# Data Analysis

**Ivica Puljak**
University of Split, FESB, Split, Croatia

Ivica.Puljak@cern.ch

---

## Starting the new era

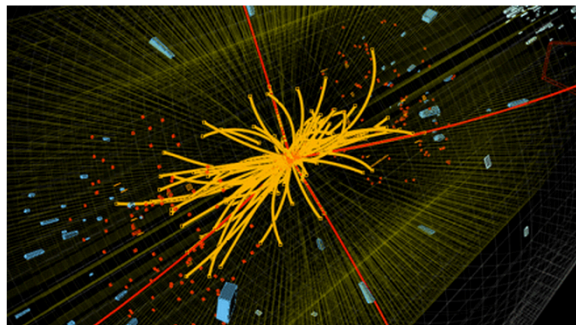- In the future, calendar of particle physics will be

*Before Higgs (BH)*      *After Higgs (AH)*
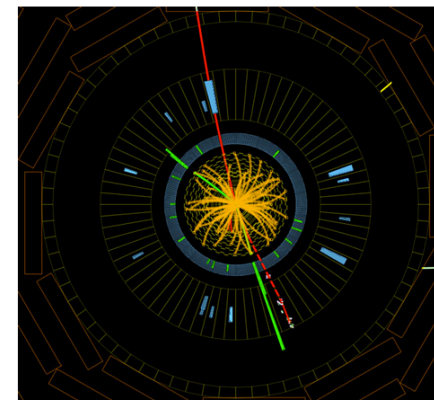
**July 4, 2012**

---

## 04.07.2012: Higgs within reach



Proton-proton collision in the CMS experiment producing four high-energy muons (red lines). The event shows characteristics expected from the decay of a Higgs boson but it is also consistent with background Standard Model physics processes (Image: CMS)

At a seminar on 4 July, the ATLAS and CMS experiments at CERN presented their latest results in the search for the long-sought Higgs boson. Both experiments see strong indications for the presence of a new particle, which could be the Higgs boson, in the mass region around 126 gigaelectronvolts (GeV).

---

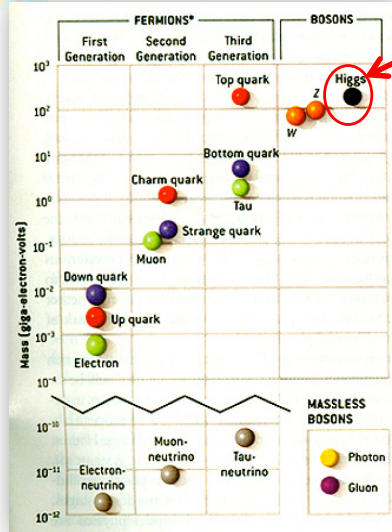## 01.08.2012: ATLAS and CMS submit Higgs-search papers



Protons collide in the CMS detector at 8 TeV, forming Z bosons which decay into electrons (green lines) and muons (red). Such an event is compatible with the decay of a Standard Model Higgs boson (Image: CMS)

The ATLAS and CMS collaborations today submitted papers to the journal *Physics Letters B* outlining the latest on their searches for the Higgs boson. The teams report even stronger evidence for the presence of a new Higgs-like particle than announced on 4 July.

# The Standard Model

**1 Missing piece: Higgs**



**Confirmed to better than 1 % uncertainty by 100's of precision measurements**

| | Measurement | Fit | $|O^{meas}-O^{fit}|/\sigma^{meas}$ |
|---|---|---|---|
| | | | 0   1   2   3 |
| $\Delta\alpha_{had}^{(5)}(m_Z)$ | $0.02758 \pm 0.00035$ | 0.02768 | |
| $m_Z$ [GeV] | $91.1875 \pm 0.0021$ | 91.1874 | |
| $\Gamma_Z$ [GeV] | $2.4952 \pm 0.0023$ | 2.4959 | |
| $\sigma_{had}^0$ [nb] | $41.540 \pm 0.037$ | 41.479 | |
| $R_l$ | $20.767 \pm 0.025$ | 20.742 | |
| $A_{fb}^{0,l}$ | $0.01714 \pm 0.00095$ | 0.01645 | |
| $A_l(P_\tau)$ | $0.1465 \pm 0.0032$ | 0.1481 | |
| $A_c$ | $0.670 \pm 0.027$ | 0.668 | |
| $A_l$(SLD) | $0.1513 \pm 0.0021$ | 0.1481 | |
| $\sin^2\theta_{eff}^{lept}(Q_{fb})$ | $0.2324 \pm 0.0012$ | 0.2314 | |
| $m_W$ [GeV] | $80.399 \pm 0.023$ | 80.379 | |
| $\Gamma_W$ [GeV] | $2.085 \pm 0.042$ | 2.092 | |
| $m_t$ [GeV] | $173.3 \pm 1.1$ | 173.4 | |

July 2010

0   1   2   3

5

---

# Higgs mass: theoretical constraints

➢ Problem: Higgs mass is free parameter

$$M_H^2 = 2\lambda v^2 \quad \ldots\ldots \quad v = 246\,\text{GeV}$$

➢ Theoretical constraints

▪ **Unitarity** (no probabilities > 1)

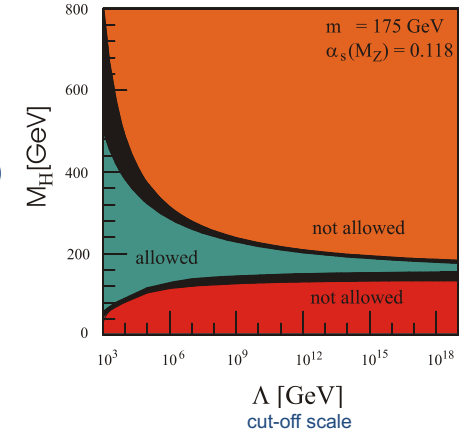$$M_H < 700 - 800\,\text{GeV}$$

▪ **Triviality**
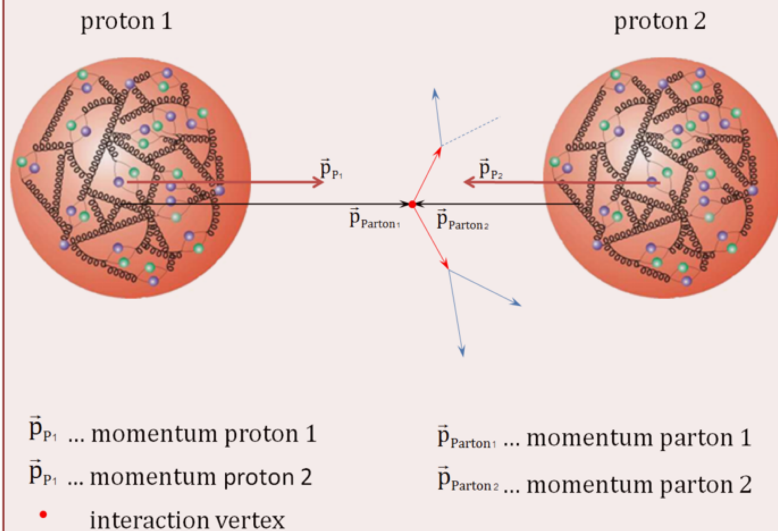(Higgs self coupling remains finite)

$$M_H^2 < \frac{4\pi v^2}{3\ln(\Lambda/v)}$$

▪ **Stability** (of vacuum)

$$M_H^2 > \frac{4m_Z^4}{\pi^2 v^2}\ln(\Lambda/v)$$



$m = 175$ GeV
$\alpha_s(M_Z) = 0.118$

not allowed

allowed

not allowed

$\Lambda$ [GeV]
cut-off scale

6

---



Interactions of constituents of the colliding protons, the so called partons (quarks, gluons)

proton 1

proton 2

$\vec{p}_{P_1}$    $\vec{p}_{P_2}$

$\vec{p}_{Parton_1}$    $\vec{p}_{Parton_2}$

$\vec{p}_{P_1}$ ... momentum proton 1

$\vec{p}_{P_1}$ ... momentum proton 2

$\vec{p}_{Parton_1}$ ... momentum parton 1

$\vec{p}_{Parton_2}$ ... momentum parton 2

• interaction vertex

---

# Collisions in LHC



**Bunches**

**Proton**

**Partons**
(quark, gluon)

**Particles**

Higgs
SUSY.....

