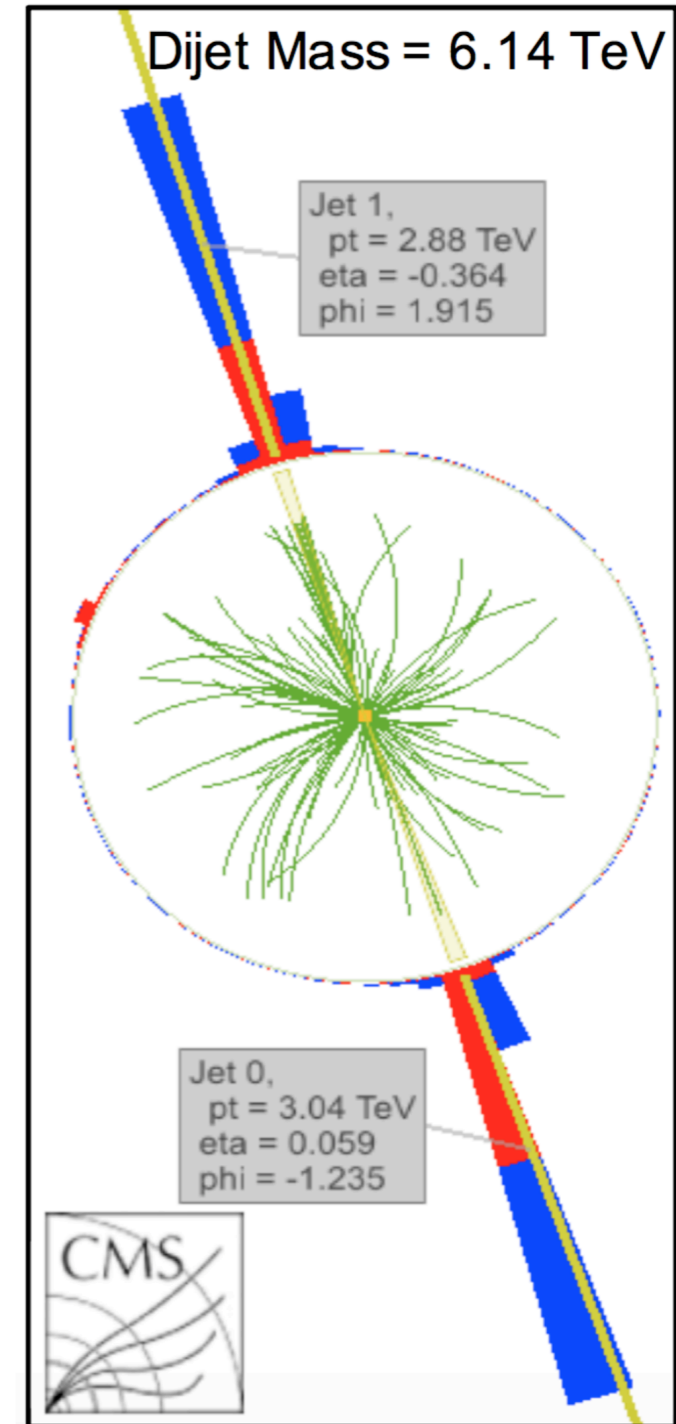
$\mu_2$
$p_T$ = 36.1 GeV
$\eta$ = 0.98

$e_2$
$p_T$ = 25.5 GeV
$\eta$ = 0.20

CMS

➡ processes are explored with many (more) final state particles

# the di-muon analysis

# the *di-muon* spectrum (X →μμ)

## 50 years of particle physics in one plot!

PRECISION

J/ψ

φ

ψ'

Y

B

ω

Z

NEW
PARTICLES

$10^5$

$10^4$

$10^3$

RARE

Direct

1

10

$10^2$

$\mu^+\mu^-$ invariant mass [GeV]

Indirect

# from detector to physics ...

**?**

# di-muon 'invariant mass' ?



particle identification
- signal in muon chambers
$\rightarrow$ it's a muon!
➠ $m = m(\mu) \sim 106 \text{MeV}/c^2$

particle trajectory
- muon chambers but especially
  the silicon tracker
➠ linear momentum, $\underline{p} \equiv (p_x, p_y, p_z)$

➠ form 4-momentum of each muon: $\mathbf{P_\mu} \equiv (E, p_x, p_y, p_z)$

➠ that of the di-muon pair $\mathbf{P_{\mu\mu}} = \mathbf{P_{\mu 1}} + \mathbf{P_{\mu 2}} = \mathbf{P_{X \rightarrow \mu\mu}}$

➠ invariant mass $\mathbf{P_{\mu\mu} \cdot P_{\mu\mu}} = \mathbf{M_{\mu\mu}}^2 = (\mathbf{M_X})^2$

# the reconstructed di-muon spectrum



feature: variable bin widths, resolution-dependent, properly normalized, doubly-log scales

# fit the data



**Dimuon Spectrum**
resonance: J/ψ
mass: 3.093 ± 0.000 GeV/c$^2$
with: 34.067 ± 0.344 MeV/c$^2$

Data
Background fit
Signal fit
Global Fit

- establish a **fit model**
  ‣ signal; Gaussian
  ‣ background: polynominal

- extract **signal parameters**
  ‣ yield ($N \pm \sigma_N$), mass ($m \pm \sigma_m$)

- estimate **systematic errors**
  ‣ does the choice of fit model affect the measured results ?
  ‣ quantify the systematic variations by employing different models

- quote **final measurements**
  ‣ $N \pm \sigma_{stat} \pm \sigma_{syst}$

- inspect quality of fit
  ‣ can model be improved?
  ‣ hint: final state radiation ($\mu \rightarrow \mu\gamma$) may distort shape

