

Data Analysis

Stefan Wayand
09. June 2016

INSTITUTE OF EXPERIMENTAL PARTICLE PHYSICS (IEKP) – PHYSICS FACULTY



Schedule for today

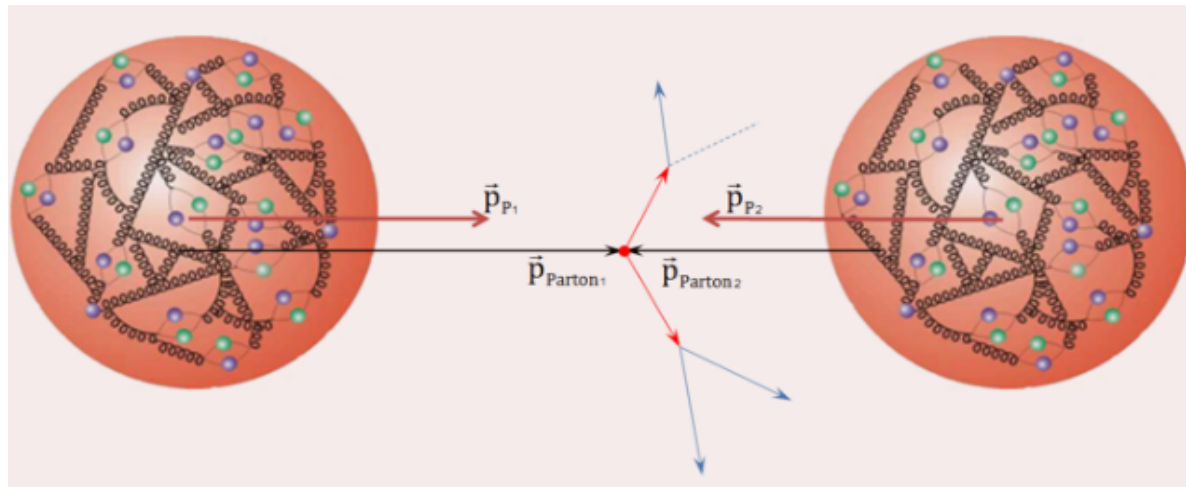
- Which objects can be identified by a particle detector ?
- What tasks are covered by the Analysis?

3 Modeling the background

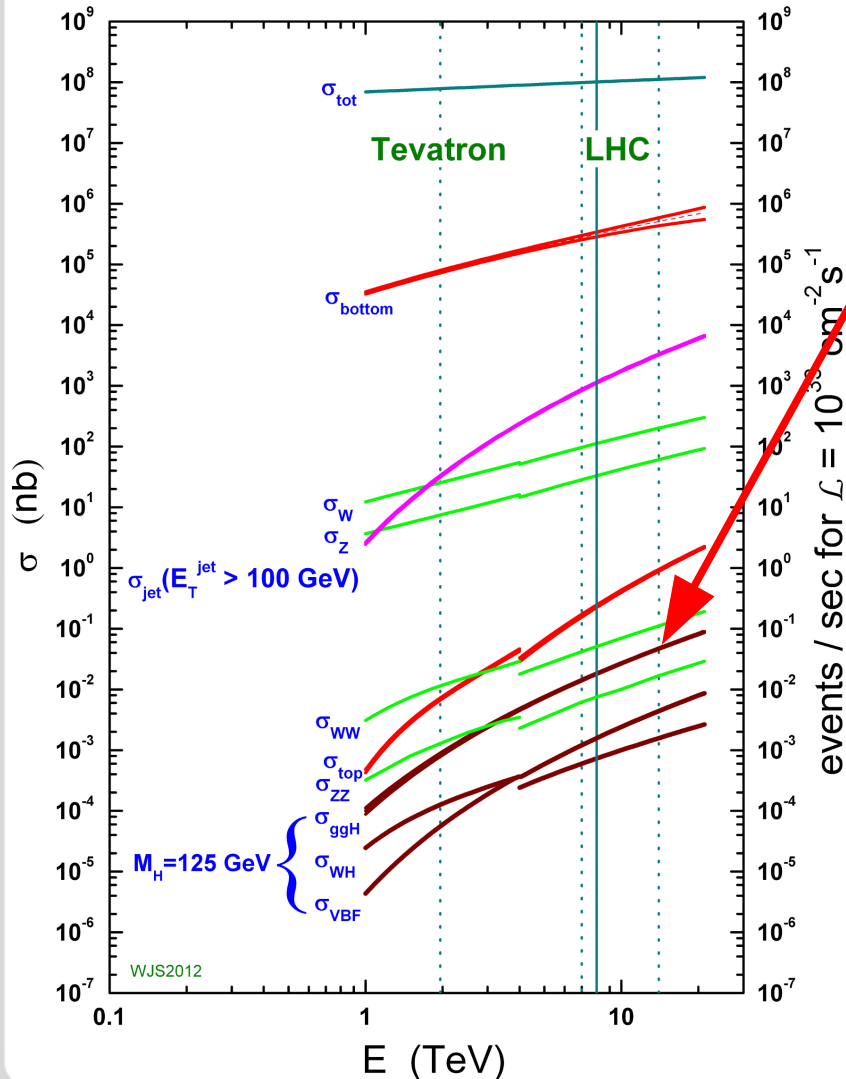
2 Techniques used to understand the reco objects

1 Basics about object reconstruction

Recap



proton - (anti)proton cross sections



Need to understand a large variety of particle physics processes to find **the Higgs**

- Understand the reconstructed objects
- Search in well defined final states ($H \rightarrow bb/\tau\tau/WW/ZZ/\gamma\gamma$). Choose your triggers
- Define the search region (optimize signal to background ratio):
cuts / shapes / MVA
- Model the background processes and estimate the signal yields

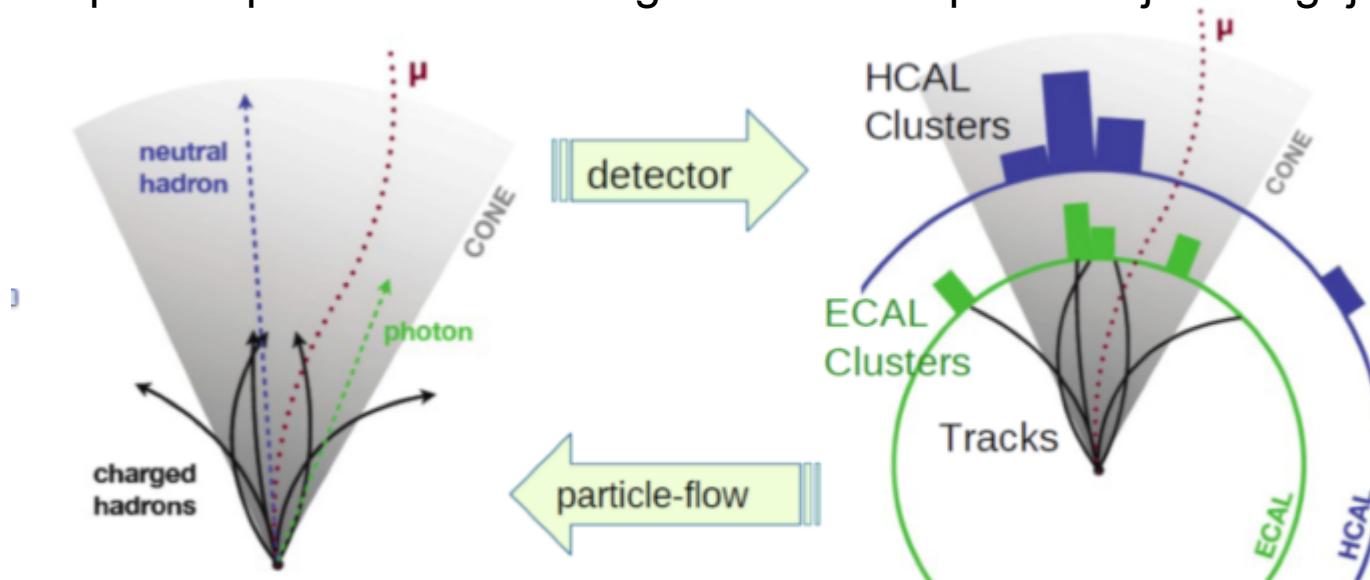
Feed into your statistical model to quantize the result

High-level reconstruction: Particle Flow

- Attempts to reconstruct and identify all particles in the event
→ need matching between calorimeter (fine granularity ECAL) and tracker
- Optimally combines information from all sub-detectors to give best four-momentum measurement of each particle type:

Charged hadrons, neutral hadrons, electrons, photons and muons

- Also improves performance for higher-level composite objects e.g. jets, MET



Reconstruction of Objects

1. combine sub-detectors to classify all stable objects, i.e. find electrons, muons, photons, hadrons. (In CMS provided by the “particle flow” algorithm)