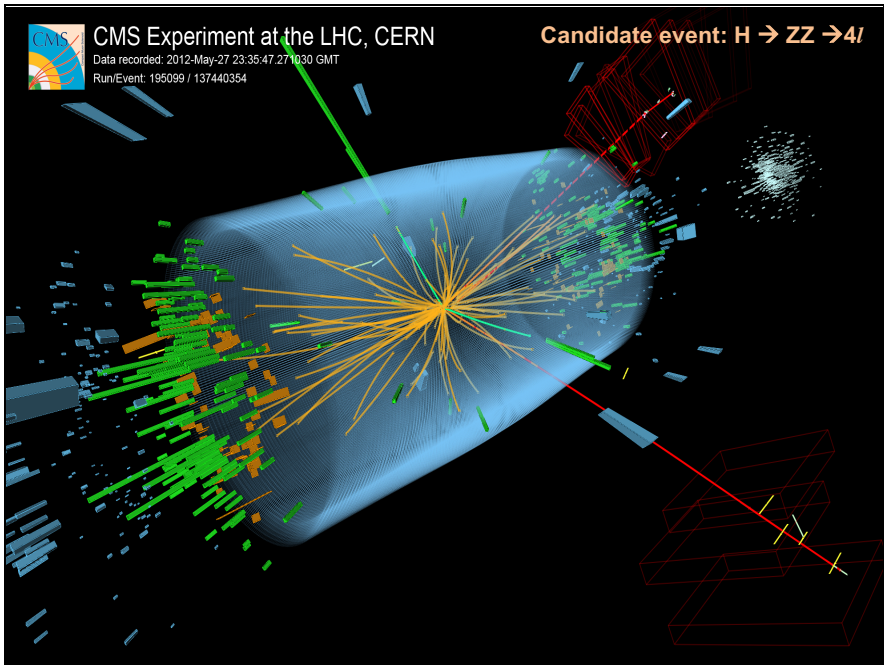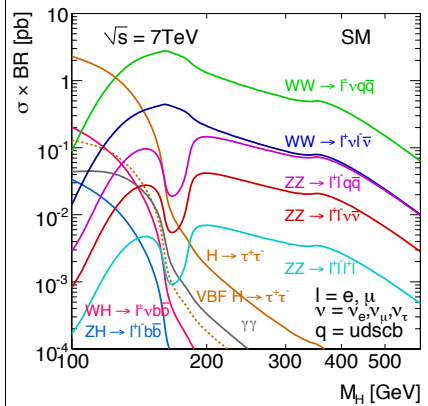| Proton - Proton | ~1300 bunches/beam |
|---|---|
| Protons/bunch | $10^{11}$ |
| Beam energy | 4 TeV ($4\times10^{12}$ eV) |
| Luminosity | $10^{34}$ cm$^{-2}$ s$^{-1}$ |
| Bunch collision frequency | 40 MHz |
| Proton collision frequency | $10^7$ - $10^9$ Hz |

**"New physics" frequency .00001 Hz**

**Event selection:**
**1 u 10 000 000 000 000**

8

# Higgs boson at LHC



g g fusion :

WW, ZZ fusion :

t t̄ fusion :

W, Z bremsstrahlung

Higgs boson ($M_H \sim 125$ GeV)
produced every ~**10 seconds**
@ $L = 5 \times 10^{33}$ cm$^{-2}$ s$^{-1}$

---

CMS Experiment at the LHC, CERN
Data recorded: 2012-May-27 23:35:47.271030 GMT
Run/Event: 195099 / 137440354
Candidate event: H → ZZ →4*l*



---

# Higgs boson: decay channels



Signal at 1 fb$^{-1}$

| Decay channel | Mass region |
|---|---|
| H → γγ | 110–150 |
| H → bb | 110–135 |
| H → ττ | 110–140 |
| H →WW →2l 2ν | 110–600 |
| H → ZZ →4l | 110–600 |
| H → ZZ →2l2τ | 180–600 |
| H → ZZ →2l2j | 226–600 |
| H → ZZ →2l2ν | 250–600 |

| $m_H$, GeV | WW→2l2ν | ZZ→4l | γγ |
|---|---|---|---|
| 120 | 127 | 1.5 | 43 |
| 150 | 390 | 4.6 | 16 |
| 300 | 89 | 3.8 | 0.04 |

The most sensitive channels for
low mass Higgs:
**H → γγ**
**H → ZZ → l⁻l⁺l⁻l⁺**

---

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

The CMS Collaboration*

**Abstract**

Results are presented from searches for the standard model Higgs boson in proton-
proton collisions at $\sqrt{s} = 7$ and 8 TeV in the CMS experiment at the LHC, using
data samples corresponding to integrated luminosities of up to 5.1 fb$^{-1}$ at 7 TeV and
5.3 fb$^{-1}$ at 8 TeV. The search is performed in five decay modes: $\gamma\gamma$, ZZ, WW, $\tau^+\tau^-$,
and bb. An excess of events is observed above the expected background, a local signif-
icance of 5.0 standard deviations, at a mass near 125 GeV, signalling the production
of a new particle. The expected significance for a standard model Higgs boson of
that mass is 5.8 standard deviations. The excess is most significant in the two decay
modes with the best mass resolution, $\gamma\gamma$ and ZZ; a fit to these signals gives a mass of
$125.3 \pm 0.4$ (stat.) $\pm 0.5$ (syst.) GeV. The decay to two photons indicates that the new
particle is a boson with spin different from one.

*This paper is dedicated to the memory of our colleagues who worked on CMS
but have since passed away.*

*In recognition of their many contributions to the achievement of this observation.*

*Submitted to Physics Letters B*

arXiv:1207.7235v1 [hep-ex] 31 Jul 2012

## Slide 13

arXiv:1207.7214v1 [hep-ex] 31 Jul 2012

**Observation of a New Particle in the Search for the Standard Model Higgs Boson with the ATLAS Detector at the LHC**

The ATLAS Collaboration

**Abstract**

A search for the Standard Model Higgs boson in proton-proton collisions with the ATLAS detector at the LHC is presented. The datasets used correspond to integrated luminosities of approximately 4.8 fb$^{-1}$ collected at $\sqrt{s} = 7$ TeV in 2011 and 5.8 fb$^{-1}$ at $\sqrt{s} = 8$ TeV in 2012. Individual searches in the channels $H \to ZZ^{(*)} \to 4\ell$, $H \to \gamma\gamma$ and $H \to WW^{(*)} \to e\nu\mu\nu$ in the 8 TeV data are combined with previously published results of searches for $H \to ZZ^{(*)}$, $WW^{(*)}$, $b\bar{b}$ and $\tau^+\tau^-$ in the 7 TeV data and results from improved analyses of the $H \to ZZ^{(*)} \to 4\ell$ and $H \to \gamma\gamma$ channels in the 7 TeV data. Clear evidence for the production of a neutral boson with a measured mass of $126.0 \pm 0.4$ (stat) $\pm 0.4$ (sys) GeV is presented. This observation, which has a significance of 5.9 standard deviations, corresponding to a background fluctuation probability of $1.7 \times 10^{-9}$, is compatible with the production and decay of the Standard Model Higgs boson.

## Slide 14

# Expectations vs measurements



Figure 4: Distribution of the four-lepton invariant mass for the ZZ → 4ℓ analysis. The points represent the data, the filled histograms represent the background, and the open histogram shows the signal expectation for a Higgs boson of mass $m_H = 125$ GeV, added to the background expectation. The inset shows the $m_{4\ell}$ distribution after selection of events with $K_D > 0.5$, as described in the text.
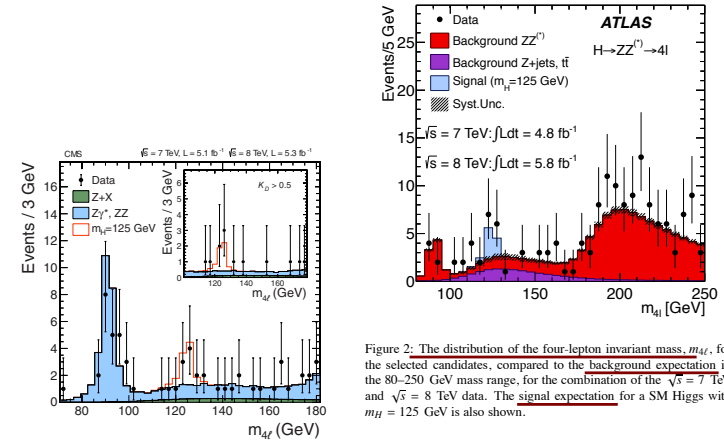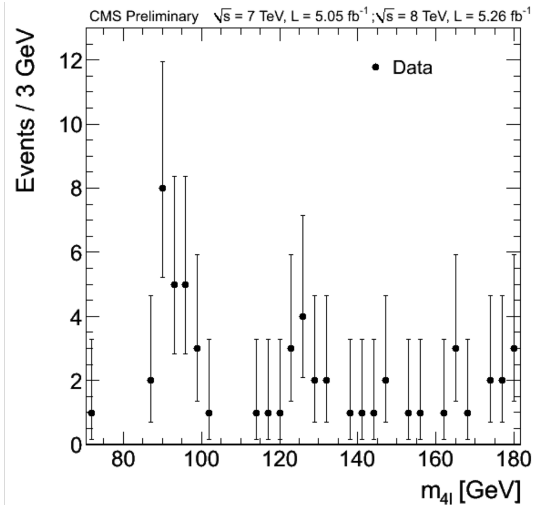
Figure 2: The distribution of the four-lepton invariant mass, $m_{4\ell}$, for the selected candidates, compared to the background expectation in the 80–250 GeV mass range, for the combination of the $\sqrt{s} = 7$ TeV and $\sqrt{s} = 8$ TeV data. The signal expectation for a SM Higgs with $m_H = 125$ GeV is also shown.

