

Analysis Methods

Inês Ochoa

Course on Physics at the LHC 2024



From collider data to fundamental physics: the role of an experimentalist*

Inês Ochoa

Course on Physics at the LHC 2024

* with a strong ATLAS bias

Introduction

- The role of an experimentalist is to piece together all the elements in the chain that links theory and data.
 - How do all these pieces come together?
- Main topics:
 - 1. Event reconstruction
 - 2. Full and fast detector simulation & data-driven correction methods
 - 3. Case study: measuring the WH cross-section
 - 1. Simulation-based backgrounds, global fit,
 - 4. Case study: searching for new physics resonances
 - 1. Data-driven background estimation methods
 - 2. Anomaly detection

What physics comes out of LHC data?

Total, fiducial* cross-sections Unfolded** differential cross-sections ATLAS Preliminary ---- Data dơ/dp_T^{4I} [fb/GeV] √s=13.6 TeV, 29 fb⁻¹ ATLAS Preliminary 🛞 Data Statistical Uncertainty Sherpa qqNLO+ggLOx1.7(→ ZZ) (*) √s = 13.6 TeV. 29 fb¹ ΖZ **Total Uncertainty** 10⊨ MATRIX qqNNLO+ggNLO(\rightarrow ZZ) (*) $ZZ \rightarrow 4I$ ---- Predicted MATRIX ga[NNLOxNLO EW]+gaNLO(\rightarrow ZZ) (*) This measurement 16.9 ± 0.7 (stat.) ± 0.8 (sys.) pb **A** *************************** Sherpa qqZZ NLO + ggZZ LO×1.7(*) 17.0^{+1.9}_{-1.4} (sys.) pb 10^{-} MATRIX qqZZ NNLO + ggZZ NLO(*)

16

14

18

total cross-section [pb]

20

12

10

(*) + Powheg gqNLO(→ ZZjj) 10^{-2} 10 Prediction/Data 67 10 10^{2} 2×10² 20 30 40 p₊^{4|} [GeV] ** correcting for detector effects



<u>STDM-2022-17</u>

17.9 ± 0.4 (sys.) pb

16.7±0.4 (sys.) pb (*) + Powheg EW ZZjj

MATRIX ggZZ NNLO×NLO.EW + ggZZ NLO(*)

*in the detector's acceptance



What physics comes out of LHC data?



Fundamental properties of particles, e.g. Higgs boson mass





7

 $Z \sim$

Н

Data

Simulation



DataSimulation

- \star an invariant mass
- Lepton reconstruction & id



DataSimulation

- \star an invariant mass
- Lepton reconstruction & id
- Calibrations, detector alignment, pile-up, much more...

How do we get to this plot?