# what's the physics process ?



production: strong force  decay: electroweak force

# what are the peaks?

**Check their measured properties from: http://pdglive.lbl.gov**

## $Z$ $\quad J = 1$

*See related reviews:*

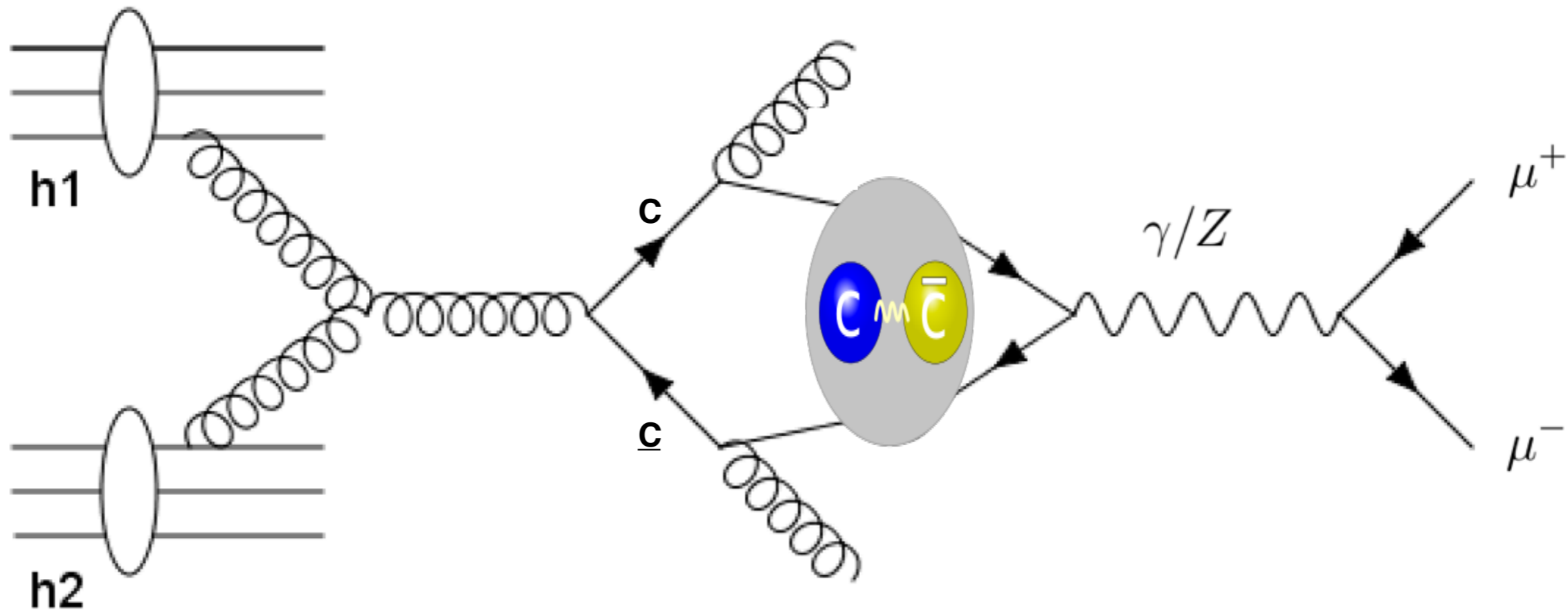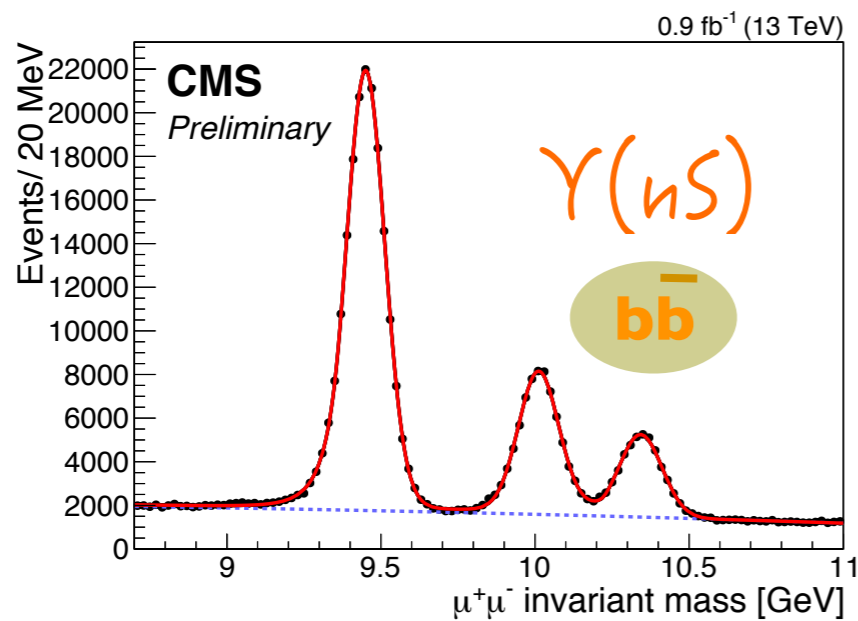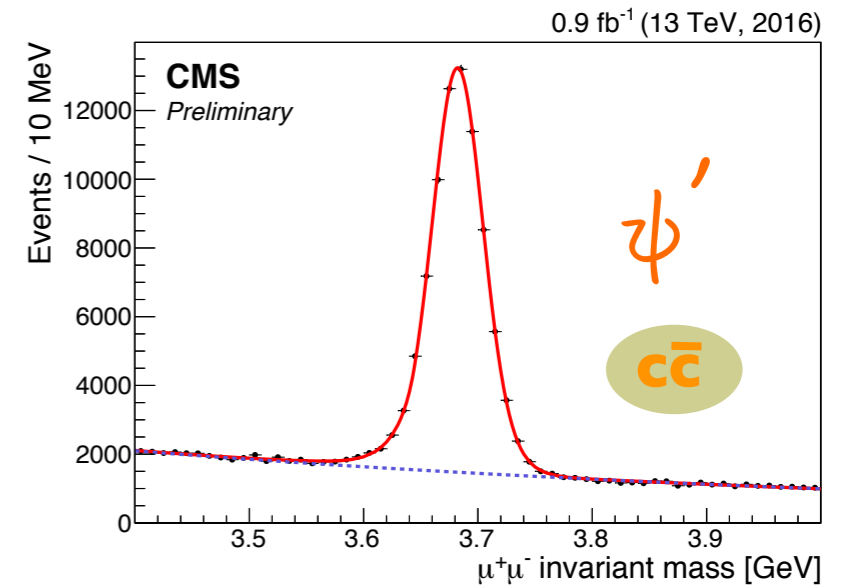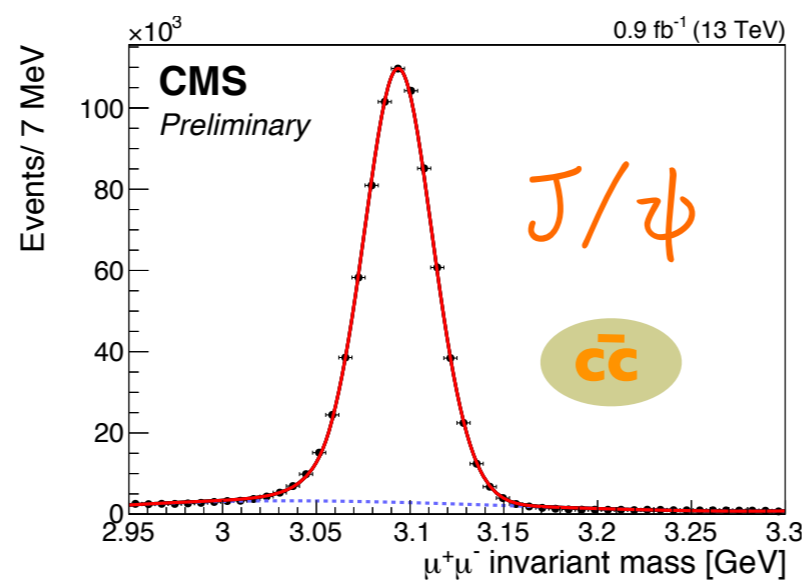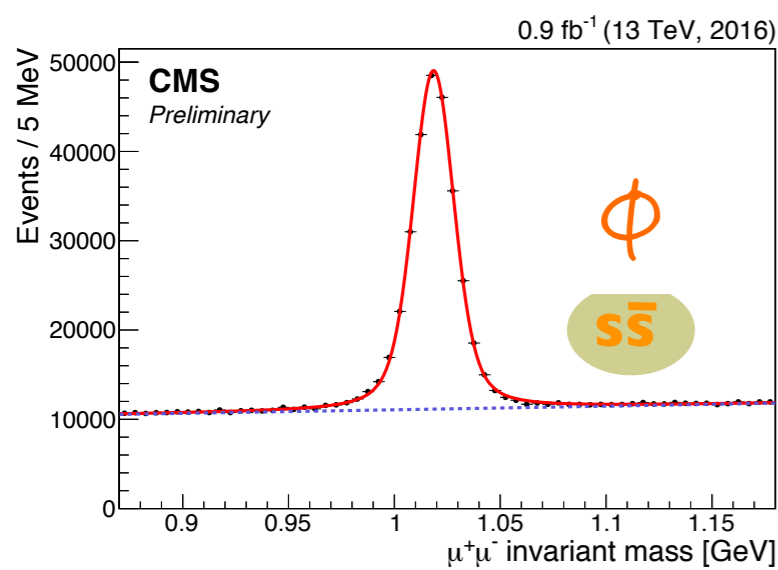$Z$ Boson — PDF

Anomalous $ZZ\gamma$, $Z\gamma\gamma$, and $ZZV$ Couplings — PDF

Anomalous $W/Z$ Quartic Couplings (QGCs) — PDF

▸ Expand all sections

| | |
|---|---|
| $Z$ MASS | $91.1876 \pm 0.0021$ GeV |

---

1      **55. $Z$ Boson**

# 55. $Z$ Boson

Revised August 2018 by M. Grünewald (University Coll. Dublin) and A. Gurtu (CERN; TIFR Mumbai).

Precision measurements at the $Z$-boson resonance using electron–positron colliding beams began in 1989 at the SLC and at LEP. During 1989–95, the four LEP experiments (ALEPH, DELPHI, L3, OPAL) made high-statistics studies of the production and decay properties of the $Z$. Although the SLD experiment at the SLC collected much lower statistics, it was able to match the precision of LEP experiments in determining the effective electroweak mixing angle $\sin^2\overline{\theta}_W$ and the rates of $Z$ decay to $b$- and $c$-quarks, owing to availability of polarized electron beams, small beam size, and stable beam spot.

The $Z$-boson properties reported in this section may broadly be categorized as:

- The standard 'lineshape' parameters of the $Z$ consisting of its mass, $M_Z$, its total width, $\Gamma_Z$, and its partial decay widths, $\Gamma(\text{hadrons})$, and $\Gamma(\ell\ell)$ where $\ell = e, \mu, \tau, \nu$;
- $Z$ asymmetries in leptonic decays and extraction of $Z$ couplings to charged and neutral leptons;
- The $b$- and $c$-quark-related partial widths and charge asymmetries which require special techniques;
- Determination of $Z$ decay modes and the search for modes that violate known conservation laws;
- Average particle multiplicities in hadronic $Z$ decay;
- $Z$ anomalous couplings.