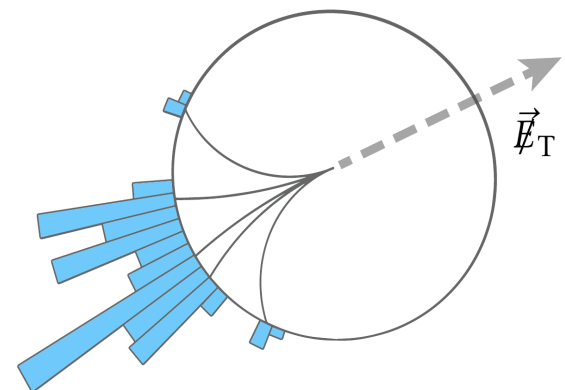
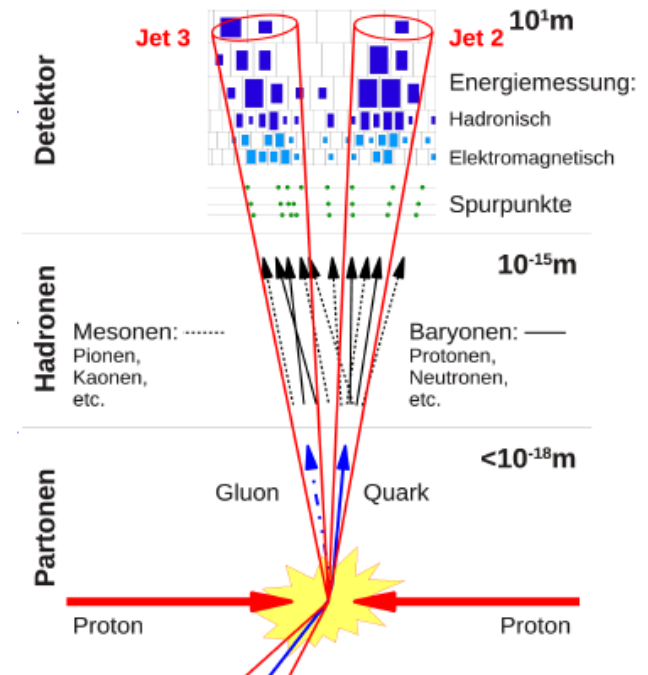
2. cluster objects into “jets” (relation between measured final state objects & hard partons) two types of algorithms:

1. **“cone”**: geometrically assign objects to the leading object
2. **sequentially combine** closest pairs of objects – different measures of “distance” exist (kT, anti-kT) with some variation of resolution parameter, which determines “jet size”

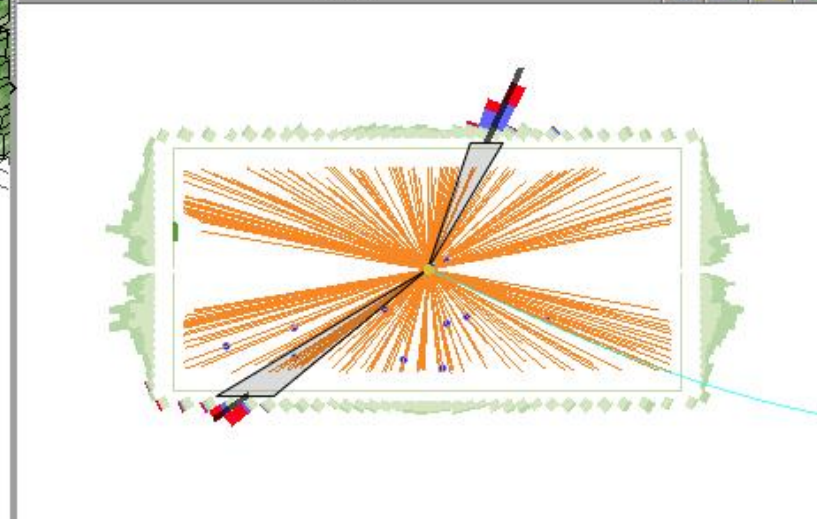
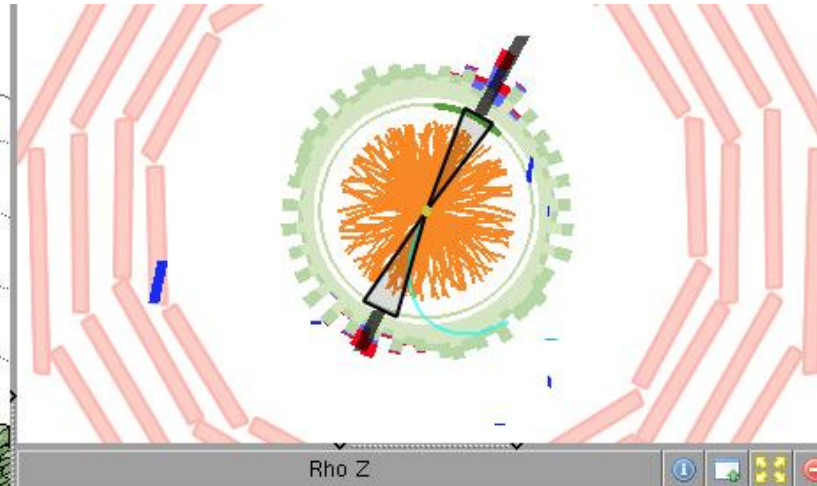
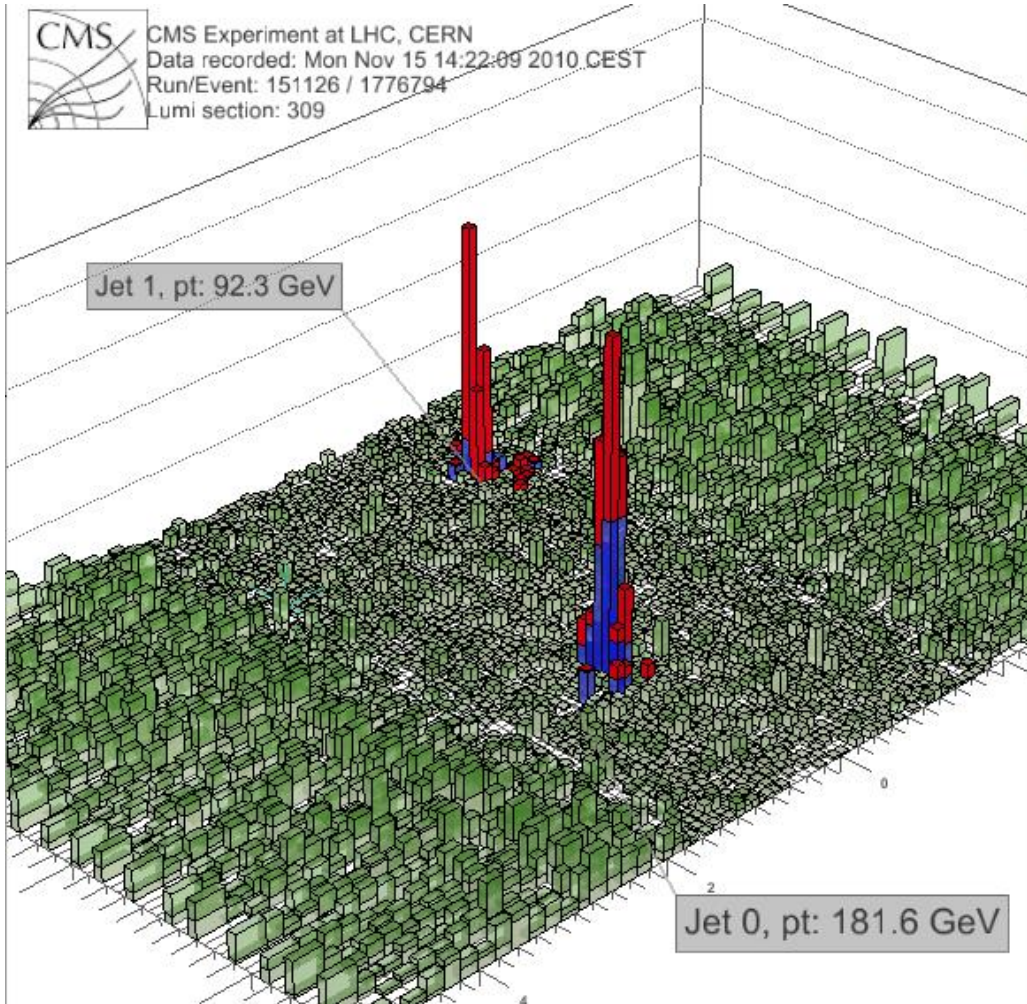
3. determine missing transverse momentum (energy) called MET:

$$p_{T\text{miss}} = - \sum_{\text{all particles}} \vec{p}_{Ti}$$

carried away by undetectable particles. In SM neutrinos, “new physics” provides more of them (e.g. dark matter)



Two-Jet event in the CMS Detector

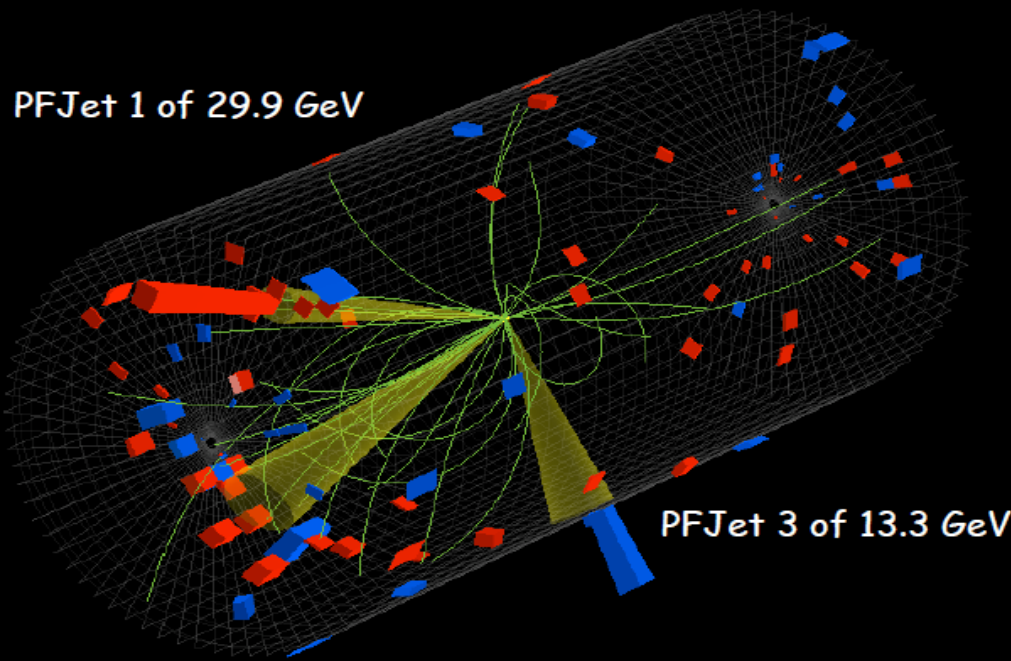


Three-Jet event in the CMS Detector



CMS Experiment at the LHC, CERN
 Date Recorded: 2009-12-14 04:21:03 CEST
 Run/Event: 124120/542515
 Candidate multijet event at 2.36 TeV