## Slide 15

# H → ZZ → $l^-l^+l^-l^+$ events distribution

## Slide 16

# H → ZZ → $l^-l^+l^-l^+$ events distribution

## Slide 17

# H → ZZ → $l^-l^+l^-l^+$ events distribution



CMS Preliminary   $\sqrt{s}$ = 7 TeV, L = 5.05 fb⁻¹ ; $\sqrt{s}$ = 8 TeV, L = 5.26 fb⁻¹

Data

Z+X

$Z\gamma^*$,ZZ

$m_H$=126 GeV

Events / 3 GeV

$m_{4l}$ [GeV]

## Slide 18

# H→γγ: Example of fitting



CMS   $\sqrt{s}$ = 7 TeV, L = 5.1 fb⁻¹ $\sqrt{s}$ = 8 TeV, L = 5.3 fb⁻¹

S/(S+B) Weighted Events / 1.5 GeV

Unweighted

Events / 1.5 GeV

$m_{\gamma\gamma}$ (GeV)

Data
S+B Fit
B Fit Component
±1σ
±2σ

$m_{\gamma\gamma}$ (GeV)

Figure 3: The diphoton invariant mass distribution with each event weighted by the $S/(S+B)$ value of its category. The lines represent the fitted background and signal, and the coloured bands represent the ±1 and ±2 standard deviation uncertainties on the background estimate. The inset shows the central part of the unweighted invariant mass distribution.

ATLAS   Data
Sig+Bkg Fit ($m_H$=126.5 GeV)
Bkg (4th order polynomial)

Events / 2 GeV

$\sqrt{s}$=7 TeV, ∫Ldt=4.8fb⁻¹
$\sqrt{s}$=8 TeV, ∫Ldt=5.9fb⁻¹   H→γγ

(a)

Events - Bkg

(b)

Data S/B Weighted
Sig+Bkg Fit ($m_H$=126.5 GeV)
Bkg (4th order polynomial)

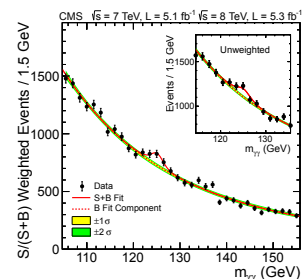Σ weights / 2 GeV

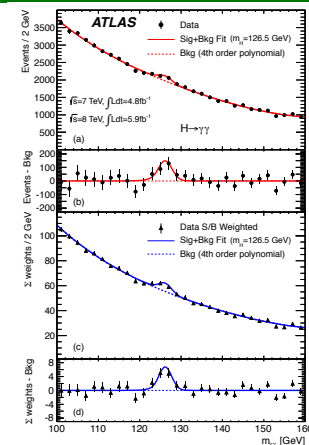(c)

Σ weights - Bkg

(d)

$m_{\gamma\gamma}$ [GeV]

Figure 4: The distributions of the invariant mass of diphoton candidates after all selections for the combined 7 TeV and 8 TeV data sample. The inclusive sample is shown in a) and a weighted version of the same sample in c); the weights are explained in the text. The result of a fit to the data of the sum of a signal component fixed to $m_H$ = 126.5 GeV and a background component described by a fourth-order Bernstein polynomial is superimposed. The residuals of the data and weighted data with respect to the respective fitted background component are displayed in b) and d).

## Slide 19

# H→bb: example of Multivariate analysis (MVA)

For the multivariate analysis, a boosted decision tree (BDT) [115, 116] is trained to give a high output value (score) for signal-like events and for events with good diphoton invariant mass resolution, based on the following observables: (i) the photon quality determined from electromagnetic shower shape and isolation variables; (ii) the expected mass resolution; (iii) the per-event estimate of the probability of locating the diphoton vertex within 10 mm of its true location along the beam direction; and (iv) kinematic characteristics of the photons and the diphoton system. The kinematic variables are constructed so as to contain no information about the invariant mass of the diphoton system. The diphoton events not satisfying the dijet selec-



CMS   $\sqrt{s}$ = 8 TeV, L = 5.0 fb⁻¹

Events/ 0.1
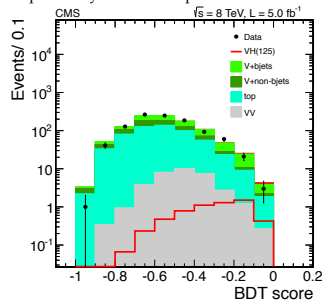
Data
VH(125)
V+bjets
V+non-bjets
top
VV

BDT score

Figure 11: Distribution of BDT scores for the high-$p_T$ subchannel of the Z($\nu\nu$)H(bb) search in the 8 TeV data set after all selection criteria have been applied. The signal expected from a Higgs boson ($m_H$ = 125 GeV), including W($\ell\nu$)H events where the charged lepton is not reconstructed, is shown added to the background and also overlaid for comparison with the diboson background.

## Slide 20

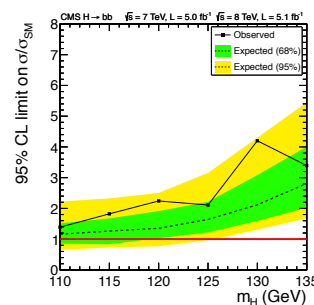# Example of limits



CMS H → bb   $\sqrt{s}$ = 7 TeV, L = 5.0 fb⁻¹ $\sqrt{s}$ = 8 TeV, L = 5.1 fb⁻¹

95% CL limit on σ/σ$_{SM}$

Observed
Expected (68%)
Expected (95%)

$m_H$ (GeV)

Figure 12: The 95% CL limit on the signal strength $\sigma/\sigma_{SM}$ for a Higgs boson decaying to two b quarks, for the combined 7 and 8 TeV data sets. The symbol $\sigma/\sigma_{SM}$ denotes the production cross section times the relevant branching fractions, relative to the SM expectation. The background-only expectations are represented by their median (dashed line) and by the 68% and 95% CL bands.

ATLAS 2011 - 2012   ± 1σ
± 2σ
Observed
Bkg. Expected

95% CL Limit on μ

$\sqrt{s}$ = 7 TeV: ∫Ldt = 4.6-4.8 fb⁻¹
$\sqrt{s}$ = 8 TeV: ∫Ldt = 5.8-5.9 fb⁻¹

(a)   CL$_s$ Limits

Local $p_0$

Sig. Expected
Observed

(b)

Signal strength (μ)

Observed
-2 ln λ(μ)<1

(c)

$m_H$ [GeV]

Figure 7: Combined search results: (a) The observed (solid) 95% CL limits on the signal strength as a function of $m_H$ and the expectation (dashed) under the background-only hypothesis. The dark and light shaded bands show the ±1σ and ±2σ uncertainties on the background-only expectation. (b) The observed (solid) local $p_0$ as a function of $m_H$ and the expectation (dashed) for a SM Higgs boson signal hypothesis (μ = 1) at the given mass. (c) The best-fit signal strength $\hat{\mu}$ as a function of $m_H$. The band indicates the approximate 68% CL interval around the fitted value.

## Slide 21

# p-value and hypothesis testing



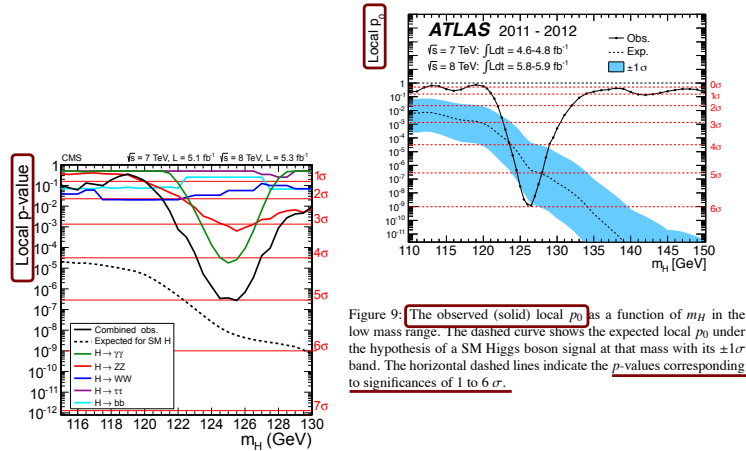Figure 9: The observed (solid) local $p_0$ as a function of $m_H$ in the low mass range. The dashed curve shows the expected local $p_0$ under the hypothesis of a SM Higgs boson signal at that mass with its $\pm 1\sigma$ band. The horizontal dashed lines indicate the p-values corresponding to significances of 1 to 6 $\sigma$.

Figure 15: The observed local p-value for the five decay modes and the overall combination as a function of the SM Higgs boson mass. The dashed line shows the expected local p-values for a SM Higgs boson with a mass $m_H$.