The effective vector and axial-vector coupling constants describing the $Z$-to-fermion coupling are also measured in $p\bar{p}$ and $ep$ collisions at the Tevatron and at HERA. The corresponding cross-section formulae are given in Section 39 (Cross-section formulae for specific processes) and Section 16 (Structure Functions) in this *Review*. In this minireview, we concentrate on the measurements in $e^+e^-$ collisions at LEP and SLC.

The standard 'lineshape' parameters of the $Z$ are determined from an analysis of the production cross sections of these final states in $e^+e^-$ collisions. The $Z \to \nu\bar{\nu}(\gamma)$ state is identified directly by detecting single photon production and indirectly by subtracting the visible partial widths from the total width. Inclusion in this analysis of the forward-backward asymmetry of charged leptons, $A_{FB}^{(0,\ell)}$, of the $\tau$ polarization, $P(\tau)$, and its forward-backward asymmetry, $P(\tau)^{fb}$, enables the separate determination of the effective vector ($\overline{g}_V$) and axial vector ($\overline{g}_A$) couplings of the $Z$ to these leptons and the ratio ($\overline{g}_V/\overline{g}_A$), which is related to the effective electroweak mixing angle

---

## $c$ $\quad I(J^P) = 0(1/2^+)$

$$\text{Charge} = \tfrac{2}{3}\, e \qquad \text{Charm} = +1$$

| | |
|---|---|
| $c$-QUARK MASS | $1.27 \pm 0.02$ GeV |
| $m_c/m_s$ MASS RATIO | $11.76^{+0.05}_{-0.10}$ |
| $m_b/m_c$ MASS RATIO | $4.58 \pm 0.01$ |
| $m_b - m_c$ QUARK MASS DIFFERENCE | $3.45 \pm 0.05$ GeV |

---

## $J/\psi(1S)$ $\quad I^G(J^{PC}) = 0^-(1^{--})$

| | |
|---|---|
| $J/\psi(1S)$ MASS | $3096.900 \pm 0.006$ MeV |
| $J/\psi(1S)$ WIDTH | $92.6 \pm 1.7$ keV (S = 1.1) |

### $J/\psi(1S)$ Decay Modes

▸ Expand all decays

| | Mode | | Fraction ($\Gamma_i/\Gamma$) | Scale Factor/ Conf. Level | P(MeV/c) |
|---|---|---|---|---|---|
| $\Gamma_1$ | hadrons | | $(87.7 \pm 0.5)\%$ | | |
| $\Gamma_2$ | virtual $\gamma \to$ hadrons | | $(13.50 \pm 0.30)\%$ | | |
| $\Gamma_3$ | $ggg$ | | $(64.1 \pm 1.0)\%$ | | |
| $\Gamma_4$ | $\gamma gg$ | | $(8.8 \pm 1.1)\%$ | | |
| $\Gamma_5$ | $e^+e^-$ | | $(5.971 \pm 0.032)\%$ | | 1548 |
| $\Gamma_6$ | $e^+e^-\gamma$ | [1] | $(8.8 \pm 1.4) \times 10^{-3}$ | | 1548 |
| $\Gamma_7$ | $\mu^+\mu^-$ | | $(5.961 \pm 0.033)\%$ | | 1545 |

▸ Decays involving hadronic resonances

▸ Decays into stable hadrons

▸ Radiative decays

▸ Dalitz decays

▸ Weak decays

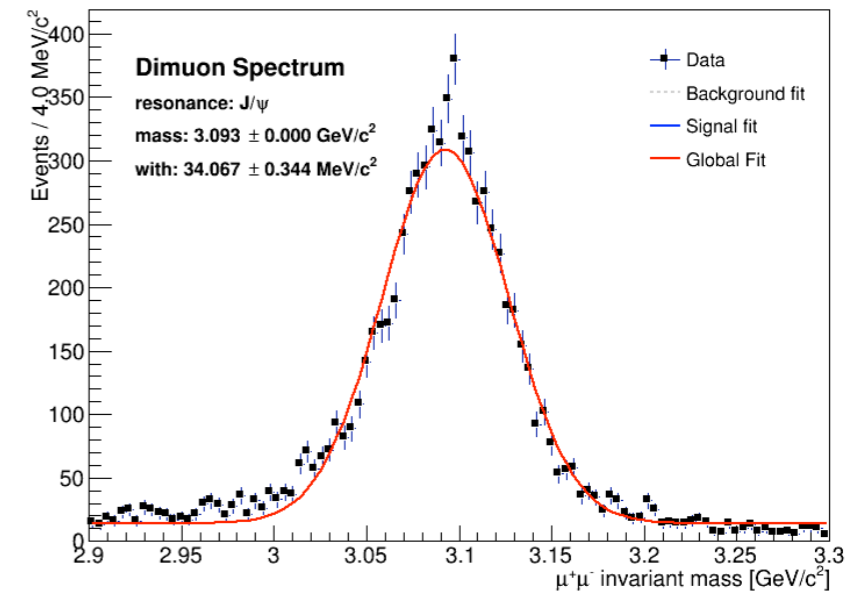▸ Charge conjugation ($C$), Parity ($P$), Lepton Family number ($LF$) violating modes

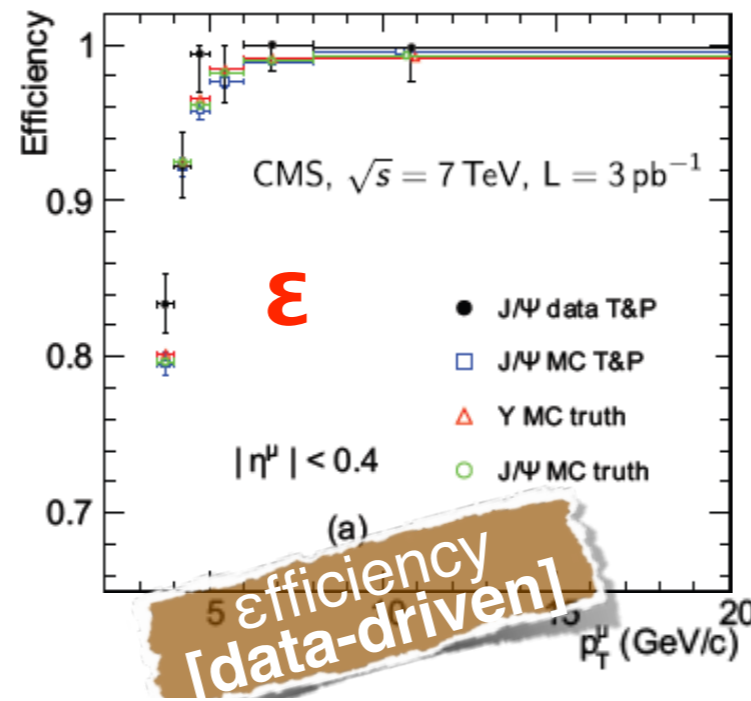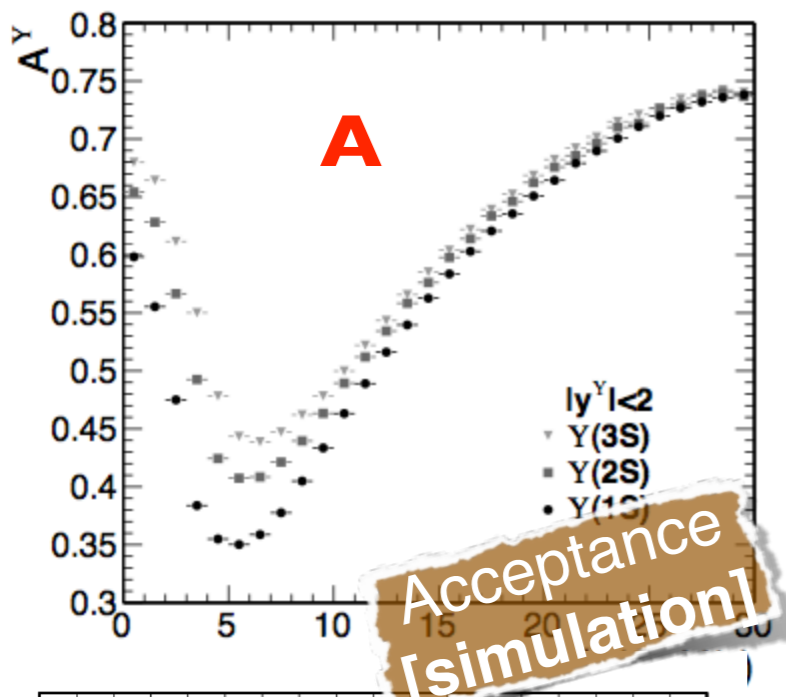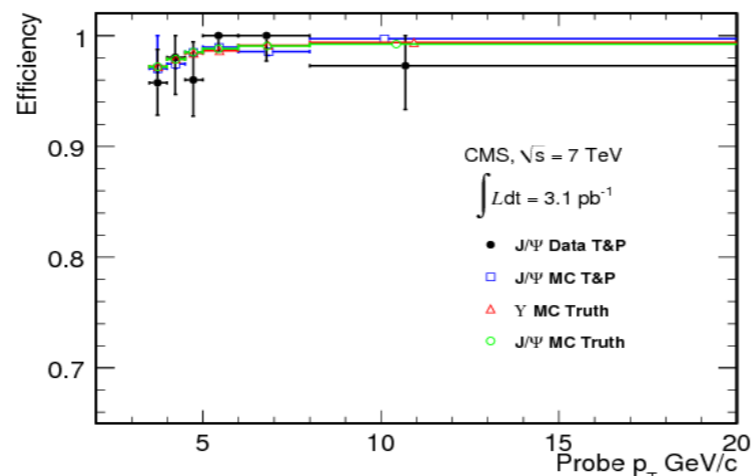▸ Other decays