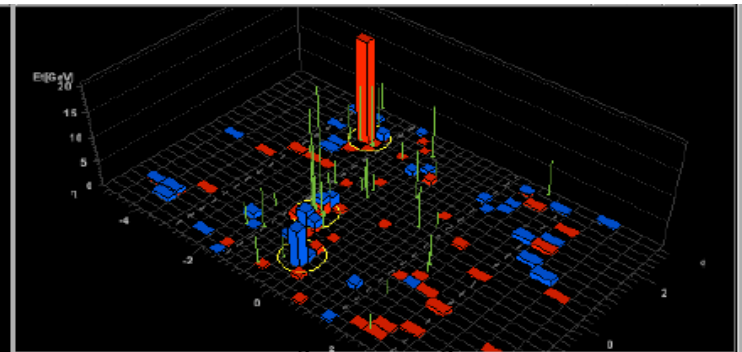
PFJet 1 of 29.9 GeV



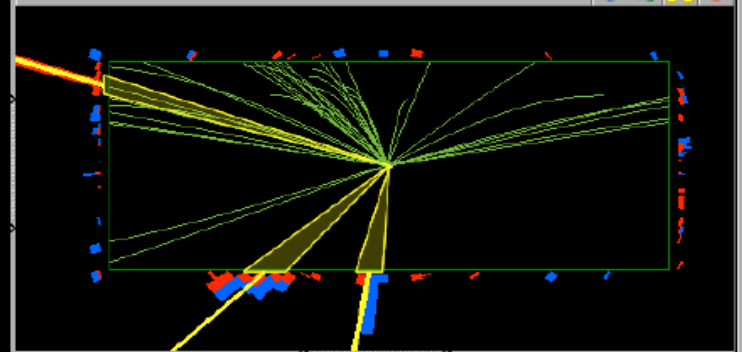
PFJet 3 of 13.3 GeV

PFJet 2 of 24.2 GeV

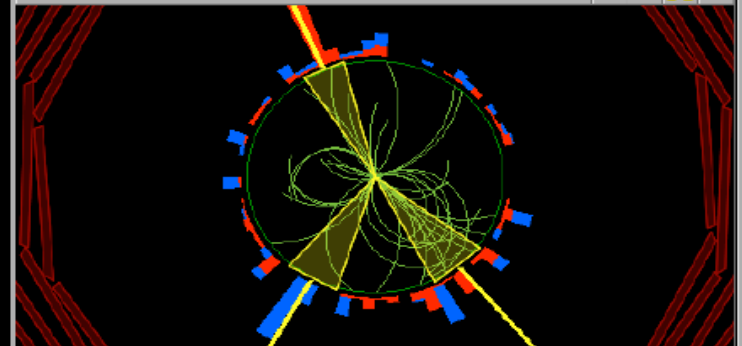
3 PFlow jets $p_T > 10$ GeV
 p_T cut on tracks displayed > 0.4 GeV



Rho Z



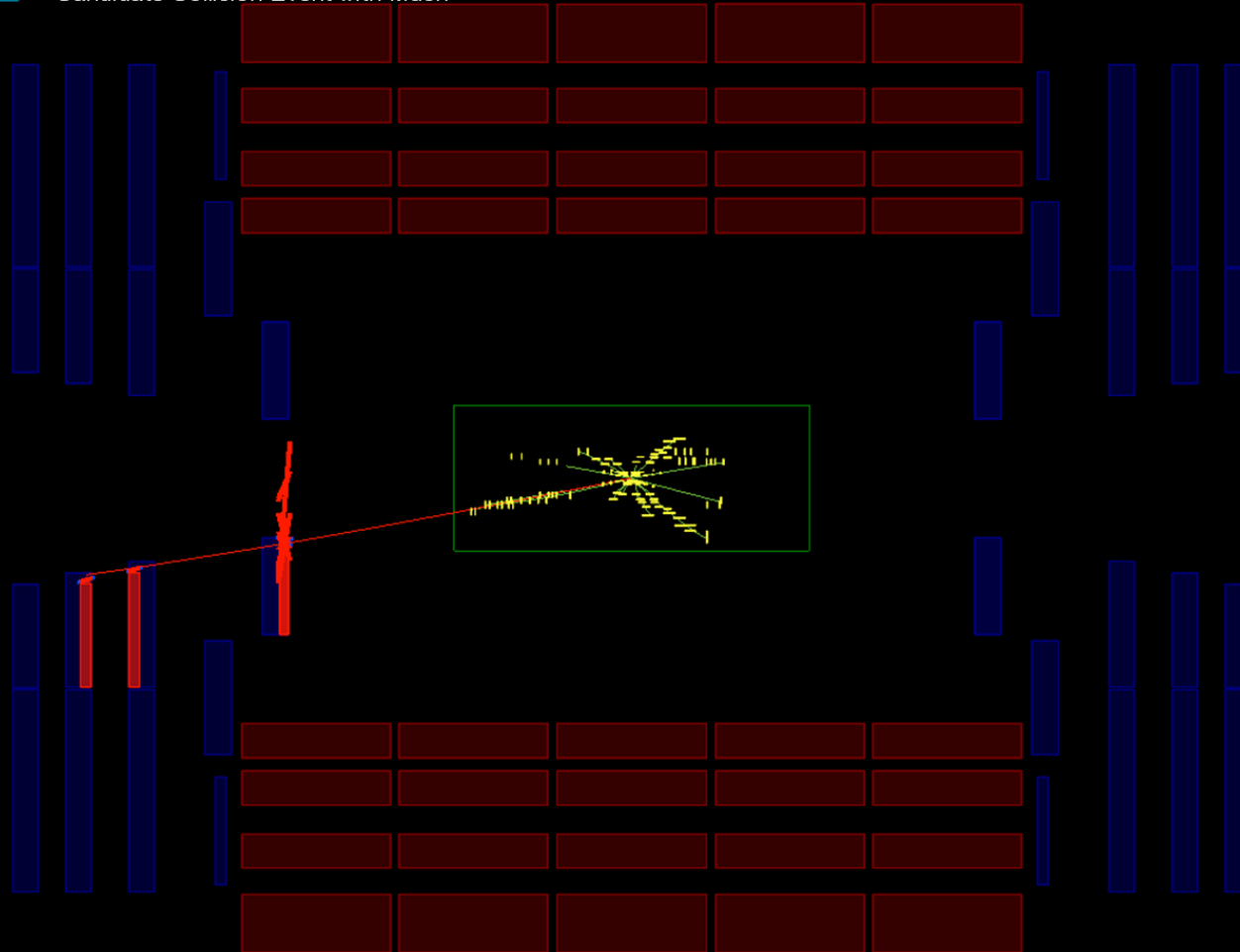
Rho Phi



Event with an end-cap muon



CMS Experiment at the LHC, CERN
Date Recorded: 2009-12-06 05:07 CET
Run/Event: 123592 / 1231789
Candidate Collision Event with Muon



Two electrons in the CMS Detector

File Edit View Window Help

Run 136033 Event 99386647 Sat May 22 07:54:30 2010 CEST
 Delay 5.0s Event filtering is OFF Lumi block id: 785

Summary View

- ECal
- HCal
- Jets
- Tracks
- Muons
- Electrons
- Vertices
- DT-segments
- CSC-segments
- Photons
- MET
- pTMet

Views

Rho Phi

3D Lego

Rho Z

Table

Collection Muons

pT	global	tracker	SA	calo	tr pt	eta	phi	matches	d0	d0 / d0Err	charge

Table

Collection pTMet

MET	phi	sumEt	mEISig
4.1	-2.169	187.6	0.302

Table

Collection Jets

Pt	eta	phi	ECAL	HCal	emf	size_eta	size_phi
43.1	-1.070	1.492	69.5	0.7	0.389	0.014	0.053
41.1	-1.802	-1.621	127.9	0.0	1.000	0.011	0.035
5.0	0.079	0.890	2.9	2.3	0.557	0.102	0.193
2.6	2.003	0.458	1.2	8.7	0.126	0.072	0.173
2.1	3.029	-1.224	5.7	15.9	0.264	0.166	0.087