The data



A pair of topquarks produced in **ATLAS**



A pair of topquarks produced in ATLAS

Trigger and Data Acquisition: a simplified picture

Data Acquisition





Raw data

Trigger system

trigger

path

data

path

Trigger and DAQ

Trigger

DAO

trigger decisions

Storage

00000004 0000001 0000c89c aa1234aa 00003227 0000001c 04000000 00793c29 00000001 0000000 00000000 50753e27 0ab16f70 00097a2b 00000000 00033dac 00000063 920117d5 00000aa8 0000008 00000000 dd1234dd 0000002d 00000009 04000000 00210000 00000002 00000000 92011d7f 0000000 ee1234ee 00000009 03010000 00210000 00033dac 920117d5 00000aa8 00000081 00000000 2003e766 2013e282 201490d2 9c122017 ef322018 9d562023 dfa22039 c2224000 2040aa82 2041c3a2 204282b3 20489082 2057efb2 205a8616 2063cce2 2066aee2 2068a0c2 20768ff7 99522077 de72207b d822400 00000000 00000000 00000002 00000015 00000001 d04326b2 dd1234dd 0000002d 00000009 00210001 00000002 00000000 92011d80 00000001 ee1234ee 00000009 03010000 00210001 920117d5 00000aa8 00000081 00000000 2004af72 2010a3f2 20128ec2 2017c212 202083c2 9ec2202 c6c22026 a3022034 afb74000 20488602 2053c7c2 20548512 95829672 2063c2e2 e512ee02 20648fb 2074a5e2 2075d5b2 207aa892 ad32207b ed72ee32 00000000 00000000 00000002 00000015 0000000 3de510d4 dd1234dd 00000031 0000009 04000000 00210002 00000002 00000000 92011d80 0000000 ee1234ee 00000009 03010000 00210002 00033dac 920117d5 00000aa8 00000081 00000000 20109ef2 2011ee42 efc22012 93222013 e2822014 97022017 e182201b e0222025 eaa22027 cab22028 80d3202 84b22035 c5c2ccb2 2036ebc2 20389672 20508002 95a22051 d3172056 9ee22057 ef42205b cee2eca2 2060ad62 2061c4a2 2063ddb7 20649542 00000000 00000000 00000002 00000019 00000001 f631054a dd1234dd 00000029 00000009 04000000 00210003 00000002 00000000 92011d80 00000001 ee1234e 00210003 00033dac 920117d5 00000aa8 00000081 00000000 2027d422 203088a 2031d692 20369542 2037ed92 20409c92 ace22044 9a822046 a9e22047 d3422048 8fb2204a 8a12204 e172205b c4872060 8f822065 ea222067 c3f24000 00000000 00000000 00000002 0000000 aeaa0e15 dd1234dd 00000039 00000009 04000000 00210004 00000002 00000000 92011d80 ee1234ee 00000009 03010000 00210004 00033dac 920117d5 00000aa8 00000081 2006af1 2017eb47 201a8e76 2025e6d2 20268fa2 a292202b dff74000 2040a152 20469122 20529182 2060aea 2061c4c2 d722d942 2063c5e2 2064a772 206aa152 206bc322 c7c22070 89d22072 8ad22073 c0b7800 c187c1a7 c1f7c227 c287c2c7 c2e7c3a7 c3c7800f c3f7c417 c497c4d7 c547c5b7 c5e7c637 c657c67 c6b7c727 c767c7a7 00000000 00000000 00000002 00000021 00000001 alfeebf3 dd1234dd 0000002d 00000009 04000000 00210005 00000002 00000000 92011d80 00000001 ee1234ee 00000009 03010000



14

Event Reconstruction

- Going from raw data to analysis objects.
- Important: data and simulation pass through the same reconstruction algorithms.
- Raw data reconstructed into:
 - Tracks
 - Calorimeter deposits
- Which are then reconstructed into "physics" objects:
 - Jets, electrons, muons, taus
 - Photons, missing transverse energy



From hits to physics: tracking



- Efficiently and precisely reconstructing charged particles:
 - Under a non-uniform magnetic field (equations of motion have to be solved numerically)
 - With hundreds to thousands of particles per event.
 - With tight CPU timing constraints.
- Used in almost every element of reconstruction.







Vertex reconstruction

Pile up removal

Jet flavour tagging

Challenge: pile-up



A Z boson decays to 2 muons in an event with 65 (!) additional pile-up collisions.

 $\sqrt{\hat{s}} = 13 \text{ TeV}$

 $<\mu>$ = mean number of interactions per crossing

Challenge: pile-up

Track p_T > 100 MeV

Track p_T > 1 GeV

Track p_T > 5 GeV



A Z boson decays to 2 muons in an event with 65 (!) additional pile-up collisions.

 $\sqrt{\hat{s}} = 13 \text{ TeV}$

 $<\mu>$ = mean number of interactions per crossing

From hits to physics: clustering

- Three-dimensional topological clustering (*topo-clustering*) of individual calorimeter cell signals.
- Algorithm sensitive to the nature of the shower producing the cluster signal:
 - EM showers are more compact, smaller intrinsic fluctuations
 - HADdronic shower have larger shower-by-shower fluctuations and are located deeper in the calorimeter.



(c) All clustered cells







From *hits* to *physics*: **clustering**

- Calorimeter topoclusters are one of the ingredients to jet
 clustering
 - **2.2.5** The anti- k_t algorithm

One can generalise the k_t and Cambridge/Aachen distance measures as [33]:

$$d_{ij} = \min(p_{ti}^{2p}, p_{tj}^{2p}) \frac{\Delta R_{ij}^2}{R^2}, \qquad \Delta R_{ij}^2 = (y_i - y_j)^2 + (\phi_i - \phi_j)^2, \qquad (10a)$$

$$d_{iB} = p_{ti}^{2p}, \qquad (10b)$$





Different techniques to handle pile-up:

- At constituent-level, e.g. subtracting lowpt constituents
- At jet-level, subtracting energy density x jet area, cut on jet timing, etc...

