

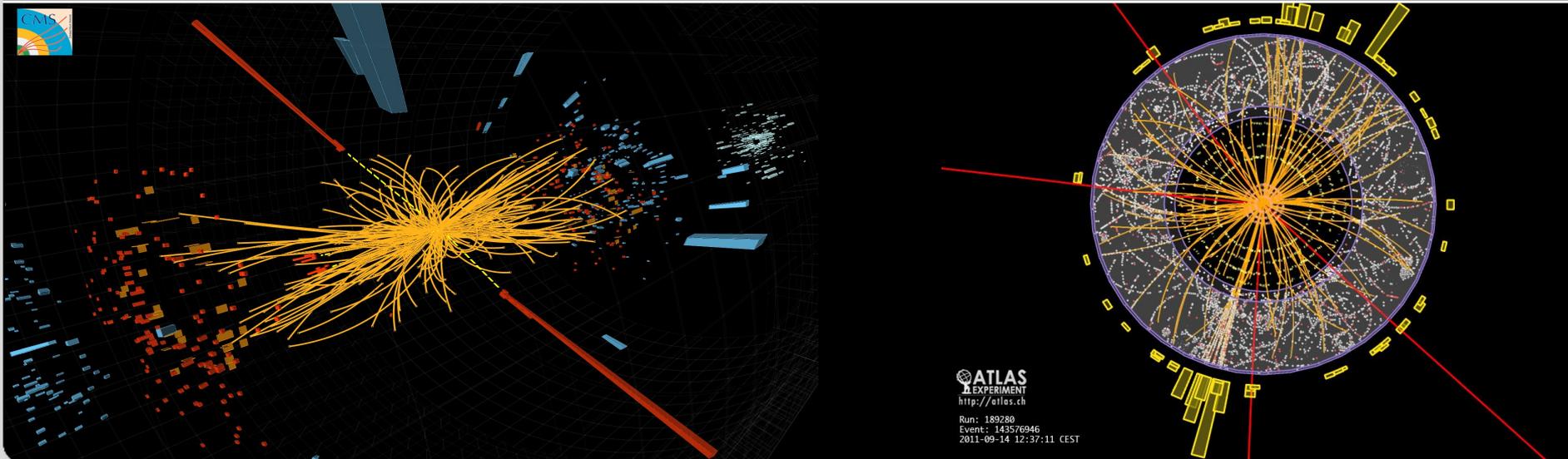
Higgs Boson Physics

Analysis Techniques

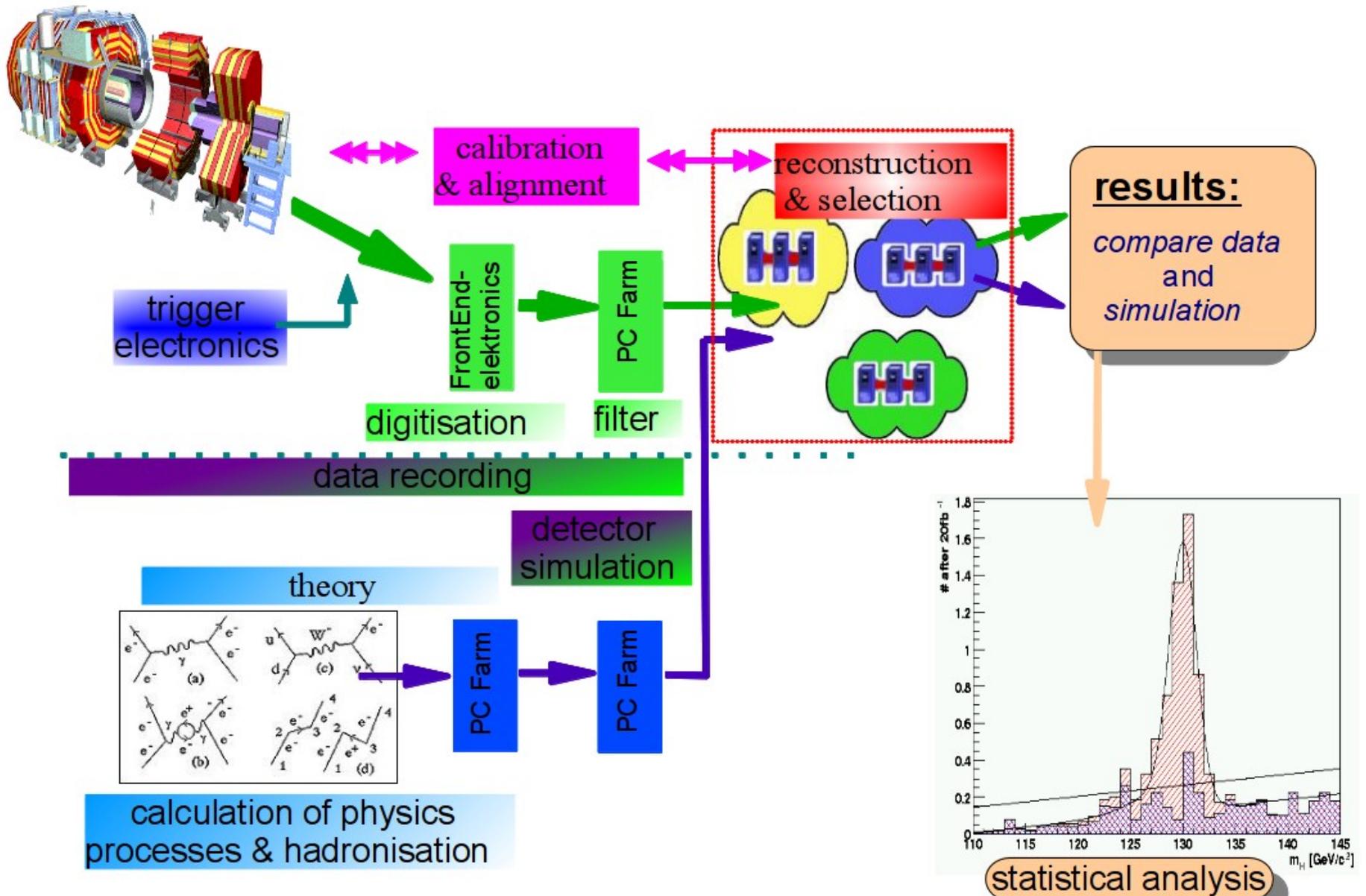
Günter Quast, Alexei Raspereza

Master-Kurs
SS 2012

Institut für Experimentelle Kernphysik



Overview: Components of Analysis Chain



Components of Analysis Chain

- **Digitizers** record data from detector cells
 - remove empty cells („zero-suppression“ and „noise reduction“)
- **Trigger and Filter** select „interesting“ events „on-line“ to be stored for „off-line“ analysis
 - (events not stored at this point are lost forever !)*
- **Reconstruction** process constructs physical objects (electrons, muons, jets, ...)
 - (this and subsequent steps can be repeated many times)*
- **Pre-selection** identifies interesting events and objects in events for further processing and analysis