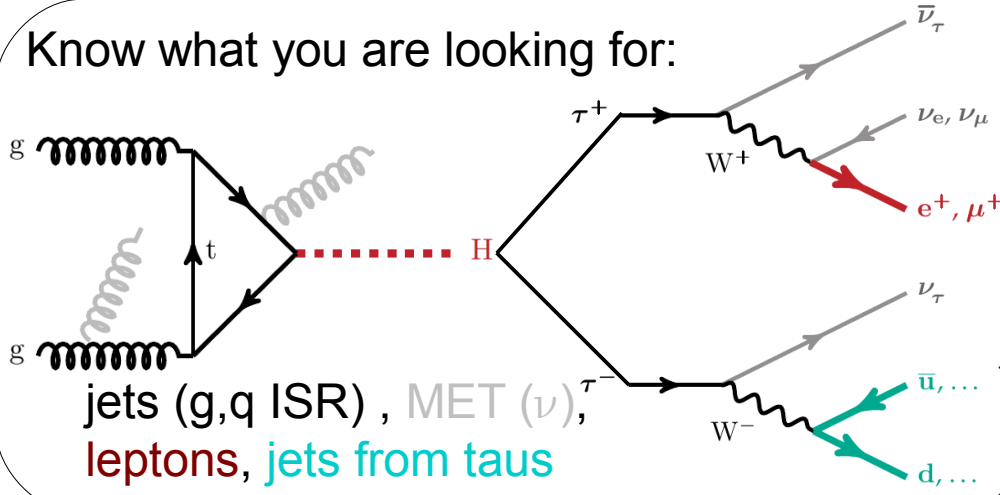
Table

Collection Electrons

pT	eta	phi	E/p	H/E	fbrem	dei	dpi	charge
46.5	-1.803	-1.607	0.988	0.000	0.872	-0.007	-0.006	1
43.1	-1.074	1.472	1.508	0.000	-0.448	0.007	-0.008	-1

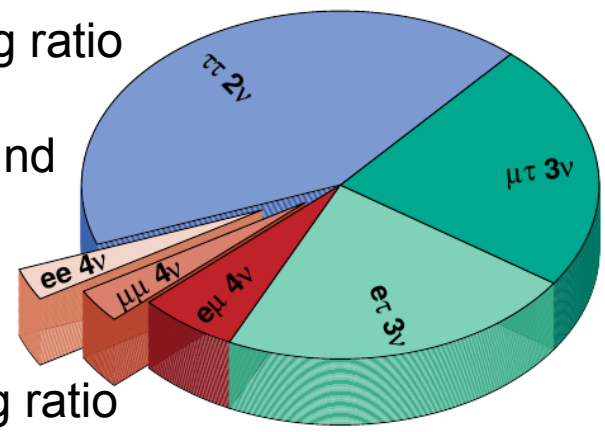
Start the Analysis

Know what you are looking for:



Decay Mode	Resonance	Branching Ratio / %
$\tau^- \rightarrow e^- \bar{\nu}_e \nu_\tau$		17.8
$\tau^- \rightarrow \mu^- \bar{\nu}_\mu \nu_\tau$		17.4
$\tau^- \rightarrow \pi^- \nu_\tau$	$\pi(140)$	11.6
$\tau^- \rightarrow \pi^- \pi^0 \nu_\tau$	$\rho(770)$	26.0
$\tau^- \rightarrow \pi^- \pi^0 \pi^0 \nu_\tau$	$a_1(1260)$	10.8
$\tau^- \rightarrow \pi^- \pi^+ \pi^- \nu_\tau$	$a_1(1260)$	9.8
$\tau^- \rightarrow \pi^- \pi^+ \pi^- \pi^0 \nu_\tau$		4.8
Other hadronic modes		1.7
All hadronic modes		64.8

- largest branching ratio
- hard to trigger
- largest background



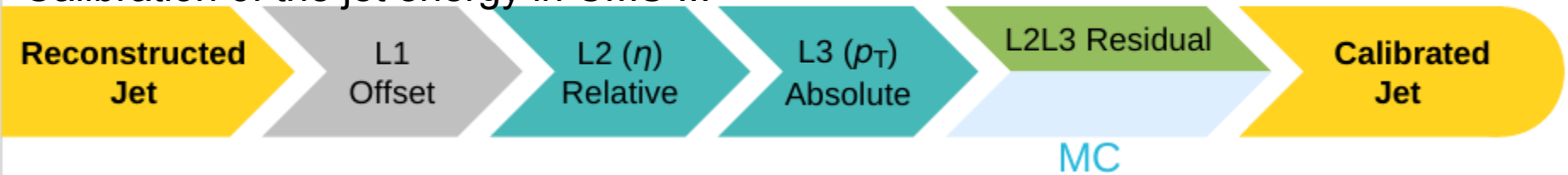
- good branching ratio
- moderate trigger thresholds

- smallest branching ratio
- smallest background (e mu)

In the final analysis all final states were considered (except ee/mu mu)

Object calibration

Calibration of the jet energy in CMS ...



... is a multi-step procedure, driven by data and MC

Level 1: offset correction for pile-up and electronic noise

Level 2: relative (η) corrections

Level 3: absolute p_T correction

MC and special balanced events

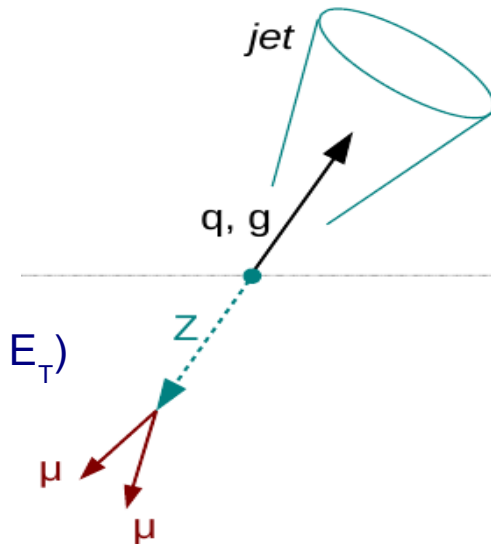
residual corrections from events with selected topology:

Level 2 residual η

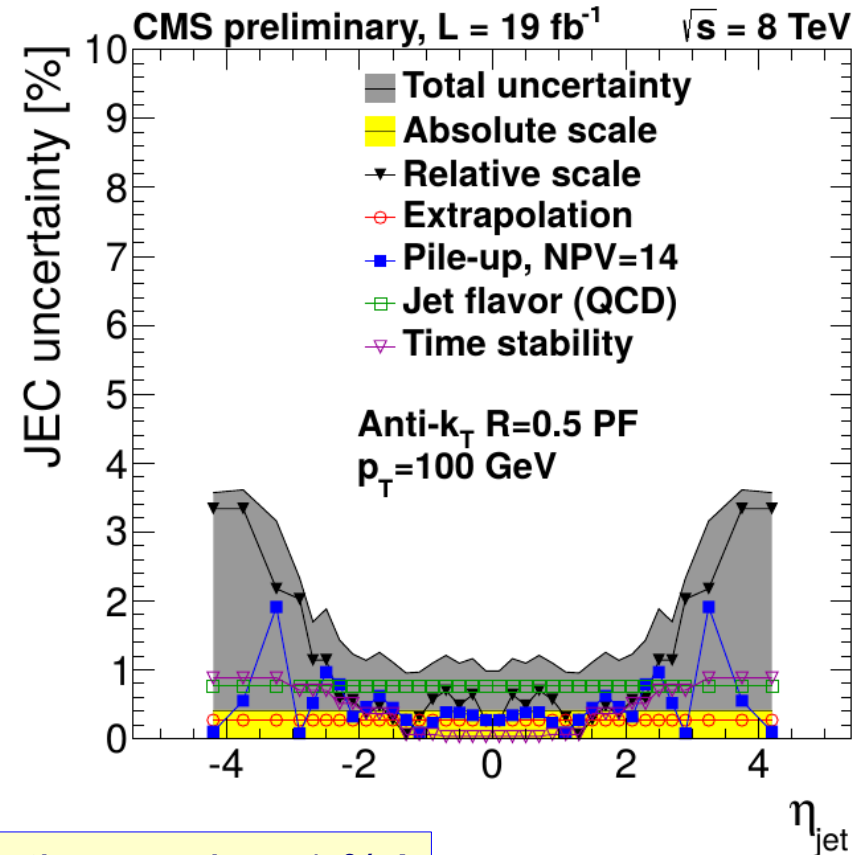
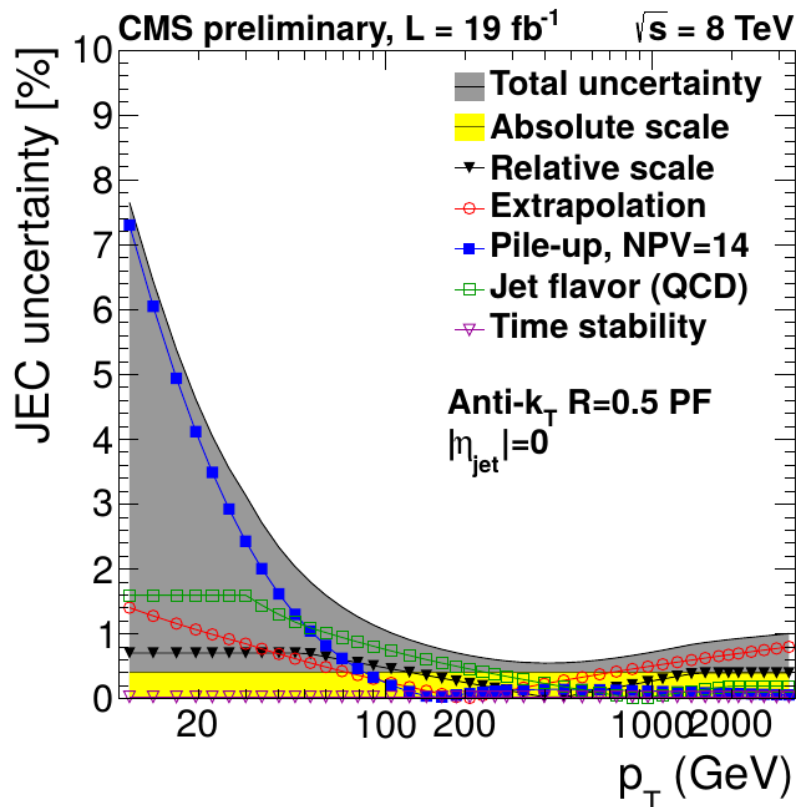
from measured di-jet events, assuming the two jets have the same E_T