## Slide 22

# Measuring properties

Asymptotically, the test statistic $-2\ln\lambda(\mu, m_H)$ is distributed as a $\chi^2$ distribution with two degrees of freedom. The resulting 68% and 95% CL contours for the $H\to\gamma\gamma$ and $H\to WW^{(*)}\to \ell\nu\ell\nu$ channels are shown in



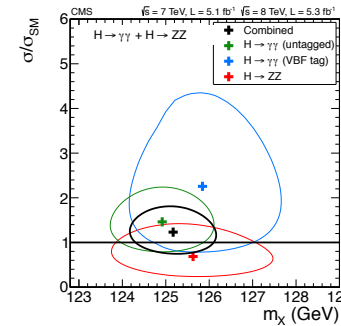Figure 11: Confidence intervals in the $(\mu, m_H)$ plane for the $H\to ZZ^{(*)}\to 4\ell$, $H\to\gamma\gamma$, and $H\to WW^{(*)}\to \ell\nu\ell\nu$ channels, including all systematic uncertainties. The markers indicate the maximum likelihood estimates $(\hat{\mu}, \hat{m}_H)$ in the corresponding channels (the maximum likelihood estimates for $H\to ZZ^{(*)}\to 4\ell$ and $H\to WW^{(*)}\to \ell\nu\ell\nu$ coincide).

Figure 17: The 68% CL contours for the signal strength $\sigma/\sigma_{SM}$ versus the boson mass $m_X$ for the untagged $\gamma\gamma$, $\gamma\gamma$ with VBF-like dijet, $4\ell$, and their combination. The symbol $\sigma/\sigma_{SM}$ denotes the production cross section times the relevant branching fractions, relative to the SM expectation. In this combination, the relative signal strengths for the three decay modes are constrained by the expectations for the SM Higgs boson.
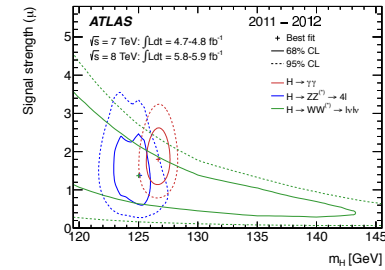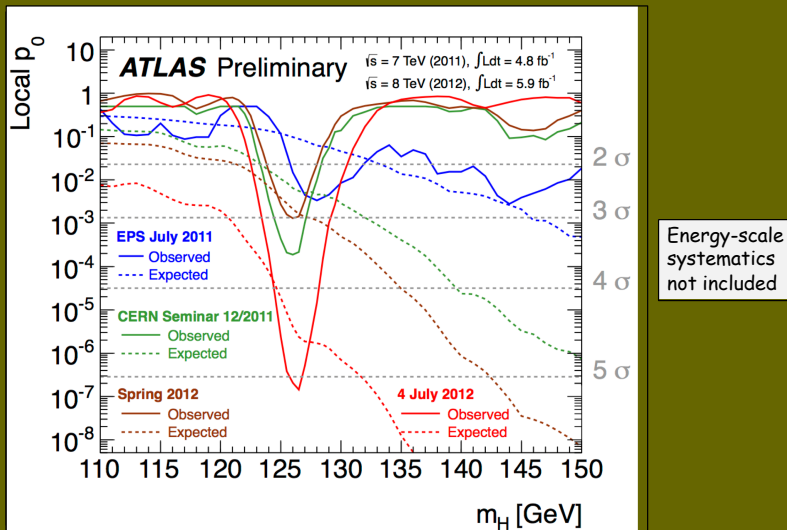
## Slide 23



Evolution of the excess with time

## Slide 24

# Conclusions of papers - ATLAS

**10. Conclusion**

Searches for the Standard Model Higgs boson have been performed in the $H\to ZZ^{(*)}\to 4\ell$, $H\to\gamma\gamma$ and $H\to WW^{(*)}\to e\nu\mu\nu$ channels with the ATLAS experiment at the LHC using 5.8–5.9 fb$^{-1}$ of $pp$ collision data recorded during April to June 2012 at a centre-of-mass energy of 8 TeV. These results are combined with earlier results [17], which are based on an integrated luminosity of 4.6–4.8 fb$^{-1}$ recorded in 2011 at a centre-of-mass energy of 7 TeV, except for the $H\to ZZ^{(*)}\to 4\ell$ and $H\to\gamma\gamma$ channels, which have been updated with the improved analyses presented here.

The Standard Model Higgs boson is excluded at 95% CL in the mass range 111–559 GeV, except for the narrow region 122–131 GeV. In this region, an excess of events with significance 5.9 $\sigma$, corresponding to $p_0 = 1.7 \times 10^{-9}$, is observed. The excess is driven by the two channels with the highest mass resolution, $H\to ZZ^{(*)}\to 4\ell$ and $H\to\gamma\gamma$, and the equally sensitive but low-resolution $H\to WW^{(*)}\to \ell\nu\ell\nu$ channel. Taking into account the entire mass range of the search, 110–600 GeV, the global significance of the excess is 5.1 $\sigma$, which corresponds to $p_0 = 1.7 \times 10^{-7}$.

These results provide conclusive evidence for the discovery of a new particle with mass $126.0 \pm 0.4$ (stat) $\pm 0.4$ (sys) GeV. The signal strength parameter $\mu$ has the value $1.4 \pm 0.3$ at the fitted mass, which is consistent with the SM Higgs boson hypothesis $\mu = 1$. The decays to pairs of vector bosons whose net electric charge is zero identify the new particle as a neutral boson. The observation in the diphoton channel disfavours the spin-1 hypothesis [140, 141]. Although these results are compatible with the hypothesis that the new particle is the Standard Model Higgs boson, more data are needed to assess its nature in detail.

# Conclusions of papers - CMS

Results are presented from searches for the standard model Higgs boson in proton-proton collisions at $\sqrt{s} = 7$ and $8\,\text{TeV}$ in the CMS experiment at the LHC, using data samples corresponding to integrated luminosities of up to $5.1\,\text{fb}^{-1}$ at $7\,\text{TeV}$ and $5.3\,\text{fb}^{-1}$ at $8\,\text{TeV}$. The search is performed in five decay modes: $\gamma\gamma$, ZZ, $W^+W^-$, $\tau^+\tau^-$, and $b\bar{b}$. An excess of events is observed above the expected background, with a local significance of $5.0\,\sigma$, at a mass near $125\,\text{GeV}$, signalling the production of a new particle. The expected local significance for a standard model Higgs boson of that mass is $5.8\,\sigma$. The global $p$-value in the search range of $115$–$130$ ($110$–$145$) GeV corresponds to $4.6\,\sigma$ ($4.5\,\sigma$). The excess is most significant in the two decay modes with the best mass resolution, $\gamma\gamma$ and ZZ, and a 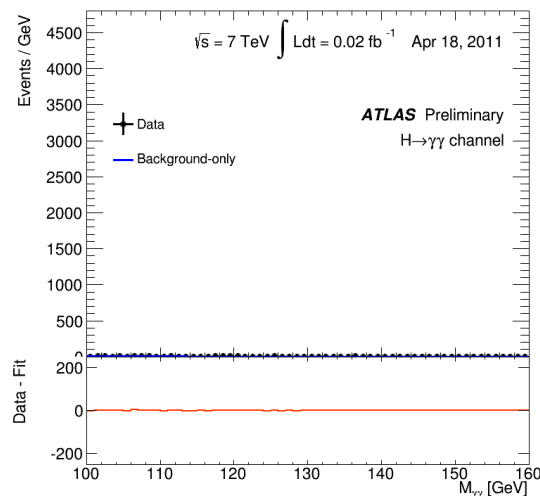fit to these signals gives a mass of $125.3 \pm 0.4\,(\text{stat.}) \pm 0.5\,(\text{syst.})\,\text{GeV}$. The decay to two photons indicates that the new particle is a boson with spin different from one. The results presented here are consistent, within uncertainties, with expectations for a standard model Higgs boson. The collection of further data will enable a more rigorous test of this conclusion and an investigation of whether the properties of the new particle imply physics beyond the standard model.

**We'll come back to this at the end of lectures**

---

# Evolution of excess: CMS H→ZZ→4$l$

---

# Evolution of excess: ATLAS H→$\gamma\gamma$

---

# Evolution of language

- **February 2012**
  - Combined results of **searches for the standard model Higgs boson** in pp collisions at sqrt(s) = 7 TeV
  - By CMS Collaboration, **Phys. Lett. B710 (2012) 26-48**
- **July 2012**
  - **Observation** of a **new boson** with a mass of 125 GeV with the CMS experiment at the LHC
  - By CMS Collaboration, **Phys. Lett. B716 (2012) 30-61**
- **December 2012**
  - Study of the Mass and Spin-Parity of the **Higgs Boson Candidate** Via Its Decays to Z Boson Pairs
  - By CMS Collaboration, **Phys. Rev. Lett. 110 (2013) 081803**
- **July 2013**
  - Measurements of **Higgs boson** production and couplings in diboson final states with the ATLAS detector at the LHC
  - By ATLAS Collaboration, **Phys. Lett. B 726 (2013) 88**
- **March 2015**
  - Combined Measurement of **the Higgs Boson** Mass in pp Collisions at sqrt(s) = 7 and 8 TeV with the ATLAS and CMS Experiments
  - By ATLAS and CMS Collaborations, **Phys.Rev.Lett. 114 (2015) 191803**

## Outline of Lecture Series

1. Introduction, Monte Carlo methods and distributions

2. Estimators and confidence intervals

3. Confidence intervals

4. Hypothesis testing

# Data Analysis

### Lecture 1: Introduction to data analysis and Monte Carlo methods

## In this lecture

- **Introduction to data analysis**
  - Confirmatory and exploratory data analysis
  - Quantitative vs graphical techniques
  - Experimental vs observational studies
  - Exploring the data