# Cross section

an effective area of interaction
unit: barn, 1b = $10^{-28}$ m$^2$ = 100fm$^2$

$$\frac{d^2\sigma(Q\overline{Q})}{dp_T\,dy}\mathcal{B}\left(Q\overline{Q}\to\mu^+\mu^-\right) = \frac{N_{fit}(Q\overline{Q})}{\mathcal{L}\cdot\mathcal{A}\cdot\epsilon\cdot\Delta p_T\cdot\Delta y}$$

**A**

**Acceptance [simulation]**

**ε**

CMS, $\sqrt{s}$ = 7 TeV, L = 3 pb$^{-1}$

- J/Ψ data T&P
- J/Ψ MC T&P
- Υ MC truth
- J/Ψ MC truth

$|\eta^\mu| < 0.4$

(a)

**εfficiency [data-driven]**

**Dimuon Spectrum**
resonance: J/ψ
mass: 3.093 ± 0.000 GeV/c$^2$
with: 34.067 ± 0.344 MeV/c$^2$

- Data
- Background fit
- Signal fit
- Global Fit

CMS, $\sqrt{s}$ = 7 TeV
$\int L\,dt$ = 3.1 pb$^{-1}$

- J/Ψ Data T&P
- J/Ψ MC T&P
- Υ MC Truth
- J/Ψ MC Truth

- **N**: fitted signal yield

- **A**: detector acceptance from simulation

- **ε**: detector reconstruction and trigger efficiencies (simulation or data-driven)
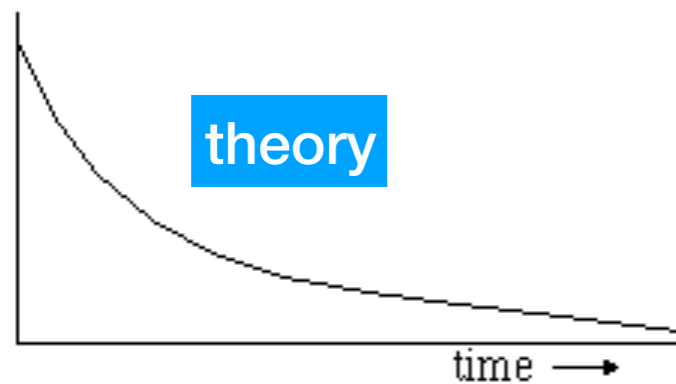
- **L**: integrated sample luminosity
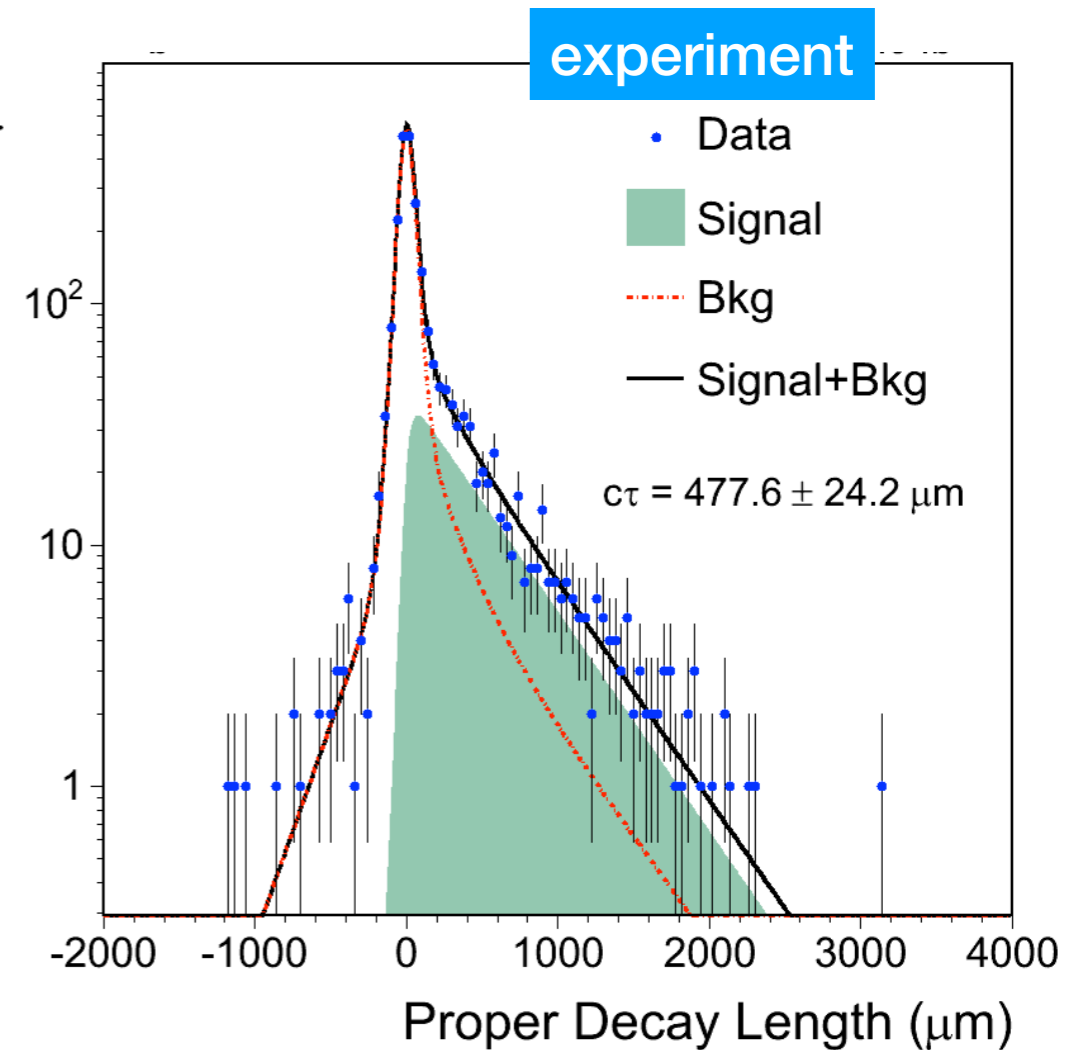
# (extra) statistics
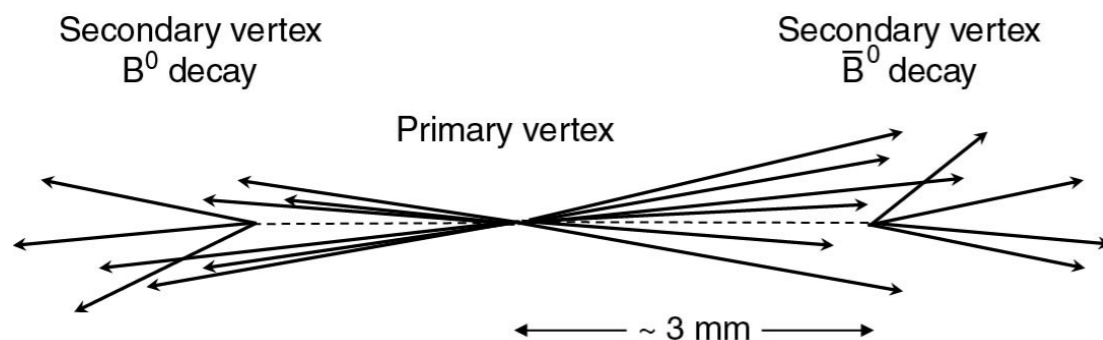
# measurement: a lifetime example