Level 3 residual p_T

from measured Z+jet & photon+jet, jet balanced by Z/ γ



Object calibration (Jets)

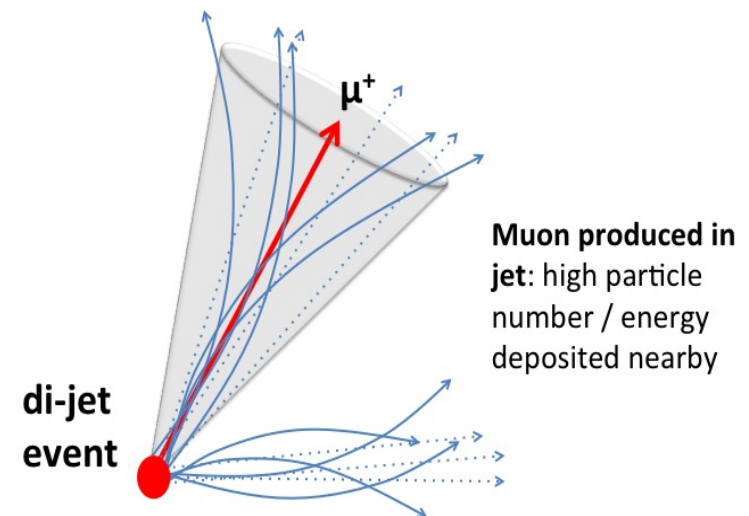
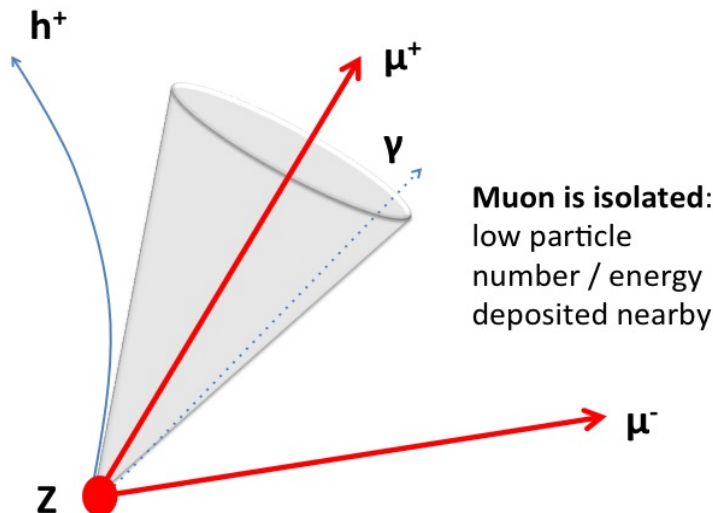


Precision of Jet energy calibration reaches 1 % !

Result is also propagated into MET which helps to improve MET resolution

Object identification and object isolation

- **Identification:** The true particle type can be ambiguous
 - “Is it an electron or a pion?” → can apply object criteria to increase purity of a particle type, e.g. small hadronic energy / EM energy → more likely to be an electron
- **Isolation:** powerful handle to reduce background from jets
 - We are often interested in leptons produced from decays of top quarks, W bosons, Z bosons, Higgs etc
 - These electroweak processes are 'clean' compared to QCD → less activity in the region around lepton direction

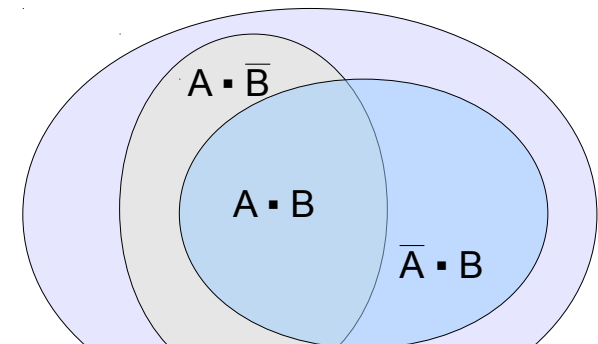


Determination of efficiencies

1. **take efficiencies from simulation** not always believable !
 check classification in simulated data vs. truth, i.e. determine
 ϵ_{MC} = fraction of correctly selected objects
 (probability to select background determined in the same way)
2. **design data-driven methods** using redundancy of at least two variables discriminating signal and background
 - **tag & probe method:**
 select very hard on one criterion, even with low efficiency,
 check result obtained by second criterion

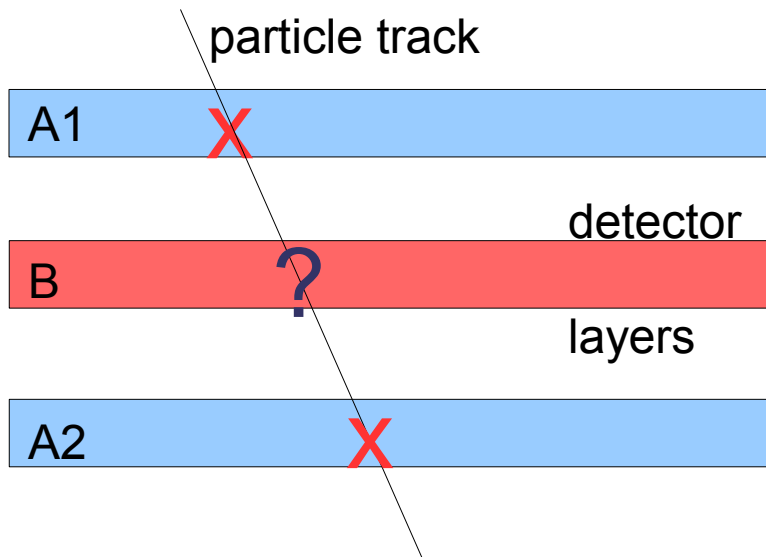
Illustration: two **independent** criteria A, B

$$\epsilon_B = \frac{n(A \cdot B)}{n(A \cdot B) + n(A \cdot \bar{B})}$$



Important: selecting on A must not affect B, i.e. A and B must be uncorrelated !

Tag and Probe: Example 1



Hits in layers A1 and A2 define valid particle track (tag)

probe hit in layer B

Coincidence of Layers A1 and A2 guarantees high purity of the tag (protects against random noise)

allows determination of efficiency of layer B

$$\Rightarrow \epsilon_B = \frac{n_B}{n_{A1 \cdot A2}}$$

Trigger efficiencies

Determination of trigger efficiencies depends on
existence of independent selection methods

Important to ensure redundancy when building trigger systems !

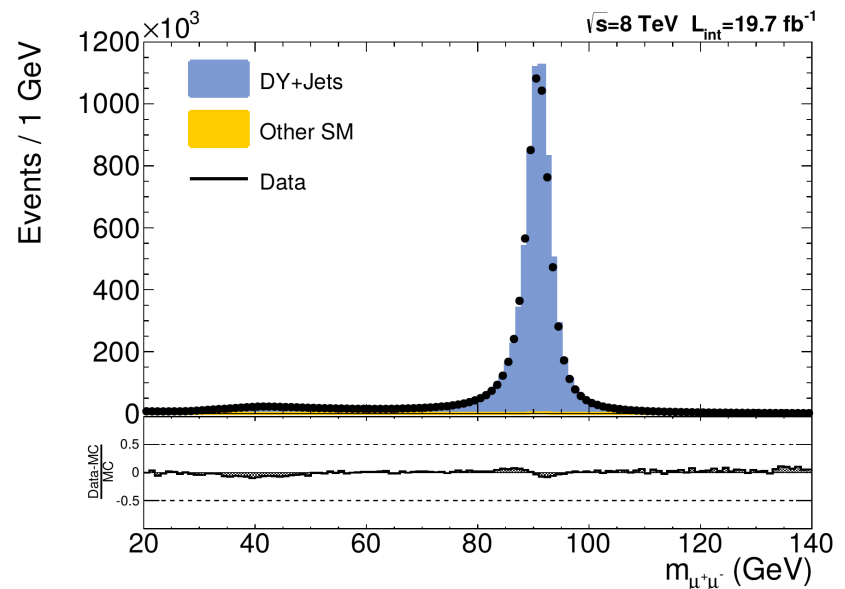
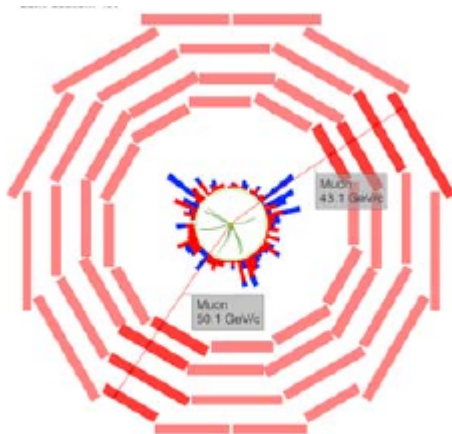
Trigger information must be stored for later use in efficiency determination !

typical methods:

- use trigger from independent sub-systems
- trigger at lower threshold (typically pre-scaled to run at acceptable rates)
 - probe higher-threshold triggers
- trigger on pairs of objects at low threshold,
 - probe higher threshold on each member of the pair
 - !!! potential bias, because higher-threshold trigger depends on same input signals as the tag !!!
- trigger only one object of a pair and use an off-line criterion to identify 2nd member of the pair and probe trigger decision on it

Tag and Probe: Example 2

- criterion A:** a tight muon/electron and one other track with tight selection on Z mass (“tag”) thus selecting $Z \rightarrow \mu\mu$ (or ee) (which is possible with very high purity) \rightarrow 2nd track also is a muon/electron with very high probability
- criterion B:** 2nd track selected by trigger (or analysis) (“probe”) allows measurement of trigger efficiency (or selection efficiency) of second muon



Statistical error on efficiency

determination of efficiencies is a clear application of **binomial statistics**:
number of successes k in n trials at probability p per trial

Binomial Distribution

$$P(k; p, n) = \binom{n}{k} p^k (1-p)^{n-k}, k = 1, \dots, n \quad \binom{n}{k} = \frac{n!}{k!(n-k)!}$$

Expectation value $E[k] = np$

Variance $V[k] = np(1-p)$

Error on efficiency: insert measured efficiency $\epsilon = k/n$ in formula for variance
 (instead of true (but unknown) selection efficiency p !)