## Data analysis, statistics and probability

- **Data analysis** is the process of transforming raw data into usable information

| RAW data | Data analysis | Usable information |
|----------|---------------|--------------------|

- Data analysis uses **statistics** for presentation and interpretation (explanation) of data
  - *Descriptive statistics*
    - Describes the main features of a collection of data in quantitative terms
  - *Inductive statistics*
    - Makes *inference* about a random process from its observed behavior during a finite period of time
- A mathematical foundation for statistics is the **probability theory**

## Confirmatory and exploratory data analysis

- **Confirmatory** data analysis = Statistical **hypothesis testing**
  - A method of making statistical decisions using experimental data
  - Two main methods
    - **Frequentist** hypothesis testing
      - Hypothesis is either true or not
    - **Bayesian** inference
      - Introduces a "degree of belief"
- **Exploratory** data analysis
  - Uses data to suggest hypothesis to test
  - Complements confirmatory data analysis
  - Main objectives:
    - Suggest hypothesis about the causes of observed phenomena
    - Asses assumptions on which statistical inference will be based
    - Select appropriate statistical tools and techniques
    - Eventually suggest further data collection

## Quantitative vs graphical techniques

- **Quantitative techniques** yield numeric or tabular output
  - Hypothesis testing
  - Analysis of variance
  - Point estimation
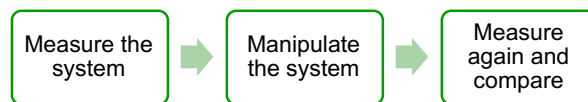  - Interval estimation
- **Graphical techniques**
  - Used for gaining insight into data sets in terms of testing assumptions, model selection, estimator selection ...
  - Provide a convincing mean of presenting results
  - Includes: graphs, histograms, scatter plots, probability plots, residual plots, box plots, block plots, biplots
  - Four main objectives:
    - Exploring the **content** of a data set
    - Finding **structure** in data
    - Checking **assumptions** in statistical models
    - **Communicate** the results of an analysis

## Experimental vs observational studies

- **Experimental studies**

| Measure the system | → | Manipulate the system | → | Measure again and compare |

  - Example: Study of whether and how much a free coffee would improve working performace of scientists in Building 40 at CERN

- **Observational studies**
  - No experimental manipulation
  - Data are gathered and analysed
  - Example:
    - Study of correlation between number of beers drunk in a pub on Wednesday evening on performance on the exam the day after
    - Be careful who pays! → see later
    - One could discuss whether to manipulate or not the system ☺

## Experiments – basic steps

| Planning | • Select subject to study<br>• Select an information source |
| Design and Building | • Design an experiment<br>• Build and test a model (f.g. MC simulation)<br>• Once happy with the model build the experiment |
| Collecting data | • Employ descriptive statistics to summarize data<br>• Suppres details<br>• Early exploratory analysis |
| **Analysing data** | • **Statistical inference**<br>• **Reach a consensus what observations tell about an underlaying reality** |
| Presenting Documenting | • Publish article and disseminate results<br>• Enjoy in the fruits of the hard work! |

## LHC experiments – basic steps

**Planning**
- Started ~ 30 years ago (Aachen 1989)
- Core teams from previous experiments UA1&2

**Design Building**
- 'Best' experimental design chosen (CMS, ATLAS, ALICE and LHCb)
- Detailed MC simulations performed before started to build

**Collecting data**
- Trigger and DAQ carefully planed and built
- MC simulation used for optimization

**Analysing data**
- **Statistical inference → a part of work done at this school too (learning methods&tools)**
- **For the consensus → let's see ☺**

**Presenting Documenting**
- Many articles published
- And first discoveries announced and published!

---

## What we (will) measure at LHC?

**Something we already know**
- At the very beggining of the LHC operation
- For example: production of W and Z bosons

**Something that (probably) exists but wasn't measured yet**
- Simply because we are exploring new energy domain
- Standard Model processes
- But surprises are always possible

**Hopefully something new but reasonably expected**
- Altought "reasonably" is not very well defined ☺
- For example we all expected to find the Higgs boson → and we did find it!
- Heavy neutrions?

**Maybe something new but less likely**
- New heavy bosons (Z', W')
- Micro black holes
- Extra dimensions

**Something completely unexpected**
- Well, it's hard to look for unexpected ☺

---

## Some of the physicists' jargon

- **Cross section ($\sigma$)**
  - A measure of 'frequency' of the physical process
  - Units: barns ($10^{-28}$ cm$^2$)
    - Typical values: femtobarns (fb), picobarns (pb)
- **Luminosity (L)**
  - Or *instantenous luminosity*
  - A measure of collisions 'frequency'
    - Typical (at Tevatron/Early LHC): $L = 10^{32}$ cm$^{-2}$s$^{-1}$
- **Integrated luminosity ($\mathcal{L} = \int Ldt$)**
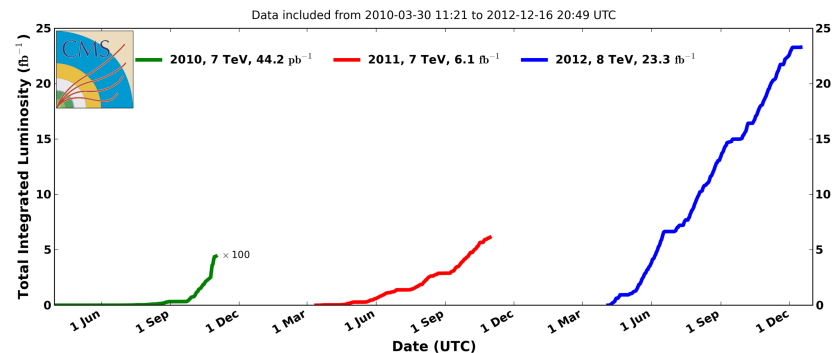  - A measure of number of accumulated collisions after a certain time period
  - Units: (cross section)$^{-1}$ …. E.g. 1 fb$^{-1}$ = 1000 pb$^{-1}$
    - Typical (Tevatron/Early LHC): few fb$^{-1}$
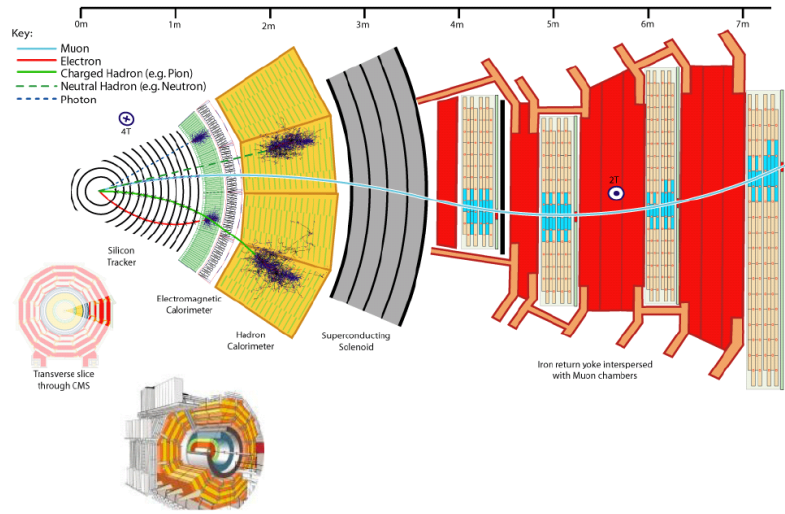- **Number of events (N)**
  - Number of (expected) events (N) after a certain time of running

$$N = \sigma \cdot \mathcal{L}$$

---

## Data collected by CMS in Run 1

## Measuring physical objects

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

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

Silicon Tracker
Electromagnetic Calorimeter
Hadron Calorimeter
Superconducting Solenoid
Iron return yoke interspersed with Muon chambers

Transverse slice through CMS

---

## Data analysis - general picture

2 Sampling a reality **Experiment**

1 **Physical phenomena** *Described by a theory*

Described by PDFs, depending on $p$ uknown parameters with true values

$$\theta^{true} = (\theta_1^{true}, \theta_2^{true}, \ldots, \theta_p^{true})$$

*For example:*

$$\theta^{true} = (m_H^{true}, \Delta m_s^{true}, \ldots, \sigma_{tot}^{true})$$

4 Data analysis

3 **Data sample** $x = (x_1, x_2, \ldots, x_N)$

For example:

$$x = (event_1, \ldots, event_N)$$

5 **Results**
- parameter estimates
- confidence limits
- hypothesis tests

In statistics $x$ is a multivariate random variable *(each event has many properties, all potential variables)*

---

## Data analysis – general picture

The main goal: **learn more about NATURE**

For example, let's suppose the TRUE state of nature is: **Higgs boson exists with the mass of $m_H$(true) = 134.26 GeV**

**Make an experiment and obtain a DATA SAMPLE**

Use data sample to examine this!