$\rightarrow$ $L(t|\sigma_t, \tau) = \dfrac{1}{\mathcal{N}} \cdot [\ \dfrac{1}{\tau}e^{-\frac{t}{\tau}}\theta(t) \otimes G(t; \sigma_t)\ ] \cdot \mathcal{E}(t)$ $+\ L(\text{Background})$

PDF normalization

theory model

t-resolution function

t-acceptance function

theory

time →

experiment

• Data
▨ Signal
-·-· Bkg
— Signal+Bkg

$c\tau = 477.6 \pm 24.2\ \mu m$

$$t = \frac{L}{\beta\gamma} = L\ \frac{M}{p} = L_{xy}\frac{M}{p_T}$$

↳ Lorentz boost factor

Secondary vertex
$B^0$ decay

Secondary vertex
$\bar{B}^0$ decay

Primary vertex

← ~ 3 mm →

Proper Decay Length (μm)

# statistics

- consider lifetime (or decay rate) distribution $\boxed{f(t; \tau) = \frac{1}{\tau} e^{-t/\tau}}$

- suppose we have n data points (measurements) $\quad t_1, \ldots, t_n$

- likelihood function $\quad \boxed{L(\tau) = \prod_{i=1}^{n} \frac{1}{\tau} e^{-t_i/\tau}}$ **maximum likelihood** method

- value of τ for which L(τ) is maximum = maximizes log-likelihood

$$\ln L(\tau) = \sum_{i=1}^{n} \ln f(t_i; \tau) = \sum_{i=1}^{n} \left( \ln \frac{1}{\tau} - \frac{t_i}{\tau} \right)$$

- can find its maximum as $\quad \dfrac{\partial \ln L(\tau)}{\partial \tau} = 0 \quad \rightarrow \quad \hat{\tau} = \dfrac{1}{n} \sum_{i=1}^{n} t_i$

- mean: $E[t] = \displaystyle\int_0^{\infty} t \frac{1}{\tau} e^{-t/\tau} \, dt = \tau \quad$ variance: $V[t] = \displaystyle\int_0^{\infty} (t - \tau)^2 \frac{1}{\tau} e^{-t/\tau} \, dt = \tau^2$

- for ML estimator $\quad E[\hat{\tau}] = E\left[ \dfrac{1}{n} \sum_{i=1}^{n} t_i \right] = \dfrac{1}{n} \sum_{i=1}^{n} E[t_i] = \tau \quad \longrightarrow \quad b = E[\hat{\tau}] - \tau = 0$

$$V[\hat{\tau}] = V\left[ \frac{1}{n} \sum_{i=1}^{n} t_i \right] = \frac{1}{n^2} \sum_{i=1}^{n} V[t_i] = \frac{\tau^2}{n} \quad \longrightarrow \quad \sigma_{\hat{\tau}} = \frac{\tau}{\sqrt{n}}$$

# maximum likelihood

- expand $\ln L\,(\theta)$ about its maximum

$$\ln L(\theta) = \ln L(\hat{\theta}) + \left[\frac{\partial \ln L}{\partial \theta}\right]_{\theta=\hat{\theta}} (\theta-\hat{\theta}) + \frac{1}{2!}\left[\frac{\partial^2 \ln L}{\partial \theta^2}\right]_{\theta=\hat{\theta}} (\theta-\hat{\theta})^2 + ..$$

- first term is $\ln L_{max}$, second term is zero, third term approximate

from information inequality:

$$\ln L(\theta) \approx \ln L_{max} - \frac{(\theta-\hat{\theta})^2}{2\widehat{\sigma^2}_{\hat{\theta}}} \qquad \ln L(\hat{\theta} \pm \hat{\sigma}_{\hat{\theta}}) \approx \ln L_{max} - \frac{1}{2} \qquad \hat{V}[\hat{\theta}] = -\left(\frac{\partial^2 \ln L}{\partial \theta^2}\right)^{-1}\Bigg|_{\theta=\hat{\theta}}$$

- error estimate: change $\theta$ away from $\hat{\theta}$ until ln L decreases by 1/2

- 1σ (68.3% CL) confidence interval

$$\ln L(\theta) = \ln L(\hat{\theta}) - \frac{1}{2} \quad \Rightarrow \quad [\hat{\theta} - \sigma_{\hat{\theta}},\ \hat{\theta} + \sigma_{\hat{\theta}}]$$

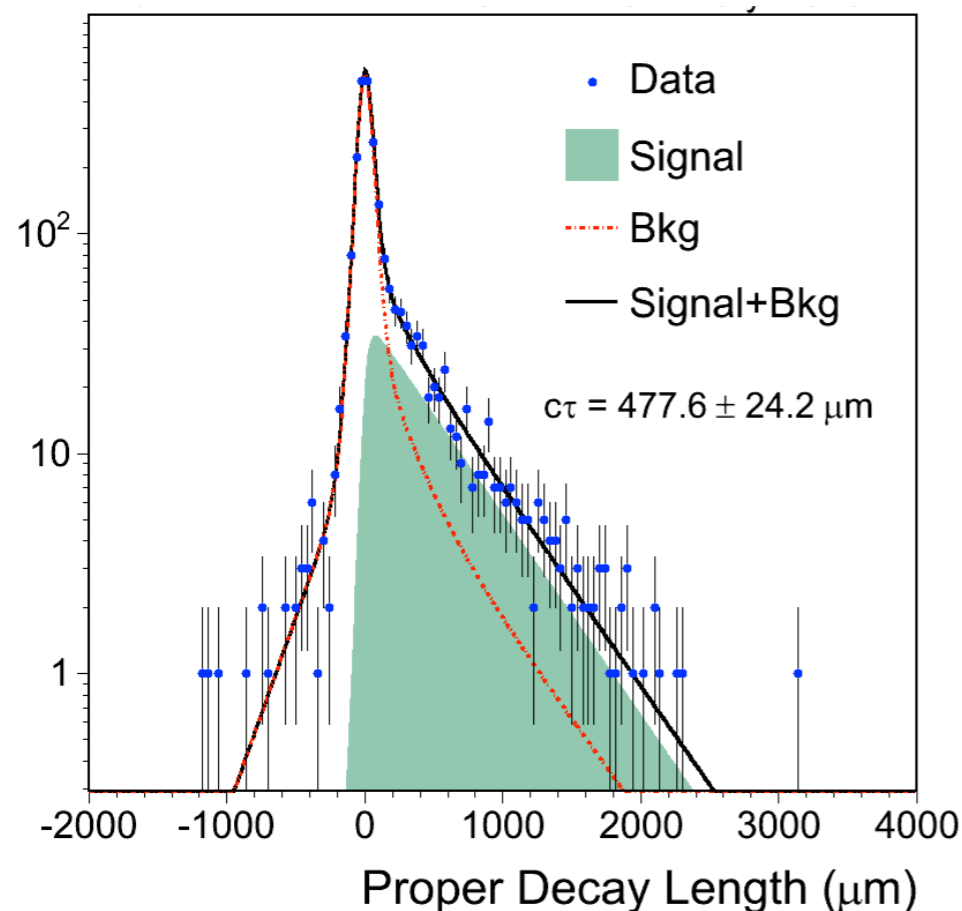$$\hat{\tau} = 0.85^{+0.52}_{-0.30}$$

# systematic uncertainties

- **statistical** uncertainty on τ is obtained from the maximum likelihood fit to the data

$$L(t|\sigma_t, \tau) \;=\; \frac{1}{\mathcal{N}} \cdot \left[\; \frac{1}{\tau} e^{-\frac{t}{\tau}} \theta(t) \otimes G(t;\sigma_t) \;\right] \cdot \mathcal{E}(t) \; + \; \text{L}(\text{Background})$$