The simulation

Physics analyses at the LHC

The power of factorisation of physics at different energy scales.

Inclusive* cross-section for the production of the **Differential partonic** final state X in the cross-section collision of hadrons h_1, h_2 $\sigma_{h_1,h_2 \to X} = \sum_{a \ b \in [a \ c]} \int dx_a \int dx_b f_a^{h_1}(x_a,\mu_F^2) f_b^{h_2}(x_b,\mu_F^2) \int d\Phi_{ab \to X} \frac{d\hat{\sigma}_{ab}(\Phi_{ab \to X},\mu_F^2)}{d\Phi_{ab \to Y}}$ **Parton distribution** functions (PDFs) at Partons a,b in the PDF factorisation scale μ_F^2 with momentum

* inclusive since no specific kinematic configuration or particle multiplicity is specified

fraction x_a, x_b

From collisions to physics results

Theory and data can be linked through precise simulations of:

- The hard scatter interaction
- The parton showering and hadronization
- The detector itself

The goal is to have a twin set of collision data to compare to the real collisions.



Detector simulation (I)

Simulation of the passage of particles through the detectors.

Particle ionisation in the trackers

Energy deposition in the calorimeters

Intermediate particle decays, radiation and scattering...

Typically done using the *GEANT4* software, taking into account:

Dense hit content in the inner trackers.

Electromagnetic and hadronic shower development.

Effect of the magnetic fields.

Complex geometry with multiple sub-detectors, support structures, cooling pipes, cables, ...



No. of steps ~ simulation time

A problem...



One of the most computational expensive steps in the entire Monte Carlo generation chain: ~40% of ATLAS resources in Rur₂2.

Detector simulation (II)

Can also be done using a *fast* simulation: Parameterising the calorimeter response to single particles (smearing 4-vectors)

New: improved methods using machine learning!



Calorimeter energy response to photons



Fast simulation provides speed gains of O(500) for calorimeter simulation!



Bringing the two together





Overlaying pile-up collisions

- Pile-up comes "for free" in our data. It needs to be modelled in our Monte Carlo as well, for a fair comparison to data.
- In ATLAS Run 3, this is done by *MC pile-up overlay*:
 - *Minimum-bias* and *single neutrino* events are generated using Pythia8.
 - GEANT4 is run on these events to simulate detector response.
 - Digitisation is then run on a combination of these events, including shifts in time to reproduce in-time and out-of-time pile-up.
 - The digitisation output of a pile-up event is then overlaid with the hard-scatter MC.
- Proton bunch structure and luminosity based on that of real data.



 $<\mu>$ = mean number of interactions per crossing 29

Taking from ATLAS Lectures on Data Processing Workflow: slides

What else?

MC WARNING

- The *real world* need to be reflected in the Monte Carlo simulation.
 - E.g. a section of the calorimeter readout dies and cannot be repaired until the detector is opened during an LHC shutdown.
 - If this impacts x% of the data, we need a representative slice of the problem in our MC.
 - But x is usually hard to know until we know how much data we will collect until the shutdown. At that point we need to *reprocess* the MC.
- Even then, MC often doesn't describe the data.
 - Improving MC (e.g. via tuning of input parameters) is an ever on-going (and time consuming task.
 - Another way to deal with inaccurate modeling is to **correct / calibrate** the MC.
 - We can correct:
 - An efficiency (event-level correction).
 - An object's energy scale or resolution (object-level correction).

Example #1: tag & probe method

• How efficiently do we identify an electron?





Use a Standard Model candle like $Z \rightarrow ee$

- ✓ We know the Z decays to one electron and a positron
- ✓ We know the Z invariant mass very well
- ✓ We "tag" one electron and study the "probe" electron

For example, the identification efficiency can be calculated as:



EGAM-2021-01

Example #1: tag & probe method

• How efficiently do we identify an electron?