**Events collected after some time of LHC running**

| Events collected |
|---|
| Event 1 |
| Event 2 |
| ... |
| Event N |

**Event 1**

| Event 1 |
|---|
| Object 1 |
| Object 2 |
| ... |
| Object k |

**If Object 1 == electron**

| |
|---|
| $p_x$ |
| $p_y$ |
| $p_z$ |
| E |
| ... |

$N \sim 100/s \times 10^7$ s/year
**$N \sim 10^9$ events per year**

Objects $\equiv$ reconstructed objects i. e. electrons, photons, jets, muons ...

---

## Analysis steps in typical LHC analysis

1. Event reconstruction
2. Event selection
3. Background estimation
4. Systematic uncertainties
5. Yields and kinematics distributions
6. Kinematic discriminant
7. Statistical analysis and results

## Signal vs background(s)

- **Signal**: an event coming from the physical process under study
  - Example: H→ ZZ→$e^+e^-e^+e^-$ (henceforth both $e^+$ and $e^-$ are 'electron')
- **Background**: any other event
  - 'Dangerous' background is any other process giving at least 4 electrons in the final state
    - But be careful: electrons seen by detector are reconstructed objects and in some cases when some other objects (f.g. jets) are misreconstructed as electrons
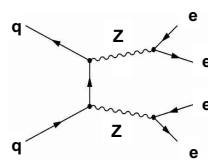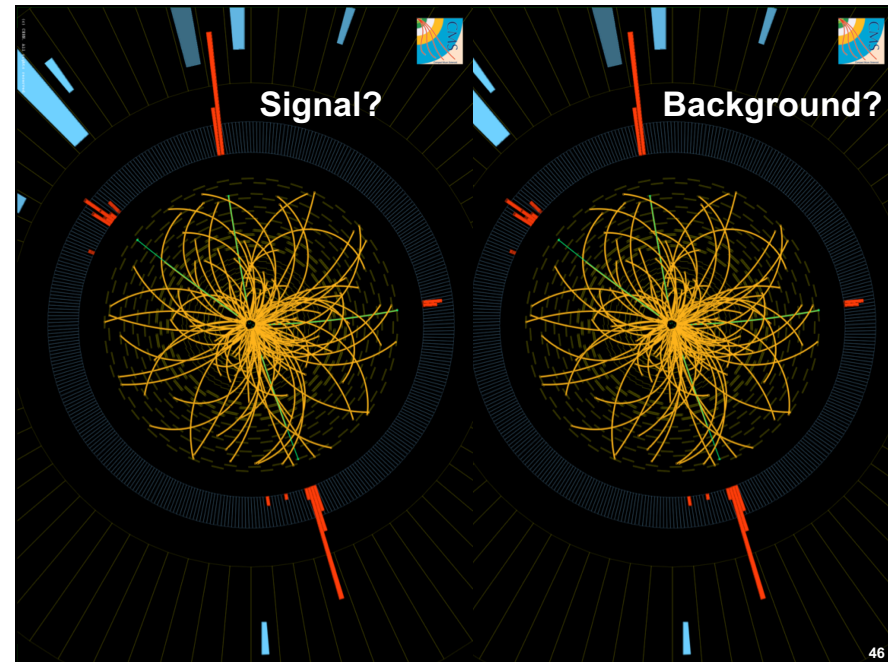  - 'Trivial' backgrounds are all other backgrounds and are easily *rejected* by a simple requirement of having at least 4 electrons in the final state

Signal: pp→H→ZZ→4e          'Dangerous' background: pp→ZZ→4e

---

**Signal?**          **Background?**

---

## Separating signal and background

- Ultimate goal of the analysis: separate as much as possible signal from background events to obtain a **reduced sample** as clean as possible
  - This is usually obtained in several steps

  [ Trigger ] ➡ [ Preselection ] ➡ [ Selection ]

  - Usually all these steps have substeps
  - More in example on the next page
- Be aware:
  - Nature is probabilistic, i.e. for a given event it'll never be possible to tell whether it's signal or background!
  - We can only make an educated guess → attribute probabilities that the observed event comes from signal or background

    ***p(event|signal)* and *p(event|background)***

- Very often we have to solve the following statistical problem: **maximum reduction of the background for a given signal acceptance**

---

## Exploring the data

- Once data are collected → exploratory data analysis
  - Heavily use of graphical techniques
- Example: **data reduction =** Preselection
  - Goal: **getting rid of all unuseful events**
  - Unusefullness is not uniquely defined:
    - We have a certain interest to keep some background events for better control and its measurement from data
  - Some numbers:
    - ~ $10^9$ events collected per year (after trigger)
    - ~ 1 MB event size on a tape (rought estimate)
    - ⇒ ~ 1 PB of data collected per year → non manageable at once
  - Interested physical processes are rare
    - F.g. just a handful (~10) H→ZZ→4e events per year
    - So be careful when choosing criteria for data reduction not to lose too many signal events

## Slide 49

### Example: H→ZZ→4e in CMS

- Very basic cuts: High Level Trigger+ ≥ 3 electrons, any charge and $p_T^{1,2,3}$ > 10, 10, 5 GeV/c
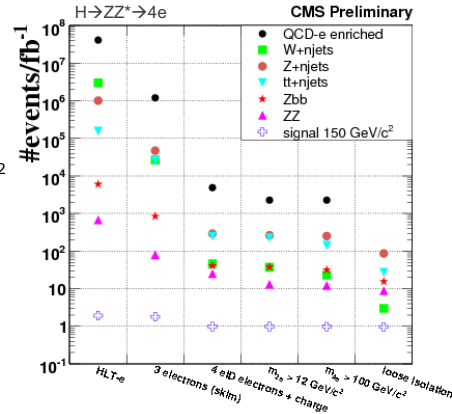
- Preselection cuts:
  - ≥ 2 *ee* pairs of identified, opposite charge and same flavor leptons with
    - $p_T$ > 5 GeV/c; $|\eta|$ < 2.5
  - At least two $m_{ee}$ > 12 GeV/c²
  - At least one $m_{4e}$ > 100 GeV/c²
  - Loose track based isolation

- After these steps
  - Some background gone
  - Some heavily reduced
  - Some still resisting

- Full **selection** needed for the final analysis



H→ZZ*→4e       CMS Preliminary

Legend: QCD-e enriched, W+njets, Z+njets, tt+njets, Zbb, ZZ, signal 150 GeV/c²

## Slide 50

### Analysis steps with some details   Just for illustration (not for exam ☺)

**Trigger**
- **Level 1 Trigger**: Electron-gamma and muon seed
- **High Level Trigger** : Multilepton paths

Triggered events

**Object reconstruction**
- **Electrons** Isolation, ID, $SIP_{3D}$   + – ...
- **Muons** Isolation, ID, $SIP_{3D}$   + – ...
- **FSR photons**

**Event selection**
- **$Z_1$**: A pair of opposite sign same flavor leptons closest to Z, with 40 < $M_{ll}$ < 120 GeV
- **$Z_2$**: Remaining lepton combination with highest $p_T$ sum and with 12 < $M_{ll}$ < 120 GeV
- **Additional**: at least one lepton with $p_T$ > 20 GeV and one with $p_T$ > 10 GeV, all opposite-charge lepton pairs with $m_{ll'}$ > 4 GeV
- **Higgs boson phase space**: $m_{4l}$ > 100 GeV

Selected events

**Analysis**
- Background estimates
- Systematics
- Kinematic discriminant
- Statistical analysis
  - Limits
  - P-values
  - Signal strength

## Slide 51

# Probability

# Random variables

## Slide 52

### Probability – basic concepts

- Definitions of probability
  - **Mathematical probability**
    - Probability is a basic and an abstract concept
  - **Frequentist probability**
    - Using only measured frequencies
  - **Bayesian probability**
    - Based on a *degree of belief*

# Mathematical probability

- Developed in 1933 by Kolmogovor in his "*Foundations of the Theory of Probability*"
- Define $\Omega$ as an exclusive set of all possible elementary events $x_i$
  - Exclusive means the occurence of one of them implies that none of the others occurs
- We define the probability of the occurency of $x_i$, $P(x_i)$ to obey the **Kolmogorov axioms**:

$$(a)\ P(x_i) \geq 0 \quad \text{for all } i$$
$$(b)\ P(x_i \text{ or } x_j) = P(x_i) + P(x_j)$$
$$(c)\ \sum_{\Omega} P(x_i) = 1$$

- From these properties more complex probability expressions can be deduced
  - For non-elementary events, i.e. set of elementary events
  - For non-exclusive events, i.e. overlapping sets of elementary events

# Frequentist probability

- Experiment:
  - N events observed
  - Out of them $n$ is of type $x$
- **Frequentist probability** that any single event will be of type $x$

$$P(x) = \lim_{N \to \infty} \frac{n}{N}$$

- Important restriction: such a probability can only be applied to repeatable experiments
  - For example one can't define a probability that it'll snow tomorrow
  - Altough this seems to be a serious problem, a job of scientist is to try to get as close as possible to repeatable experiments and produce reproducible results
- Frequentist statistics is often associated with the names of *Jerzy Neyman* and *Egon Pearson*

# Bayesian probability

- Based on a concept of "degree of belief"
- An operational definition of belief is based on coherent bet by Finneti
  - **What's amount of money one 's willing to bet based on her/his belief on the future occurence of the event**
- Bayesian inference uses Bayes' formula for conditional probability:

$$P(H \mid D) = \frac{P(D \mid H) P(H)}{P(D)}$$

- $H$ is a **hypothesis**, and $D$ is the **data**.
- $P(H)$ is the **prior probability** of $H$: the probability that $H$ is correct before the data $D$ was seen.
- $P(D|H)$ is the **conditional probability** of seeing the data $D$ given that the hypothesis $H$ is true. $P(D|H)$ is called the **likelihood**.
- $P(D)$ is the **marginal probability** of $D$.
  - $P(D)$ is the prior probability of witnessing the data $D$ under all possible hypotheses
- $P(H|D)$ is the **posterior probability**: the probability that the hypothesis is true, given the data and the previous state of belief about the hypoth.