- **systematic** uncertainty quantifies any uncertainty in the procedure going from the raw data to a published result
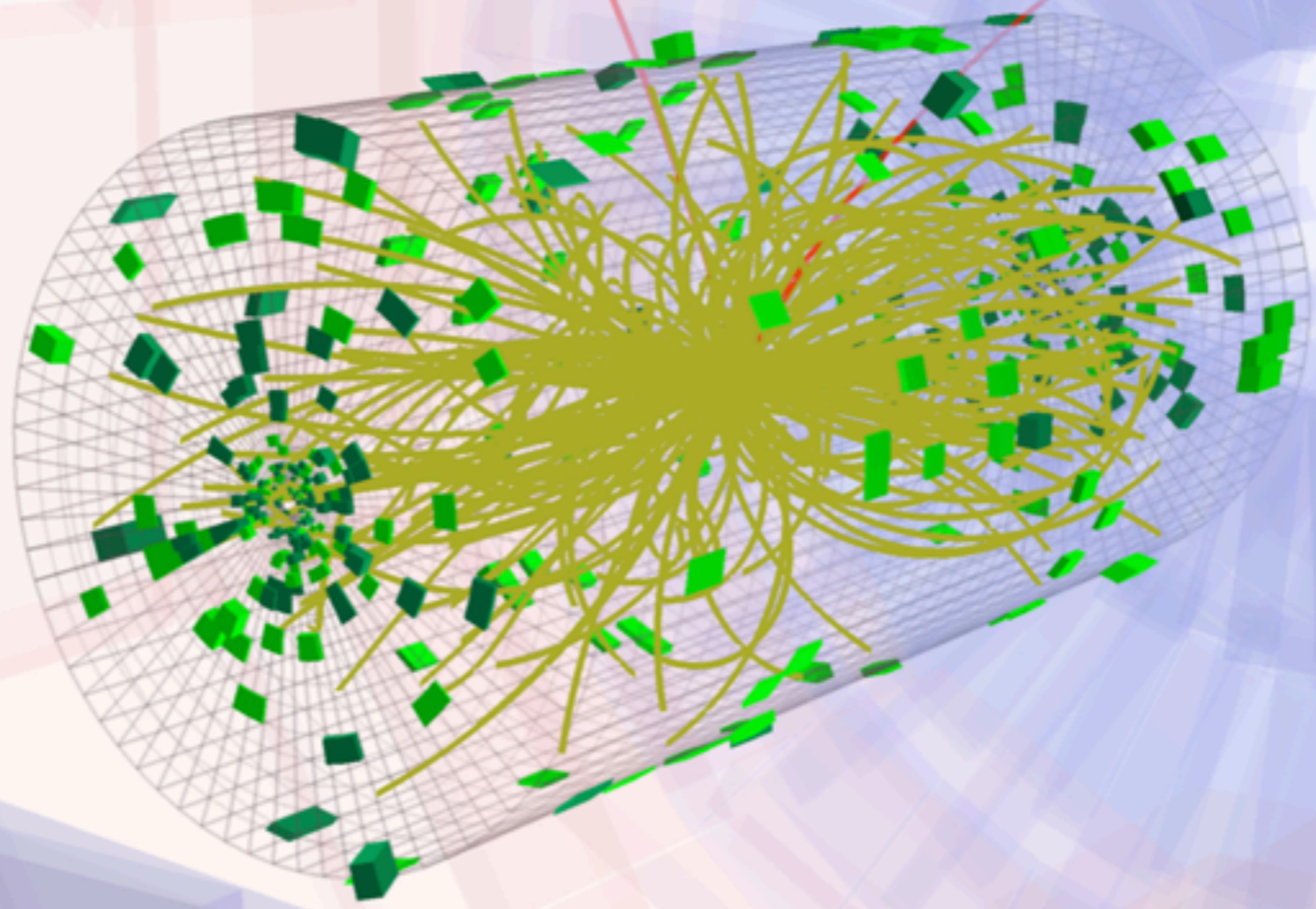


cτ = 477.6 ± 24.2 μm

**Sources of systematic error:**
- resolution calibration $\sigma_t$
- resolution function, $G(t)$
- t-efficiency $\epsilon(t)$
- background model
- …

*Inputs from Simulation and Data-driven*

$$c\tau = 477.6 \pm 24.2\,(stat.) \pm 17.6\,(syst.)\ \mu m$$

# (extra) ingredients of a physics measurement

B_s → μμ

the 'golden' rare decay

# searching for an *ultra-rare* decay: B→μμ

## 1. ONLINE SELECTION (TRIGGER)



**trigger paths**
- $\psi'$
- $J/\psi$
- $B_s$
- $\Upsilon$
- low $p_T$ double muon
- high $p_T$ double muon

**Dimuon Trigger**

- L1 Hardware Trigger
  - $p_T > 3$ GeV (few kHz)
- HLT Full tracking and vertexing
- HLT $B_s \to \mu\mu$
  - Leading and sub-leading μ $p_T > 3,4$ (4,4) GeV $|\eta_{\mu\mu}| < 1.8$ (1.8 $< |\eta_{\mu\mu}| < 2.2$)
  - $p_T (\mu\mu) > 5$ (4.8-6) GeV
  - 4.8 $< m(\mu\mu) < 6.0$ GeV
  - $P(\chi^2/dof) > 0.5\%$

# *searching for an ultra-rare decay:* B→μμ

## 1. online selection (trigger)

## 2. blind the data (avoid bias)

**SIGNAL REGION BLINDED**



analysis procedure and event selection
developed without inspecting the data
in region where signal is expected