- Electrons are identified with tracks and EM topo-clusters.
- Id efficiency shown as a function of electron η:
 - Also studied as a function of electron $\mathsf{E}_{\mathsf{T}},$ pile-up...
- Data and simulation have different efficiencies, an approximately 5% effect.
 - Weights or scale-factors are derived as a function of η, E_T, ...
 - We *reweight* the simulation to achieve the same efficiencies as in the data.

Example #2: smearing the MC

• How well do we measure the momenta of muons?



- Muons are typically reconstructed using the ATLAS inner detector and the muon spectrometer.
- · Each detector has its own momentum resolution:



• We *smear the MC* (depending on detector region) to reproduce the muon momentum resolution and scale of data at high precision.

Example #2: smearing the MC

• How well do we measure the momenta of muons?





Example #3: measuring tag rates, fake rates

• How do we identify b-jets?

- b-jets contain the decay particles of longlived b-hadrons and some additional particles
- Key properties:
 - Relatively large b-hadron mass ~5 GeV
 - Significant b-hadron lifetime ~1.5 ps
- This leads to **unique characteristics** that distinguish them from light (u,d,s,g) and to a lesser extent charm (c) jets:
 - A secondary vertex
 - Tracks with large impact parameters
 - Leptons from the b-hadron decay



Example #3: measuring tag and mis-tag rates

• State-of-the-art b-tagging in ATLAS


• How efficiently do we tag a b-jet?



- Use a highly-enriched sample of top-pair events to:
 - Measure the jet flavour composition.
 - Measure the b-tagging efficiency vs jet pT.
- Invariant mass of each of the top systems:
 - $m_{t1} = m_{j1,\ell}$
 - $m_{t2} = m_{j2,\ell}$
- Real top-pair events will have $m_{j1,\ell}, m_{j2,\ell}$ distributions with an upper limit around the top-quark mass of 172.5 GeV
 - In practice smaller due to the undetected neutrino.



• How efficiently do we tag a b-jet?



• And how often do we tag a light or a c-jet instead (*mis-tag*)?



• Scale-factors to correct light mis-tag rate in MC as a function of jet transverse momentum.

40

The grid

1.1.1.19

The grid

- Processing data and simulation poses huge computing, storage and analysis challenges.
- We rely on the World LHC Computing Grid (WLCG), and international organisation of computing centres.
 - Tier-0: the CERN Data Centre where O(100) PB of data are stored on magnetic tapes.
 - Tier 1: 14 large data centres for intensive computing tasks and secondary storage.
 - Tier 2: ~160 smaller processing centres, like universities or labs that can provide storage and computing power for specific analysis tasks.



The grid

- Processing data and simulation poses huge computing, storage and analysis challenges.
- We rely on the World LHC Computing Grid (WLCG), and international organisation of computing centres.
 - Tier-0: the CERN Data Centre where O(100) PB of data are stored on magnetic tapes.
 - Tier 1: 14 large data centres for intensive computing tasks and secondary storage.
 - Tier 2: ~160 smaller processing centres, like universities or labs that can provide storage and computing power for specific analysis tasks.

An ATLAS example:

- DAOD	129 PB
- AOD	80.2 PB
- HITS	51.4 PB
- RDO	17.9 PB
- EVNT	13.0 PB
- RAW	11.9 PB
- ESD	2.27 PB
- DESD	1.52 PB
- log	1.43 PB
- no_name	1.24 PB
- TXT	1.00 PB
- HIST	886 TB
- DRAW	787 TB
- user	697 TB
- NTUP	263 TB

Let's do an analysis...

Physics motivation

• Precise measurements of the cross-section and decays of Higgs boson as a test of Standard Model predictions and probe of New Physics.



Signal characterization



SM Higgs Boson decay modes

Decay channel	Branching ratio
$H ightarrow \gamma \gamma$	2.27×10^{-3}
$H \rightarrow ZZ$	2.62×10^{-2}
$H \to W^+ W^-$	2.14×10^{-1}
$H \to \tau^+ \tau^-$	6.27×10^{-2}
$H ightarrow b ar{b}$	5.84×10^{-1}
$H \to Z \gamma$	1.53×10^{-3}
$H \to \mu^+ \mu^-$	2.18×10^{-4}

~	W boson decay modes	Fraction (Γ_i/Γ)	
\Rightarrow	$\ell^+ \nu$	[<i>b</i>]	(10.86± 0.09) %
	$e^+ u$		(10.71 \pm 0.16) %
	$\mu^+ \nu$		$(10.63\pm~0.15)~\%$
	$ au^+ u$		$(11.38\pm~0.21)~\%$
	hadrons		(67.41± 0.27) %



Signal characterisation



Let's explore other options later in the talk...