## Example: Who will pay the next round?

You meet an old fried at Göttingen in a pub. He proposes that the next round should be payed by whichever of the two extracts the card of lower value from a pack of cards.

This situation happens many times in the following days. What is the probability that your friend cheats if you end up paying *wins* consecutive times[2]

You assume:
- $P(cheat) = 5\%$ and $P(honest) = 95\%$. (Surely an old friend is an unlikely cheater …)
- $P(wins|cheat) = 1$ and $P(wins|honest) = 2^{-wins}$

Bayesian solution:

$$P(cheat|wins) = \frac{P(wins|cheat)P(cheat)}{P(wins|cheat)P(cheat) + P(wins|honest)P(honest)}$$

$$P(cheat|0) = \frac{1 P(cheat)}{1 P(cheat) + 2^{-0} P(honest)} = \frac{0.05}{0.05 + 0.95} = 5\%$$

$$P(cheat|5) = \frac{1 P(cheat)}{1 P(cheat) + 2^{-5} P(honest)} = \frac{0.05}{0.05 + 0.03} = 63\%$$

[2]Adapted from G. D'Agostini, *Bayesian Reasoning in High-Energy Physics: Principles and Applications*, CERN-99-03, 1999

## Example: Learning by experience

The process of updating the probability when new experimental data becomes available can be followed easily if we insert

- $P(cheat) = P(cheat|wins - 1)$ and $P(honest) = P(honest|wins - 1)$, where $wins - 1$ indicate the propability assigned after *the previous win*
- $P(wins = 1|cheat) = P(win|cheat) = 1$ and $P(wins = 1|honest) = P(win|honest) = \frac{1}{2}$

Iterative aplication of the Bayes formula for $P(cheat|wins) =$

$$\frac{P(win|cheat)P(cheat|wins - 1)}{P(win|cheat)P(cheat|wins - 1) + P(win|honest)P(honest|wins - 1)}$$
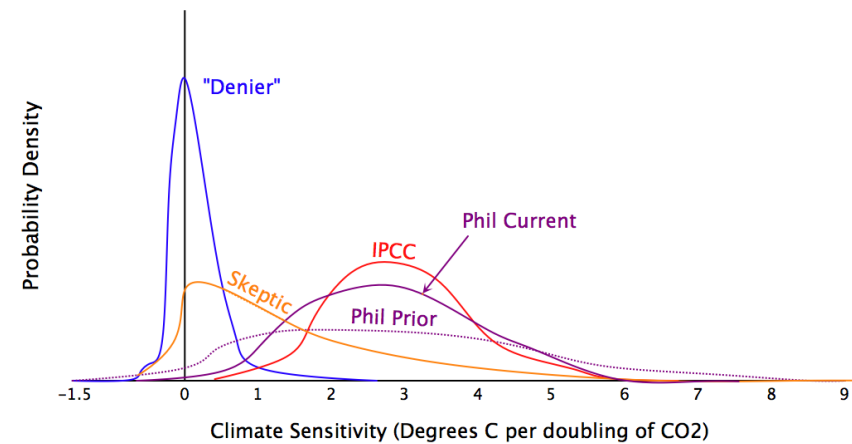
$$= \frac{P(cheat|wins - 1)}{P(cheat|wins - 1) + \frac{1}{2}P(honest|wins - 1)}$$

| P(cheat) % | P(cheat|wins) wins=5 | 10 | 15 |
|---|---|---|---|
| 1 | 24 | 91 | 99.7 |
| 5 | 63 | 98 | 99.94 |
| 50 | 97 | 99.9 | 99.997 |

When you learn from the experience, your conclusions no longer depend on the initial assumptions.

---

## Example: Priors and posteriors – expressing degree of belief

Phil is learning from experience:



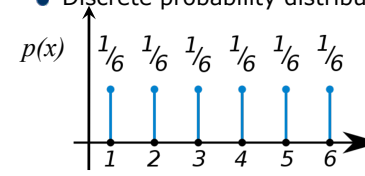(From discussion of climate change on Andrew Gelman's blog.)

---

## Random variables

- **Random event**: event having more than one possible outcome
  - Each outcome may have associated probability
  - Outcome not predictable, only the probabilities known
- Different possible outcomes may take different possible numerical values $x_1, x_2, \ldots \rightarrow$ **random variable** $x$
  - The corresponding probabilities $P(x_1), P(x_2), \ldots$ form a **probability distribution**
- If observations are **independent** the distribution of each random variable is unaffected by knowledge of any other observation
- When an experiment consists of $N$ repeated observations of the same random variable $x$, this can be considered as the single observation of a random vector $\boldsymbol{x}$, with components $x_1, \ldots, x_N$

---

## Random variables: discrete

- Rolling a die:
  - Sample space = {1,2,3,4,5,6}
  - Random variable $x$ is the number rolled

$$x = \begin{cases} 1 & \text{if a 1 is rolled} \\ 2 & \text{if a 2 is rolled} \\ 3 & \text{if a 3 is rolled} \\ 4 & \text{if a 4 is rolled} \\ 5 & \text{if a 5 is rolled} \\ 6 & \text{if a 6 is rolled} \end{cases}$$

  - Discrete probability distribution

# Random variables: continuous

- A spinner
  - Can choose a real number from [0,2n]
  - All values equally likely
  - $X$ = the number spun
  - Probability to select any real number = 0
  - Probability to select any **range** of values > 0
    - Probability to choose a number in [0,n] = 1/2
  - Now we say that probability **density** p(x) of x is *1/2n*
  - Probability to select a number from any range $\Delta x$ is $\Delta x/2n$
  - More general

$$P(A < x < B) = \int_A^B p(x)dx$$

# Probability density function

- Let $x$ be a possible outcome of an observation and can take any value from a continuous range
- We write $f(x;\theta)dx$ as the probability that the masurement's outcome lies betwen $x$ and $x + dx$
- The function $f(x;\theta)dx$ is called the **probability density function (PDF)**
  - And may depend on one or more parameters $\theta$
- If $f(x;\theta)$ can take only **discrete values** then $f(x;\theta)$ is itself a **probability**
- The p.d.f. is always normalized to unit area (unit sum, if discrete)
- Both $x$ and $\theta$ may have multiple components and then written as vectors
- If $\theta$ is unkown we may wish to estimate its value from a set of measurements of $x$ → Parameter estimation in Lecture 3

# Cumulative and marginal distributions

- **Cumulative distribution function, CDF**
  - For every real number $Y$, the CDF of $Y$ is equal to the probability that the random variable $x$ takes a value less or equal to $Y$

$$F(Y) = P(x \le Y) = \int_{x_{min}}^{Y} f(x)dx$$

  - If $x$ restricted to $x_{min} < x < x_{min}$ then $F(x_{min}) = 0$, $F(x_{max}) = 1$
  - $F(x)$ is a monotonic function of $x$
- **Marginal density function**
  - Is the projection of multidimensional density
  - Example: if $f(x,y)$ is two-dimensional PDF the marginal density $g(x)$ is

$$g(x) = \int_{y_{min}}^{y_{max}} f(x,y)dy$$

- **Introduction to Monte Carlo method**
  - Monte Carlo techniques
  - Monte Carlo in HEP

# Monte Carlo



## Why simulations?

- Design studies
  - optimize detectors before building them
  - estimation of performances, costs ...

- Development of reconstruction algorithmes
  - exploring different algorithmes,
  - tuning parameters
  - optimizing analysis

- Simulation is a good way to save money!

- But, there is even a "deeper" reason!