"box opening" only later,
at final analysis stages

# *searching for an ultra-rare decay:* B→μμ

1. ONLINE SELECTION (TRIGGER)

2. BLIND THE DATA (AVOID BIAS)
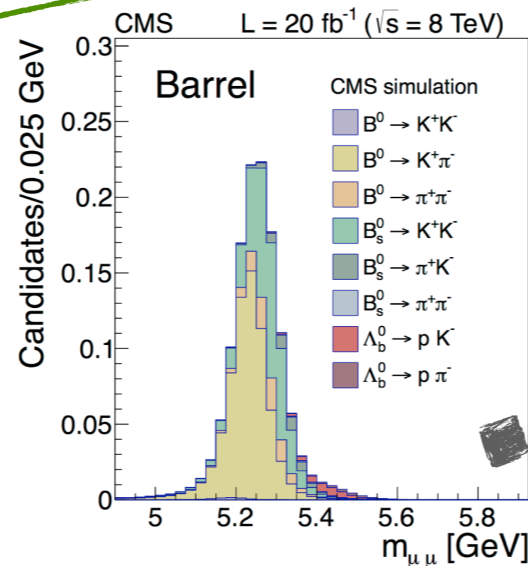
3. MULTIVARIATE SELECTION

# searching for an *ultra-rare* decay: B→µµ

1. ONLINE SELECTION (TRIGGER)

2. BLIND THE DATA (AVOID BIAS)

3. MULTIVARIATE SELECTION

4. FIT THE DATA (LIKELIHOOD)

Fit the data accounting for the various signal and background components
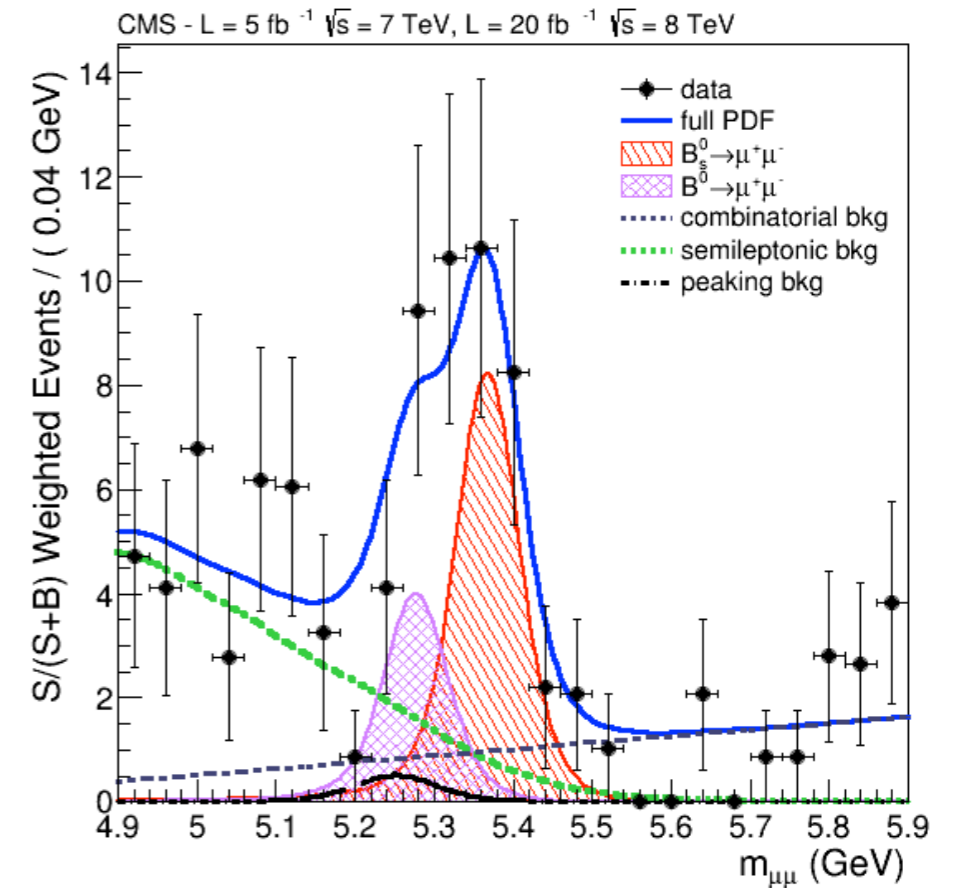
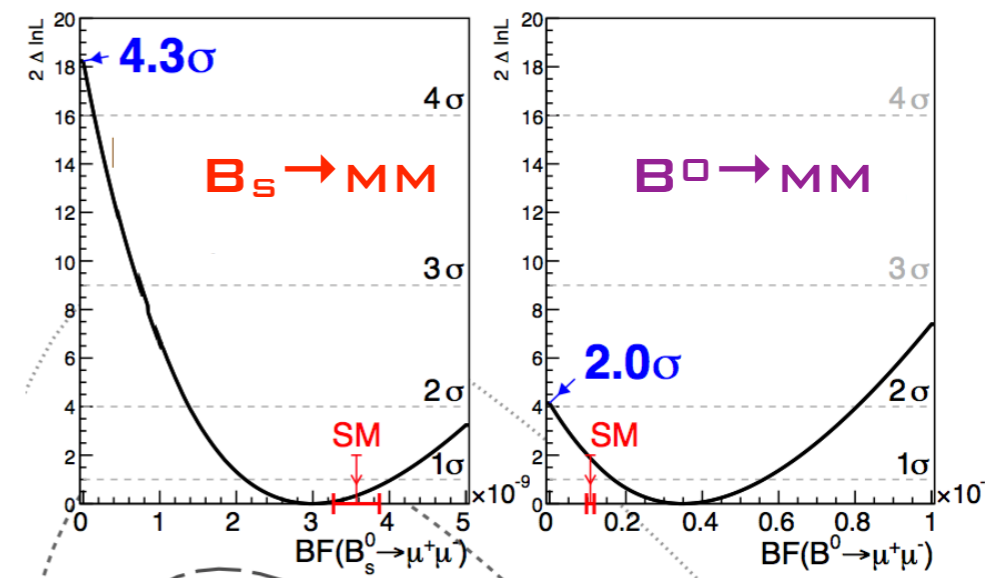SEMILEPTONIC BKG

COMBINATORIAL BKG

PEAKING BKG

SIGNAL 1: $B_s$→MM

SIGNAL 2: $B^0$→MM

# searching for an *ultra-rare* decay: B→μμ

1. **ONLINE SELECTION (TRIGGER)**

2. **BLIND THE DATA (AVOID BIAS)**

3. **MULTIVARIATE SELECTION**
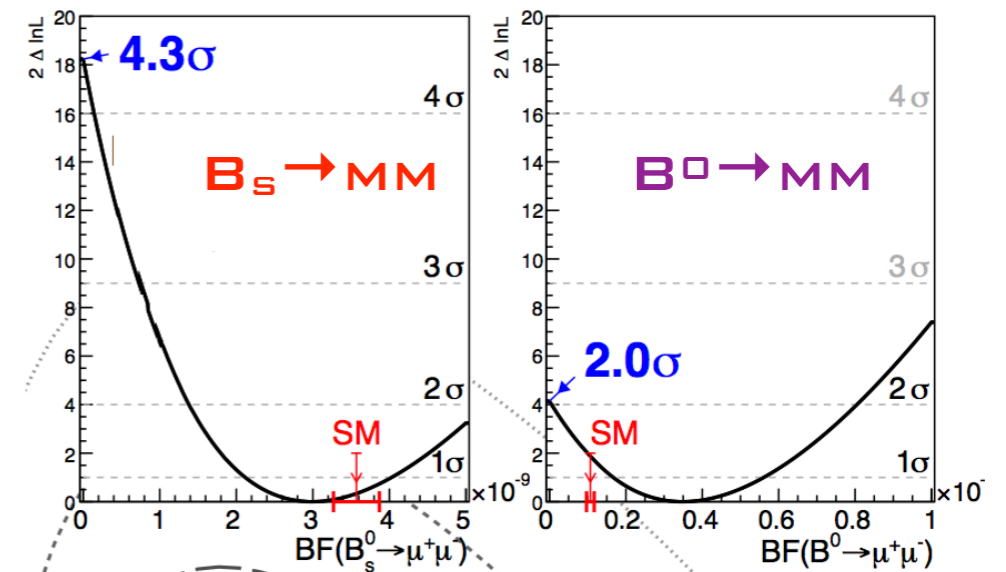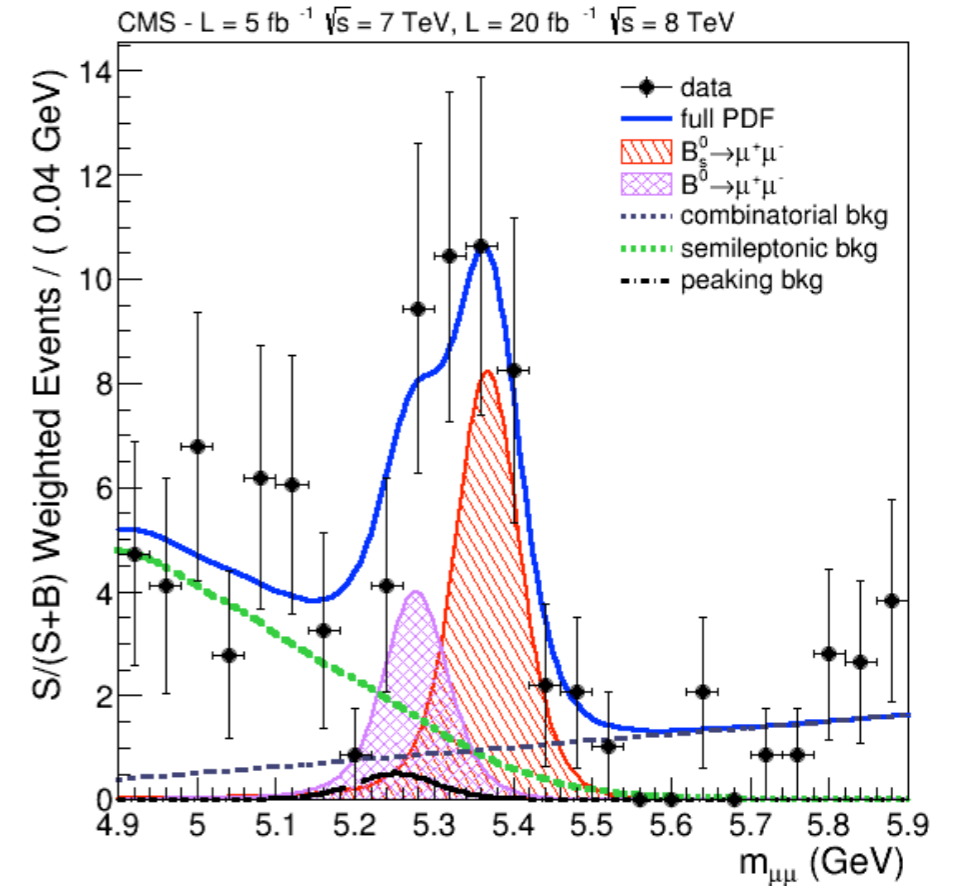
4. **FIT THE DATA (LIKELIHOOD)**