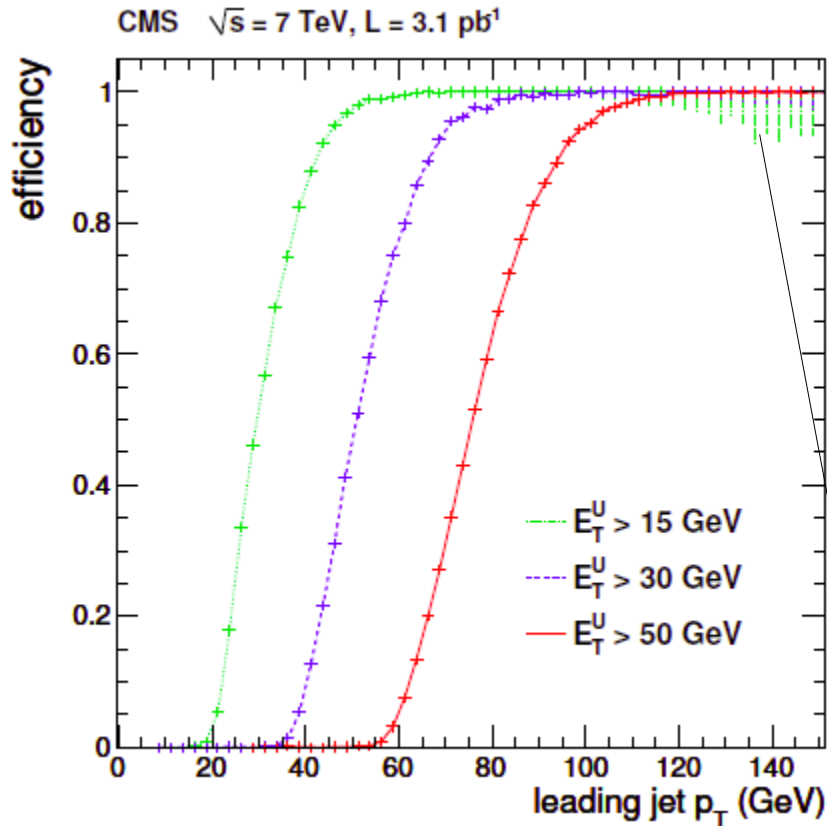
$$\rightarrow \sigma_\epsilon = \frac{\sqrt{\epsilon(1-\epsilon)}}{\sqrt{n}}$$

if this is not justified due to very small statistics, a more sophisticated method of "interval estimation" is needed to specify a confidence range on the measured efficiency:

→ Clopper-Pearson method

Typical “turn-on” curves of trigger efficiencies

(calorimeter jet trigger on transverse energy of jets, CMS experiment)



Remarks:

- efficiency at 100% only far beyond “nominal” threshold
- trigger efficiencies vary with time (depend on “on-line” calibration constants)
- to be safe and independent of trigger efficiencies, analyses should use cuts on reconstructed objects that are tighter than trigger requirements

2nd remark: errors determined as 68% confidence interval by application of Clopper-Person method per bin; this explains the (counter-intuitive) large uncertainties on the >15 GeV trigger at high p_T : there were just no events observed where trigger was inefficient.

LESSON: sophisticated methods are not always plausible !

More complicated observables

Calculate **derived quantities** from objects,

– transverse momentum or energy, *at hadron colliders where rest system of an interaction is boosted along z direction*

$$\vec{p}_T = \sum_i \vec{p}_{T_i}$$

$$E_T^2 = \sum_i m_i^2 + p_{T_i}^2$$

– missing transverse momentum, from all particles in an event, *assuming total transverse momentum of zero in each event, measures effects of invisible particles (neutrinos in the SM, but there are others in extended theories)*

$$\vec{p}_{T \text{ miss}} = - \sum_{\text{all particles}} \vec{p}_{T_i}$$

– “transverse mass” $M_T^2 = E_T^2 - p_T^2$

More complicated observables

Calculate **derived quantities** from objects,

– transverse momentum or energy, *at hadron colliders where rest system of an interaction is boosted along z direction*

$$\vec{p}_T = \sum_i \vec{p}_{T_i}$$

$$E_T^2 = \sum_i m_i^2 + \vec{p}_{T_i}^2$$

– missing transverse momentum, from all particles in an event, *assuming total transverse momentum of zero in each event, measures effects of invisible particles (neutrinos in the SM, but there are others in extended theories)*

$$\vec{p}_{T \text{ miss}} = - \sum_{\text{all particles}} \vec{p}_{T_i}$$

– “transverse mass” $M_T = \sqrt{2 \cdot E_T^{\text{miss}} \cdot p_T^{Wlep} (1 - \cos \Delta\phi)}$

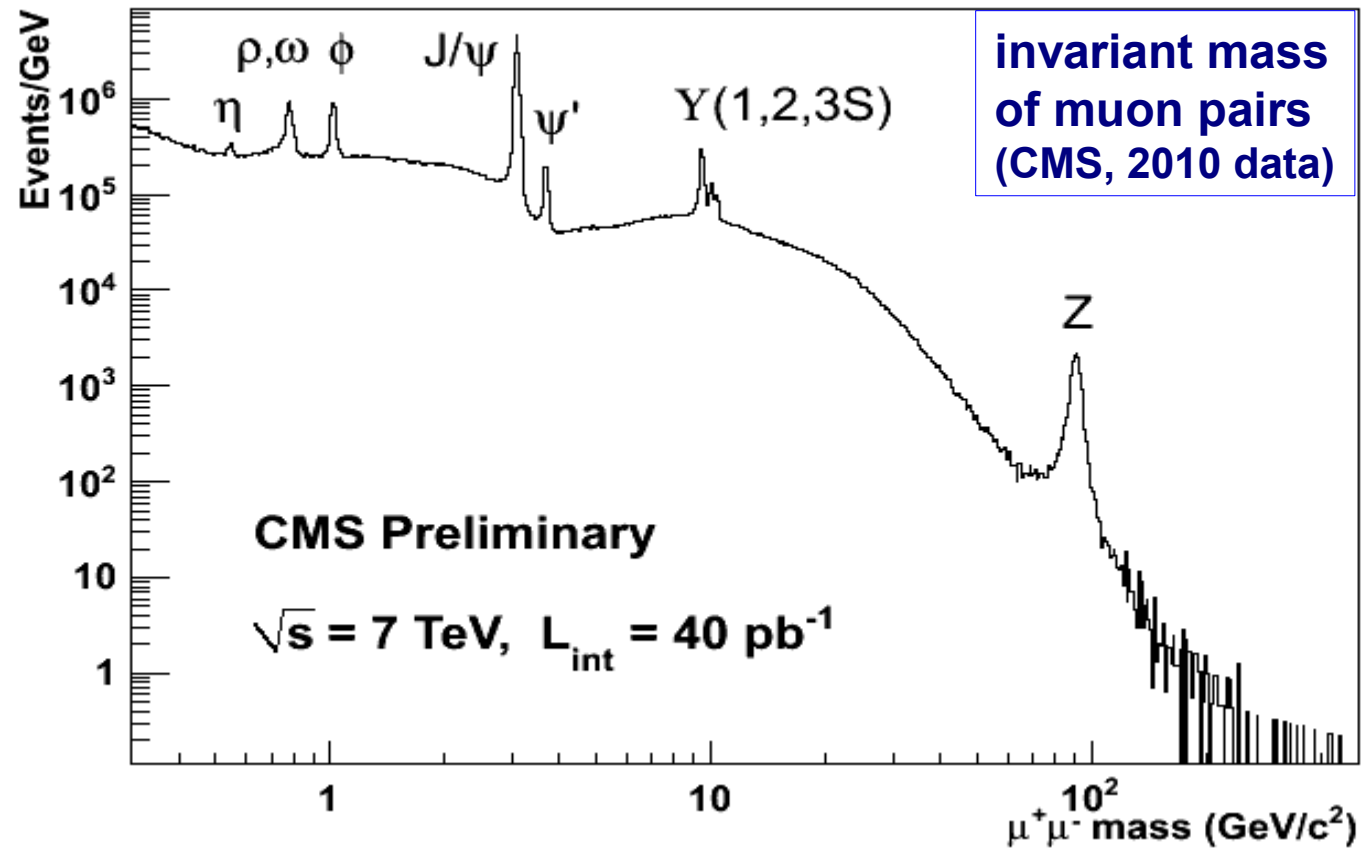
– event shape variables (for QCD analyses) *to classify jet topologies*

– all kinds of “classifiers” using MVA techniques *for object or event classification*

Invariant mass

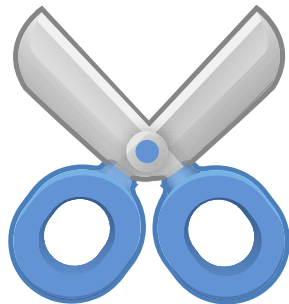
60 years of particle physics in only one year:

Example of a very simple selection:
just the invariant mass of muon pairs in events with one muon trigger



Event Selection

CUTS



lepton p_T, η

lepton isolation

muon identification

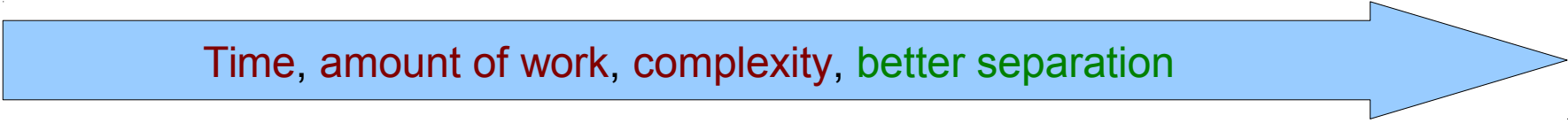
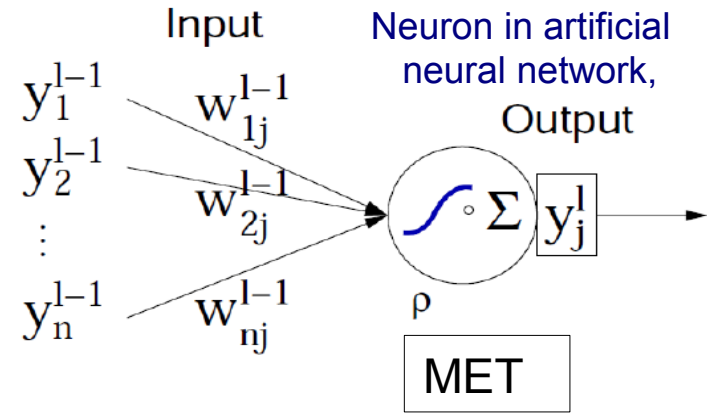
tau identification

electron identification

Number of objects
(e, μ , τ , jets)

invariant mass off di-tau system

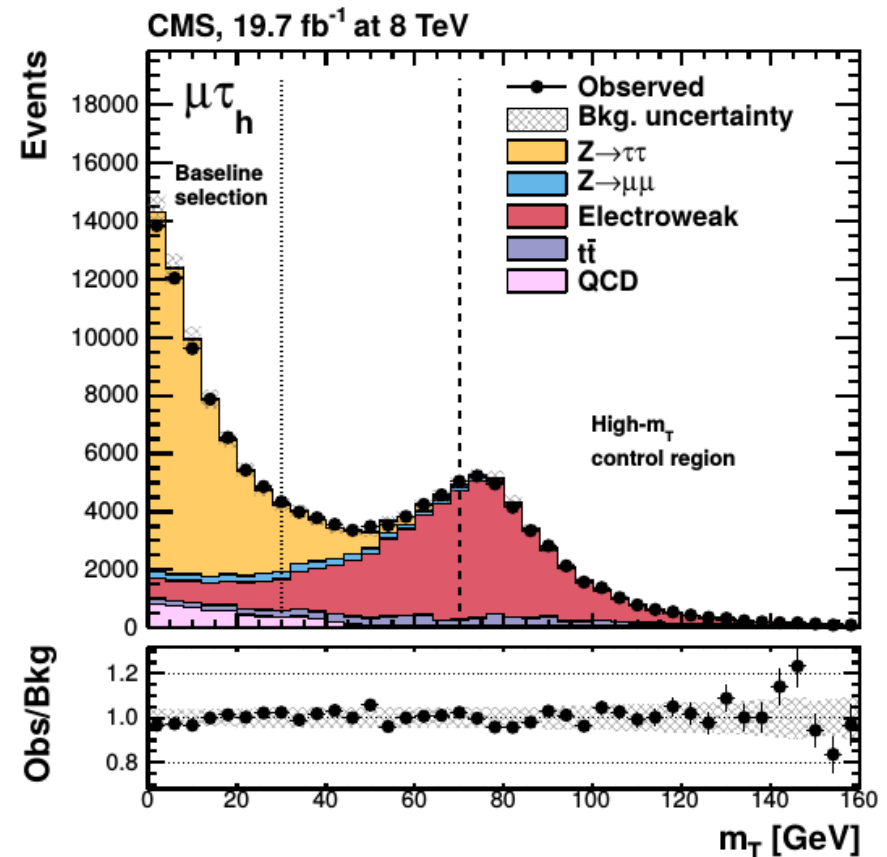
Multivariate Analysis (MVA)
e.g. decorrelated likelihood, artificial neural networks, boosted decision trees



Need to understand the efficiencies on signal and background, the uncertainties and possible correlations

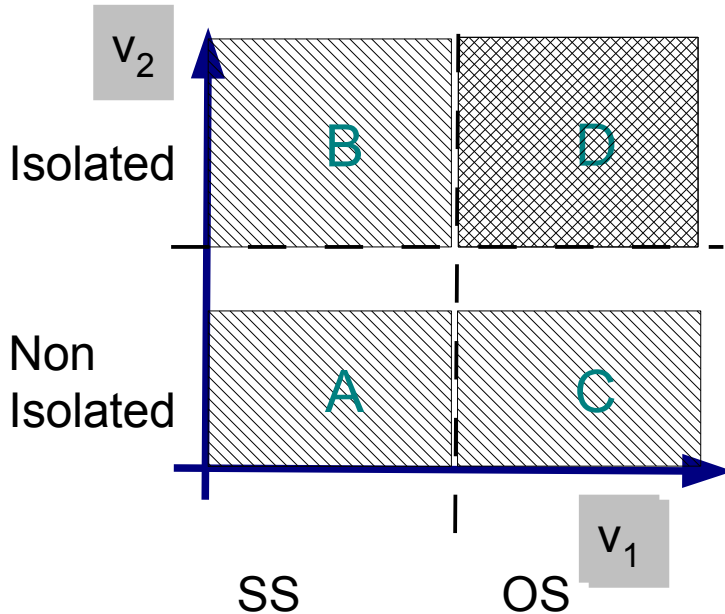
Modeling of Background: part I

- shape take from MC
- extrapolation from “side band” assuming “simple” background shape or by taking background shape from simulation
- event counting in background regions, extrapolation under signal assuming (simple) model
- fit of signal + background model to the observed data



Modeling of Background: part II

– ABCD – Method ...



Assumptions:

- two independent variables v_1 and v_2 for background
- signal only in region D

$$\rightarrow n_D^{bkg} = n_C \frac{n_B}{n_A}$$

... a **data driven estimate** of background under a signal

Example: Take the ratio of same-sign (A) and opposite-sign (B) non isolated (invert isolation criteria) leptons to predict the amount of QCD fakes.

– **more advanced methods** exist to **exploit two uncorrelated variables** to predict the background shape under a signal, see e.g. “sPlot method” in ROOT.

Modeling of Background: part III

Hybrid events: data + Monte Carlo: $Z \rightarrow \tau\tau$ background in the $H \rightarrow \tau\tau$ search

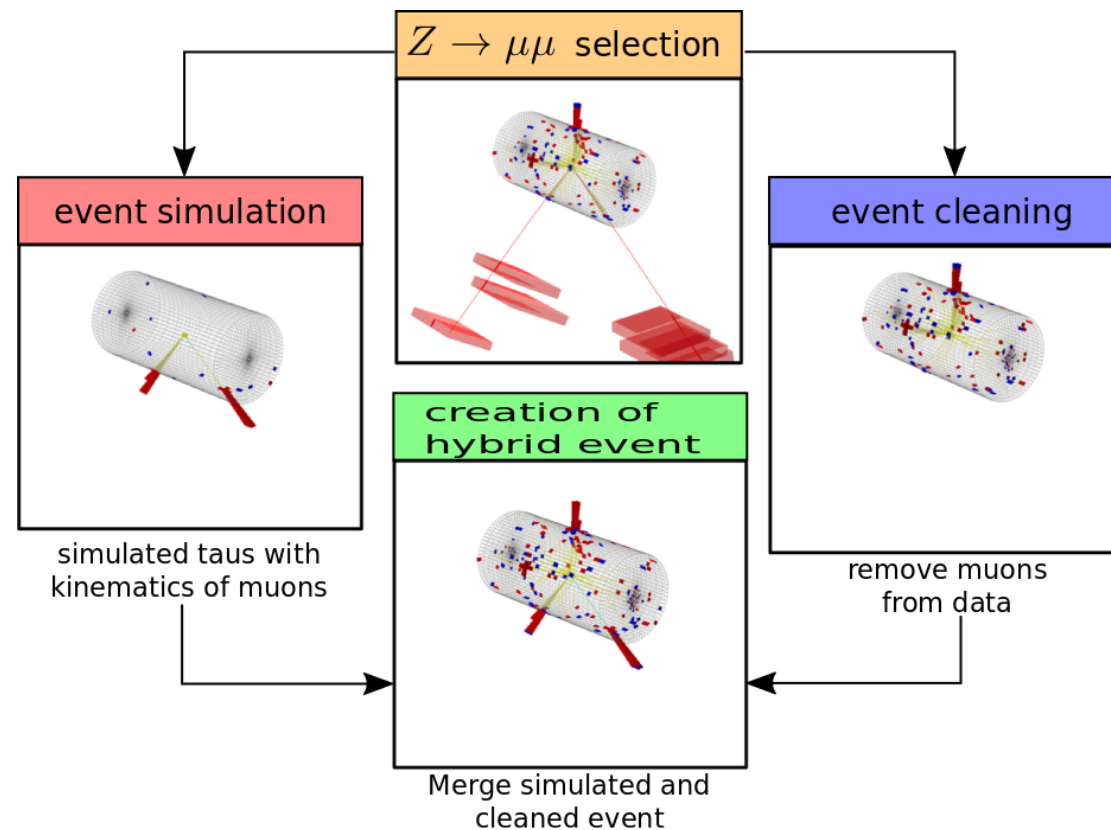
- $H \rightarrow \mu\mu$ has very low branching ratio, hence there is no $H \rightarrow \mu\mu$ under $H \rightarrow \mu\mu$
- $Z \rightarrow \mu\mu$ and $Z \rightarrow \tau\tau$ are very similar (lepton universality of weak decay)

idea:

replace real μ in $Z \rightarrow \mu\mu$ events with simulated τ to model Z background under H signal