- **Background:** it is crucial to correctly estimate the expected background and its uncertainty.
- Common strategy (for many backgrounds):

1. Use Monte Carlo estimate (yields and shapes) during analysis optimisation.

2. Use data to correct and constrain MC estimate.*

* when there is an appropriate **control region**.



☑ What are the (dominant) backgrounds? How can we reduce them?



 $ho_{
m C}$ Can take advantage of different invariant mass of Z and H (if mass resolution allows it) 50 50

☑ What are the (dominant) backgrounds? How can we reduce them?





☑ What are the (dominant) backgrounds? How can we reduce them?



 \square Event selection: regions with high signal efficiency \implies signal region(s)



 $m_J = mass of$ the large-

 \blacksquare Event selection: regions with high background purity \Rightarrow control region(s)



Top CR: events with 1 extra b (outside the large-R jet)

- Background from non-prompt / fake leptons:
 - Non-prompt: from semi-leptonic decay of hadrons or photon conversions.
 - Fake leptons from misidentified jets.
 - Very challenging to model these processes in simulation:
 - Depend strongly on details of physics simulation, often in nonperturbative regions.
 - Depend on modeling of material composition and response.
 - Very low probability for hadronic jets to fake a lepton, yet multi-jet cross section is huge and simulating this effect would be prohibitive.







Non-prompt leptons can be reduced

by requiring *isolated* leptons.

ANA-EGAM-2019-01

CR

MC

 m_T^W

- ☑ Background from non-prompt / fake leptons:
 - Use data-driven methods!
 - E.g. template method:
 - Extract a background *template* from a control region enriched in multi-jet events.
 - Built by inverting the lepton isolation and missing transverse energy requirements.
 - Assumption: shape SR = shape CR
 - Determine its normalisation in fits to the W transverse mass distribution in the signal region.



$$m_T^W = \sqrt{2p_T^{\ell} E_T^{miss}(1 - \cos \Delta \phi(\ell, E_T^{miss}))}$$

☑ Analysis specific: improvements to the invariant mass resolution



- Correct b-jets semi-leptonic decays with muon four-vector.
- ✓ Correct for missing energy from neutrinos.
- ✓ In the ZH(IIbb) channel, a kinematic fit.

In the end, a 42% improvement.



 $m_{J} = m_{bb} = mass of$

the large-radius jet (Higgs candidate)

57

☑ All together now!

 The data, the simulated and datadriven backgrounds, as well as the Higgs boson signal go into a likelihood fit of the signal and control regions, considering theoretical and experimental uncertainties.





☑ All together now

• In this case, a strength parameter of the signal is measured:

$$\mu_{VH} = \frac{\sigma_{\text{meas}}}{\sigma_{SM}} = 0.72^{+0.39}_{-0.36}$$

• We take advantage of the diboson peak for validation, before *unblinding* the data.



- Let's make it more interesting...
 - ✓ Make it a resonance: $W' \to WH$
 - ✓ Make it all hadronic





W', Z': heavier versions of the W and Z bosons

Signal characterization:

- 2 large-radius jets: 1 boosted H→bb jet and 1 boosted
 W→qq jet
- A resonant peak above the multijet background, ~TeV scale





Dominant background: multi-jet production

- Huge cross-section!
- Tagging of boosted Higgs and boosted W bosons rejects a lot of background, but what remains is tricky and expensive to simulate precisely.



We do a fully data-driven background estimation 😭

- In other words, we *interpolate* or *extrapolate* from a background dominated CR into a SR.
- We use data directly (in some cases using MC but only in defining regions or checking assumptions).