5. **STATISTICAL SIGNIFICANCE**

is the observed excess a genuine signal, or just a fluctuation of the background?
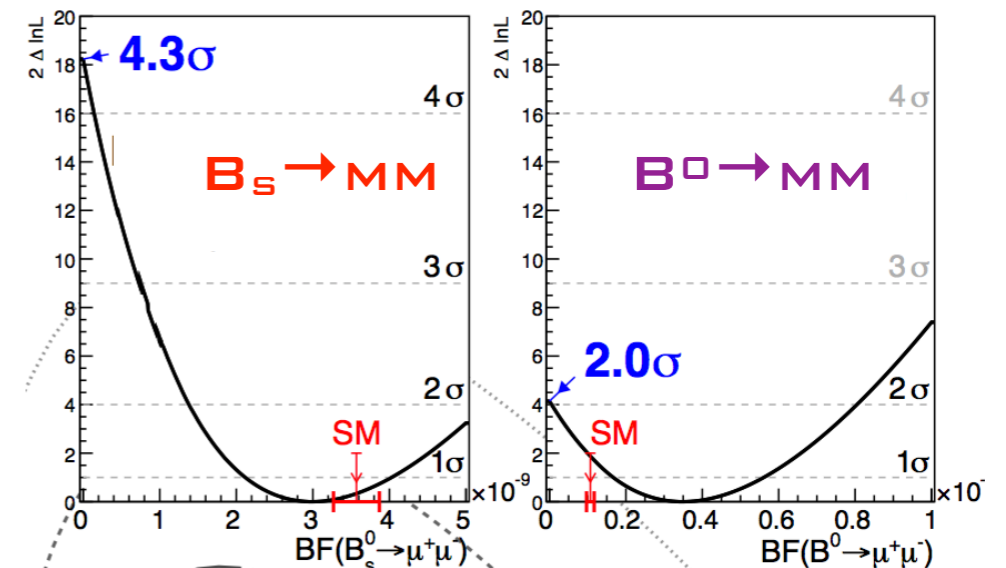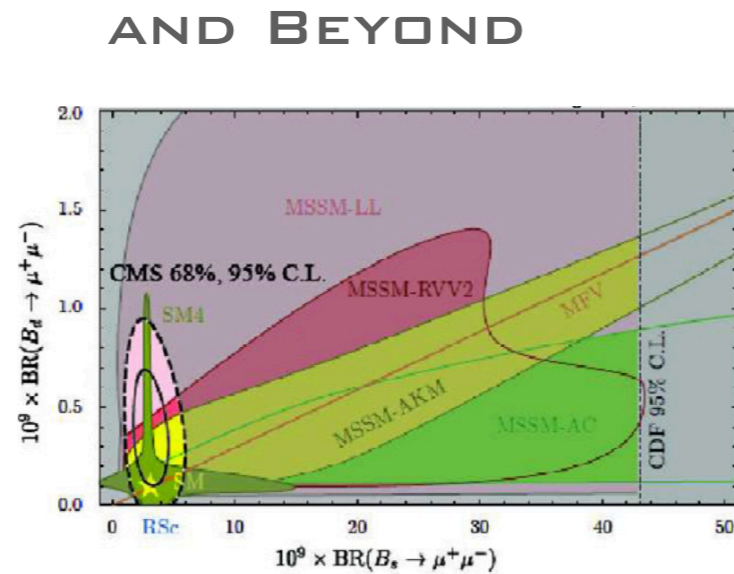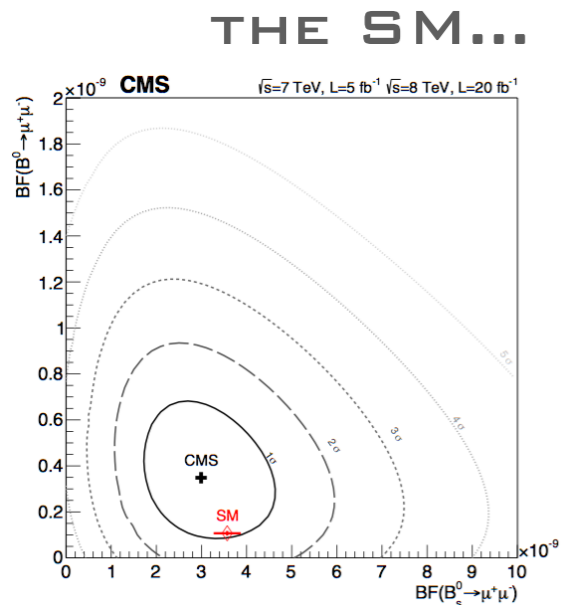
# searching for an *ultra-rare* decay: B→μμ

1. ONLINE SELECTION (TRIGGER)

2. BLIND THE DATA (AVOID BIAS)

3. MULTIVARIATE SELECTION

4. FIT THE DATA (LIKELIHOOD)

5. STATISTICAL SIGNIFICANCE

6. EXTRACT MEASUREMENT



$$BR(B_S \to \mu\mu) = \left(3.0^{+0.9}_{-0.8} \text{ (stat)}^{+0.6}_{-0.4} \text{ (syst)}\right) \times 10^{-9}$$
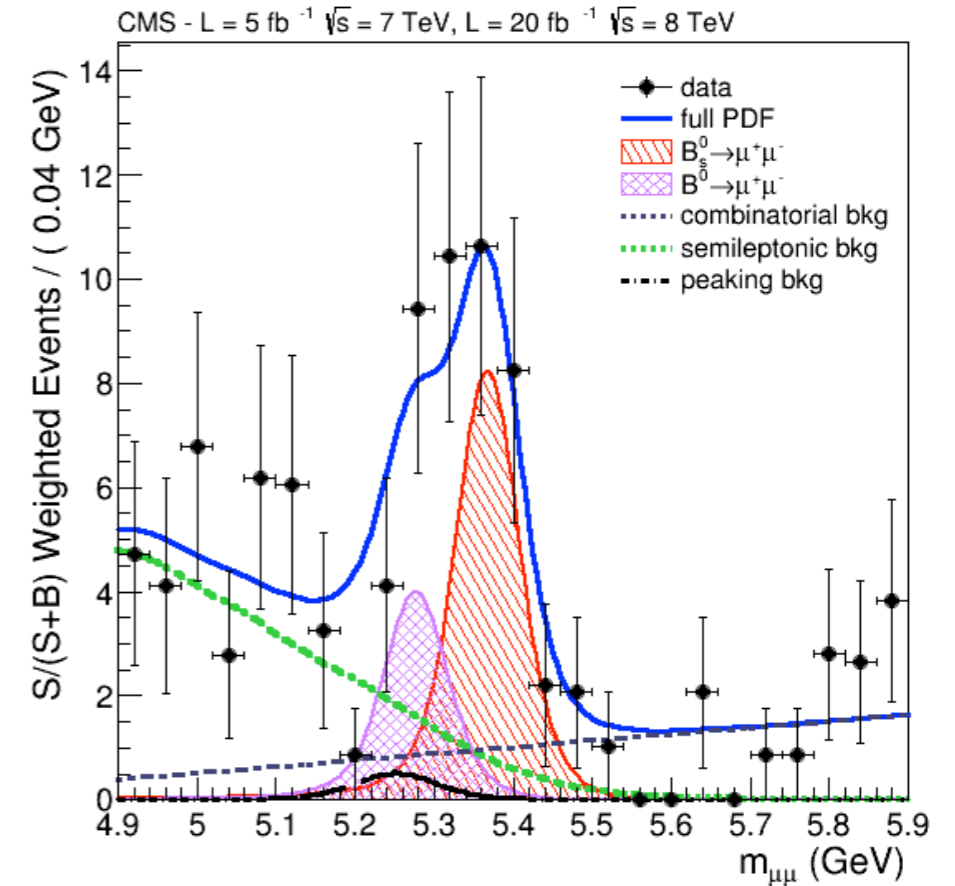
# *searching for an* *ultra-rare* *decay:* **B→μμ**

1. ONLINE SELECTION (TRIGGER)

2. BLIND THE DATA (AVOID BIAS)

3. MULTIVARIATE SELECTION

4. FIT THE DATA (LIKELIHOOD)

5. STATISTICAL SIGNIFICANCE

6. EXTRACT MEASUREMENT

7. COMPARE TO THEORY



THE SM... AND BEYOND





$$BR(B_S \rightarrow \mu\mu) = \left(3.0^{+0.9}_{-0.8}\,(\text{stat})^{+0.6}_{-0.4}\,(\text{syst})\right) \times 10^{-9}$$