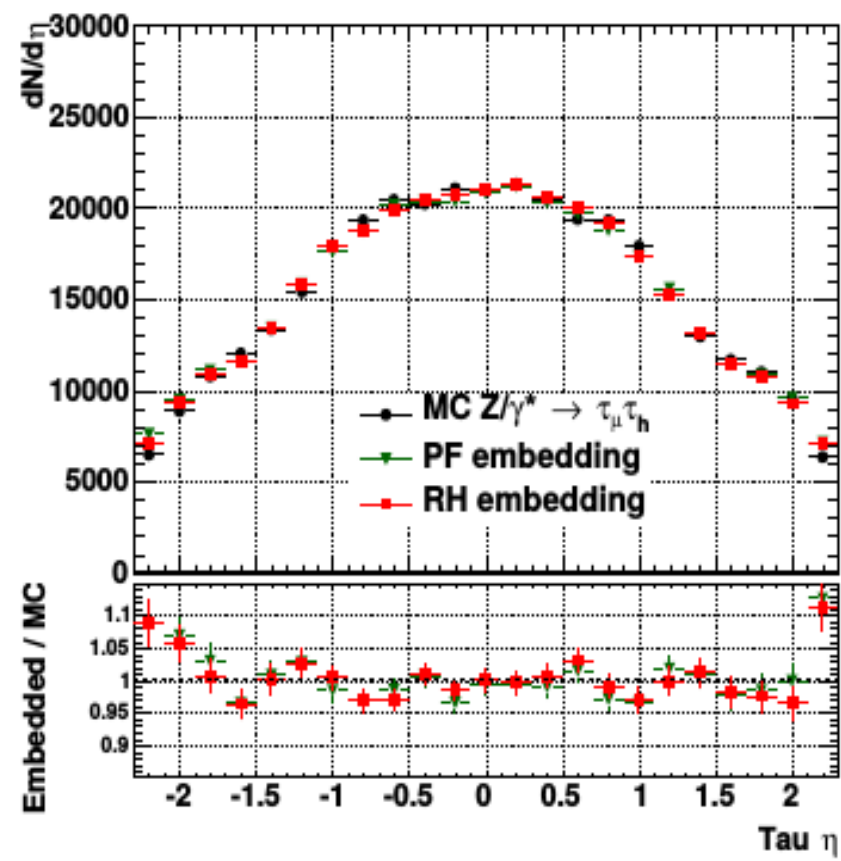
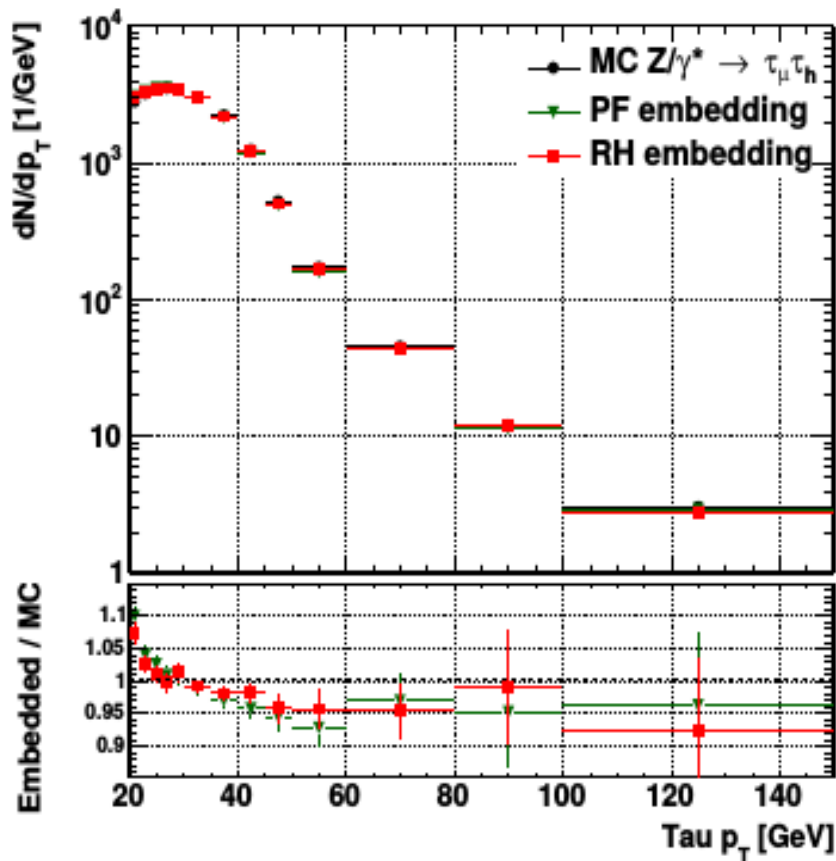
advantages:

- non-leptonic part of event is from real data, esp. important in presence of pile-up
- leptonic part can be well and easily modeled
- important cross check of full simulation via MC



“Closure Test”

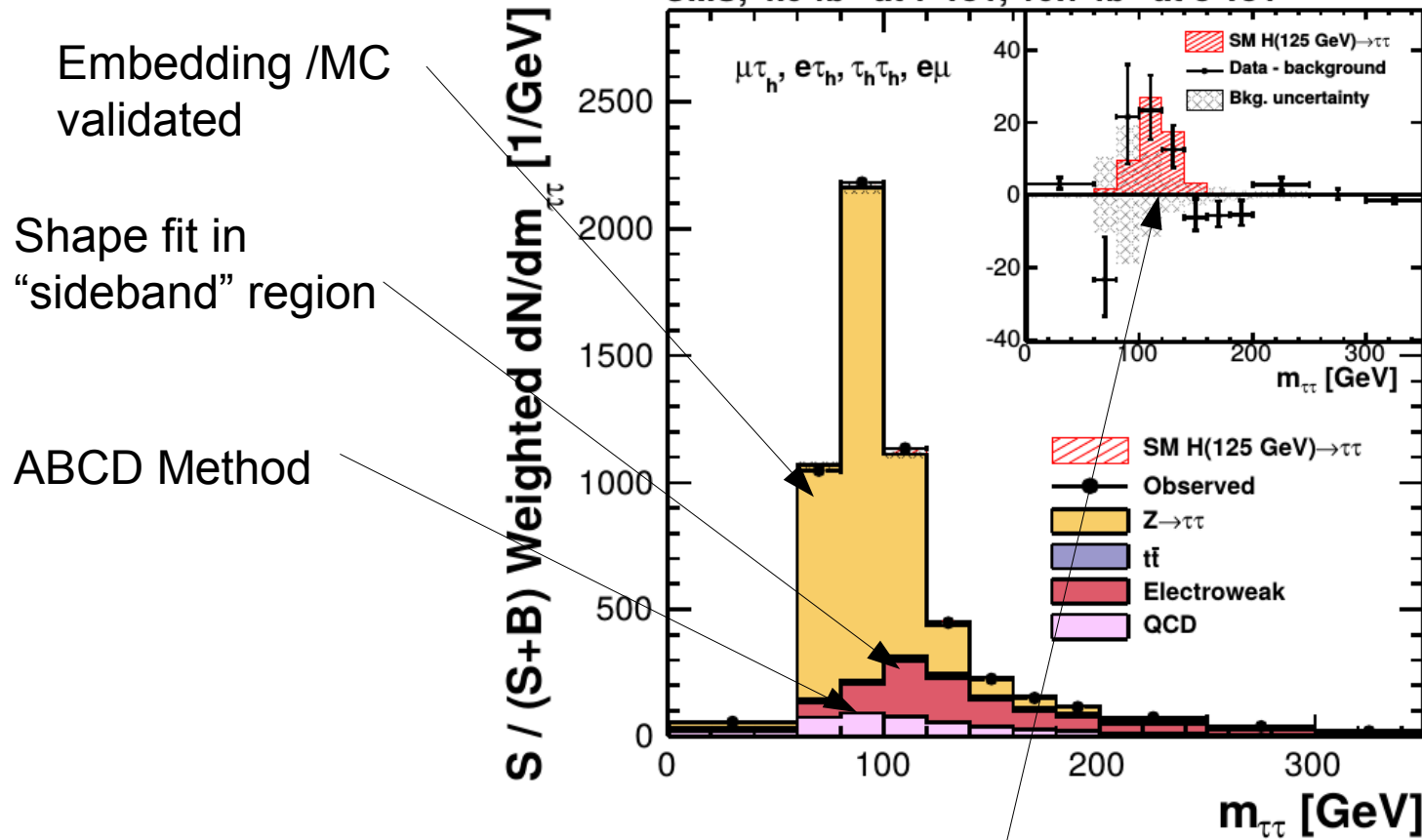
demonstrate that method works on simulated events



from PhD thesis Armin Burgmeier, Karlsruhe - DESY, June 2014

Summary and Outlook

CMS, 4.9 fb⁻¹ at 7 TeV, 19.7 fb⁻¹ at 8 TeV



Coming next:
statistical analysis of rare signals