Example #1: the ABCD method

- 1. Pick two observables *f* and *g* which are:
 - Approximately statistically independent for the background.
 - Effective discriminators of signal vs background.
- 2. Apply thresholds on these observables to define 4 regions:
 - A: signal region
 - B, C, D: background regions
- 3. If **f** and **g** are independent then the background in A can be predicted from the other three regions:

$$N_A = \frac{N_B N_C}{N_D}$$

 N_i = number of events in region i





See arXiv:2007.14400 for a ML+ABCD method

Example #2: multi-dimensional reweighting with ML

- Let's say we extrapolate the background from control region to signal region.
- We cross-check modelling in validation region and observe discrepancies.
- Then, do a **reweighting**: use one sample with distribution $p_{CR}(x)$ to model sample $p_{VR}(x)$, via a density ratio r(x)

$$p_{VR}(x) = r(x)p_{CR}(x)$$

How do we determine r(x)?



Example #2: multi-dimensional reweighting with ML

• *Likelihood-ratio trick:* a classification model (NN, BDT, ...) trained to discriminate between samples A and B can also estimate their probabilities.

$$r(x) = \frac{p_A(x)}{p_B(x)} = \frac{p(x|A)}{p(x|B)} = \frac{p(A|x)}{p(B|x)} = \frac{h(x)}{1 - h(x)}$$

h(x) is the classifier output



Example #2: multi-dimensional reweighting with ML

Before BDT reweighing After BDT reweighting





When do announce a discovery?

• To find out if our data is compatible with the presence of new physics, we can compute the probability for our observation with no signal present.

- Suppose we observe \mathbf{n}_b background events and \mathbf{n}_s signal events.
- Suppose $n=n_b+n_s$ is distributed according to a Poisson distribution with mean s+b: $D(m,a,b) = \frac{(s+b)^n}{a^{-(s+b)}}$

$$P(n;s,b) = \frac{(s+b)^n}{n!}e^{-(s+b)}$$

• If b=0.5 and n=5, do we claim discovery?

p-value = $P(n \ge 5; b = 0.5, s = 0) = 1.7 \times 10^{-4}$

☑ Interpreting the results in the absence of an excess:





HVT model A: $g_F = -0.55$, $g_H = -0.56$ HVT model B: $g_F = 0.14$, $g_H = -2.9$ HVT model C: $g_F = 0$, $g_H = 1$ *small-radius (large-radius) jets are used in resolved (boosted) events [†] with $\ell = u$, e

71

Excluded mass range [TeV]



m₅ [GeV]

R → VV (semi-leptonic)

 $H^{\pm} \rightarrow W^{\pm}Z$

Eur. Phys. J. C 80 (2020)

ATLAS-CONF-2022-005

Eur. Phys. J. C 81 (2021) 3

 $H \rightarrow ZZ \rightarrow 4I + Ibv$

JHEP 06 (2021) 146

72


m_a [GeV]

m₅ [GeV]

Bonus: anomaly detection

Weak supervision

CERN seminar

DiJet Mass

mixed sample 2 dN/dm_{res} sample background Α **CWoLa** mixed B (jet) signal Classification Without Labels C (jet) mres Features of B, C events SB SR SB #

a ion bels

Mixed Sample 2

Mixed Sample 1

Bonus: anomaly detection

https://arxiv.org/abs/2105.09274v1

https://arxiv.org/abs/2306.03637



75



Stay tuned for Run 3 and beyond!



Backup

e de trop

The LHC datasets



 If one measures the same quantity x many times, always using the same method, and if all sources of uncertainty are small and random, then the results will be distributed around the true value x_{true} determined by the mean of all measurements, and in accordance with the normal, or bell-shaped, curve:



Approximately 68% of measurements will fall within 1σ below or above the true value.

 Alternatively, if one makes this measurement only once (with the same method), there is a 68% probability that the result will be in the range xtrue ±σ.



- Most measurements obey normal distribution statistics.
- If one observes a large deviation from expectation, one can quantify the level of disagreement based on the probability of such observation:



The significance of this particular measurement being in disagreement with expectation is $>3\sigma$.

- Most measurements obey normal distribution statistics.
- If one observes a large deviation from expectation, one can quantify the level of disagreement based on the probability of such observation:



Note: 50 corresponds to 0.000057% probability!

Ingredients for boosted boson tagging



Ingredients for boosted boson tagging