## Why simulations?

EXPERIMENT  ← Is consistent with →  THEORY

Are measured by / Is compared with

Predicts

OBSERVABLES

From M. Liendl, Experiment Simulation – CSC 2008

## Example: classical mechanics



Galilei
Brahe  observations
Kepler

Newton: mechanics and gravitation

$$|F| = \frac{GM_AM_B}{|r_{AB}|^2}$$

EXPERIMENT  ← Is consistent with →  THEORY

Are measured by / Is compared with

Predicts

OBSERVABLES

## Slide 69

# Monte Carlo method

- **Monte Carlo methods** (MCMs) are a class of computational algorithms that rely on repeated **random sampling** to compute their results
  - MCMs use random or pseudo-random numbers
  - MCMs tend to be used when it is unfeasible or impossible to compute an exact result with a deterministic algorithm
- The term "**Monte Carlo method**" was coined in the 1940s by physicists working on nuclear weapon projects in the Los Alamos National Laboratory
- Generally MCMs are used in
  - *Studying systems with a large number of coupled (interacting) degrees of freedom*
    - such as fluids, disordered materials, strongly coupled solids, and cellular structures
  - *Modeling phenomena with significant uncertainty in inputs*
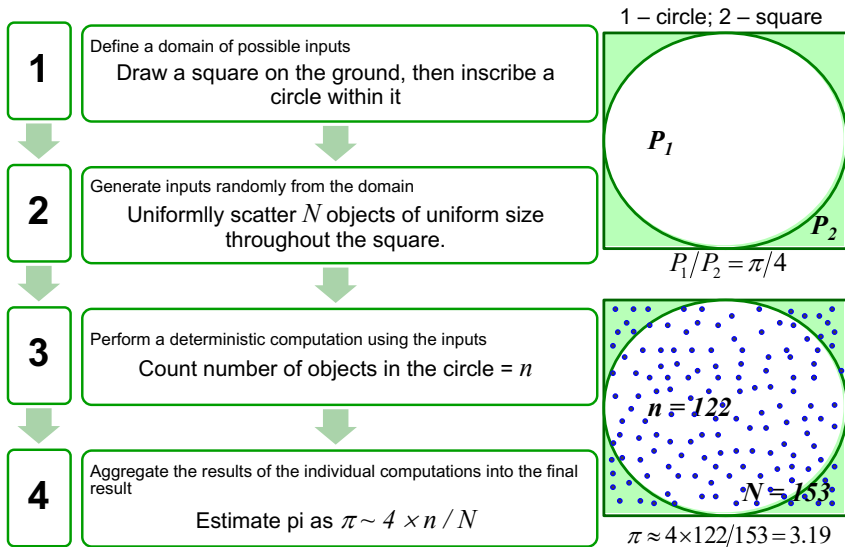    - such as the calculation of risk in business
  - *Evaluation of definite integrals*
    - particularly multidimensional integrals with complicated boundary conditions

## Slide 70

# Monte Carlo Methods – usual pattern

| | |
|---|---|
| **1** | Define a domain of possible inputs |
| **2** | Generate inputs randomly from the domain |
| **3** | Perform a deterministic computation using the inputs |
| **4** | Aggregate the results of the individual computations into the final result |

## Slide 71

# Example 1: estimating $\pi$

| | |
|---|---|
| **1** | Define a domain of possible inputs<br>Draw a square on the ground, then inscribe a circle within it |
| **2** | Generate inputs randomly from the domain<br>Uniformlly scatter $N$ objects of uniform size throughout the square. |
| **3** | Perform a deterministic computation using the inputs<br>Count number of objects in the circle = $n$ |
| **4** | Aggregate the results of the individual computations into the final result<br>Estimate pi as $\pi \sim 4 \times n / N$ |

1 – circle; 2 – square

$P_1$

$P_2$

$P_1/P_2 = \pi/4$

$n = 122$

$N = 153$

$\pi \approx 4 \times 122/153 = 3.19$

## Slide 72

# Example 2: integration

- Analytical solution

  f(x) = x

  A

  $\int x \cdot dx = 1/2 \cdot x^2$

  A = 1/2· (b$^2$ - a$^2$)

- Deterministic algorithm

  f(x) = x

  $1/2 \cdot (\Delta x)^2$

  n .. number of steps

  $\Delta x = (b - a) / n$

  A ~ $\sum$ f(a+i·$\Delta$x)·$\Delta$x

## Monte Carlo solution

$f(x) = x$

$f(b) = b$

$y_i$    $p_i$

$A$

$a$    $x_i$   $b$

$B = (b-a) \cdot b$

**n pairs** of random numbers:
$p_i = (x_i, y_i),\ i = 1...n$

   $x_i$ ... equally distributed in [a,b]

   $y_i$ ... equally distributed in [0,b]

set a counter c=0

for each pair $p_i$:

   **c = c+1** only if $p_i$ within A

**A ~ c/n·B = c/n·(b-a)·b**

**Statistical estimation** of the integral. The better, the more pairs $p_i$

---

## Random number generation

- **Physical methods**
  - "true" random numbers from "unpredictable" process
    - Example: dice, coin flopping, roulette → still in use
  - TRUE random numbers from random atomic or subatomic physical phenomena
    - Example: radioactive decay, amplitude of noise in radio
- **Computational methods**
  - Pseudo-random number generators create long runs (for example, millions of numbers long) with good random properties but eventually the sequence repeats
    - Example: Linear congruential generator   $X_{n+1} = (aX_n + b)\bmod m$
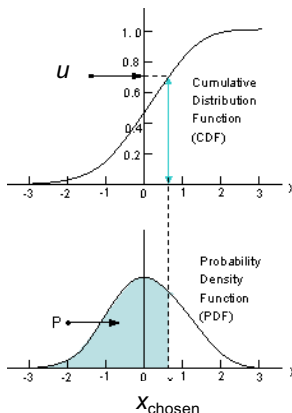- **Generation from a probability distribution f(x)**
  - Generate random numbers distributed according to the f(x)
  - Method involves transforming an uniform random number in some way
    - Examples: inversion method, acceptance-rejection method

---

## Inversion method

- Let $x$ be a random variable whose distribution can be described by the cumulative distribution function $F(x)$.
- **We want to generate values of $x$ which are distributed according to this distribution.**
- Method:

  1. Generate a random number from the standard uniform distribution; call this $u$.

  2. Compute the value $x$ such that $F(x) = u$; call this $x_{chosen}$.

  3. Take $x_{chosen}$ to be the random number drawn from the distribution described by $F$.

$u$

Cumulative Distribution Function (CDF)

Probability Density Function (PDF)

P

$x_{chosen}$

---

## Acceptance-rejection method

- It generates sampling values from an arbitrary PDF function $f(x)$ by using an instrumental distribution $h(x)$ for which we know how to sample
  - under the only restriction that $f(x) < Ch(x)$ where $C > 1$
- Usualy used in cases where the form of $f(x)$ makes sampling difficult
- Algorithm

  1. Sample $x$ from h(x) and $u$ from $U(0,1)$

  2. Check whether or not $u < f(x)$ / Ch(x).

  **NO**    **YES**

  3. Accept $x$ as a realization of $f(x)$

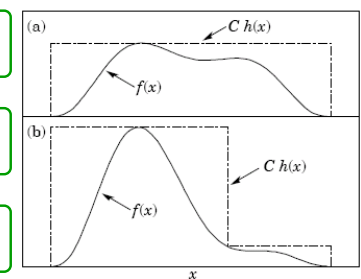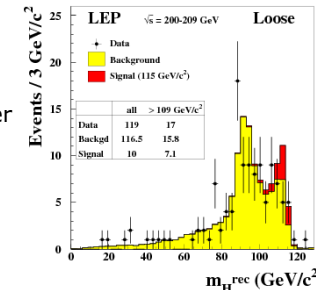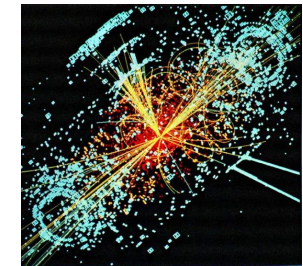(a)   $C\,h(x)$   $f(x)$

(b)   $C\,h(x)$   $f(x)$

$x$

Figure from PDG

# MC methods in engineering

- **Wireless network planing**
  - Various scenaria depending on: number of users, users' location, services users want to use
  - MC used to generate users and their states, so that network performace can be evaluated and optimized
- **Computer graphics**
  - MC methods efficient in production of photorealistic images of virtual 3D models
  - Application in video games, computer generated films, special effects in cinema
- **Wind power engineering**
  - From measured distributions of wind speeds MC generates single values for wind power sistem performace evaluation and optimization

# Monte Carlo in HEP

- Monte Carlo methods are widely used in High Energy Physics
  - **Theoretical calculations**
    - Total cross sections, differential cross sections distributions ...
  - **Event generation**
    - Distribute events according to expected probabilites
    - Many event generators on the market: f.g. PYTHIA, HERWIG
  - **Detector simulation**
    - Passage of particles through the matter
    - GEANT
  - **Data analysis**
    - Background predictions (if not measured from data)
    - Signal predictions
    - Final results

# Monte Carlo in HEP

Adopted from T. Sjöstrand, CERN Academic Training Lectures 2005

| Simulation | 'Real life' |
|---|---|
| **Event Generation** | **Collisions** |
| *Tools:* **MC** generators (PYTHIA, ...) | *Tools:* Accelerator (LHC, Tevatron ...) |
| *Output:* final state particles | *Output:* final state particles |
| **Detector simulation** | **Data acquisition** |
| *Tools:* **MC** simulators (GEANT) | *Tools:* Detectors (CMS, ATLAS,...) |
| *Output:* simulated detector response | *Output:* detector response |

| **Event reconstruction** |
|---|
| *Tools:* Detectors' software packages (custom made; **MC** used in algorithms) |
| *Output:* reconstructed physical objects (electrons, muons, jets ...) |

| **Data analysis** |
|---|
| *Tools:* Statistics (ROOT, ...; **MC** used in algorithms; f.g. **Toy MC**) |
| *Output:* new knowledge (parameter/interval estimates, hypothesis tests, article, talks ...) |