- **Analysis** compares measured distributions with theoretical expectations

Theory

Experiment

- **theoretical calculation** of production cross sections
- **hadronisation** of quarks and gluons into jets

• **Detector simulation**

same reconstruction, selection and analysis steps
for **simulated events** as for **real events**

Monte-Carlo Generators

want to understand

\mathcal{L}_{int} → final states

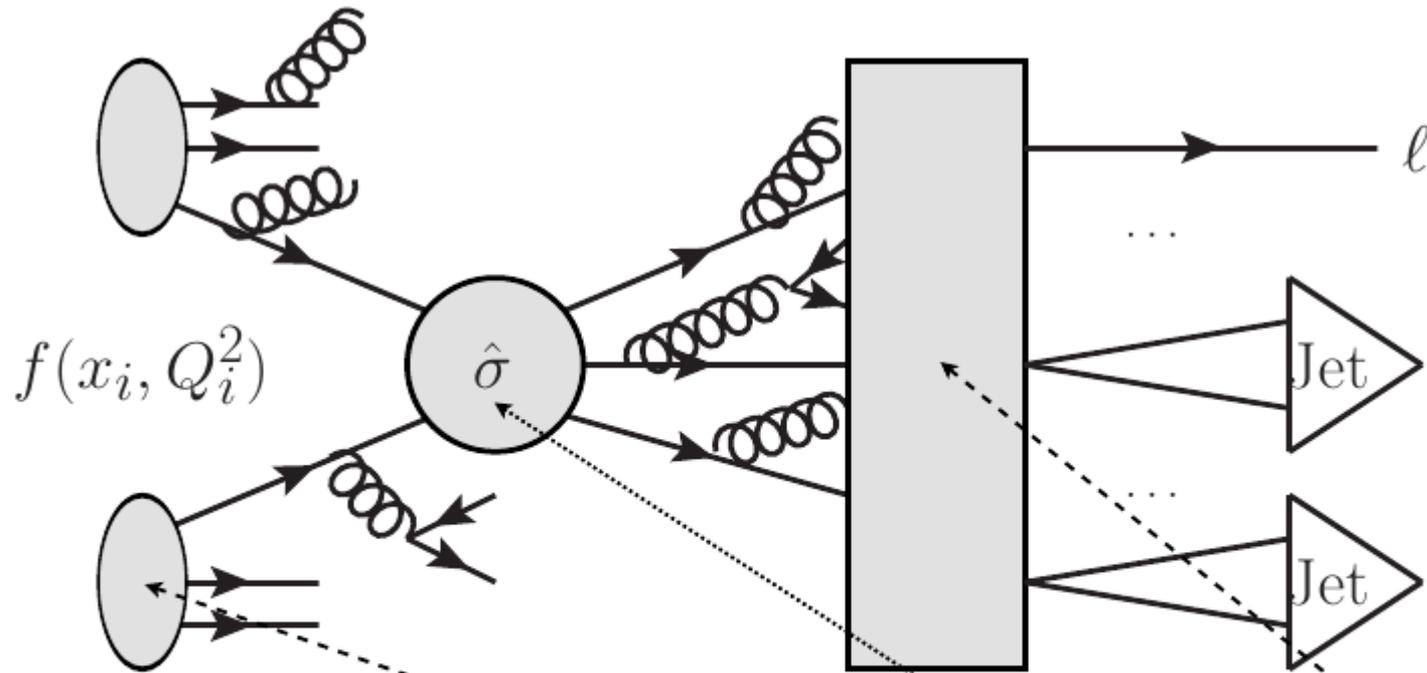
and predict measurable quantities

$\frac{\partial \sigma}{\partial O_i}$ = differential cross section

O_i : production angles of final state particles,
momenta of final state particles,
invariant mass of (groups of) final state particles;
...

Calculation of Cross sections

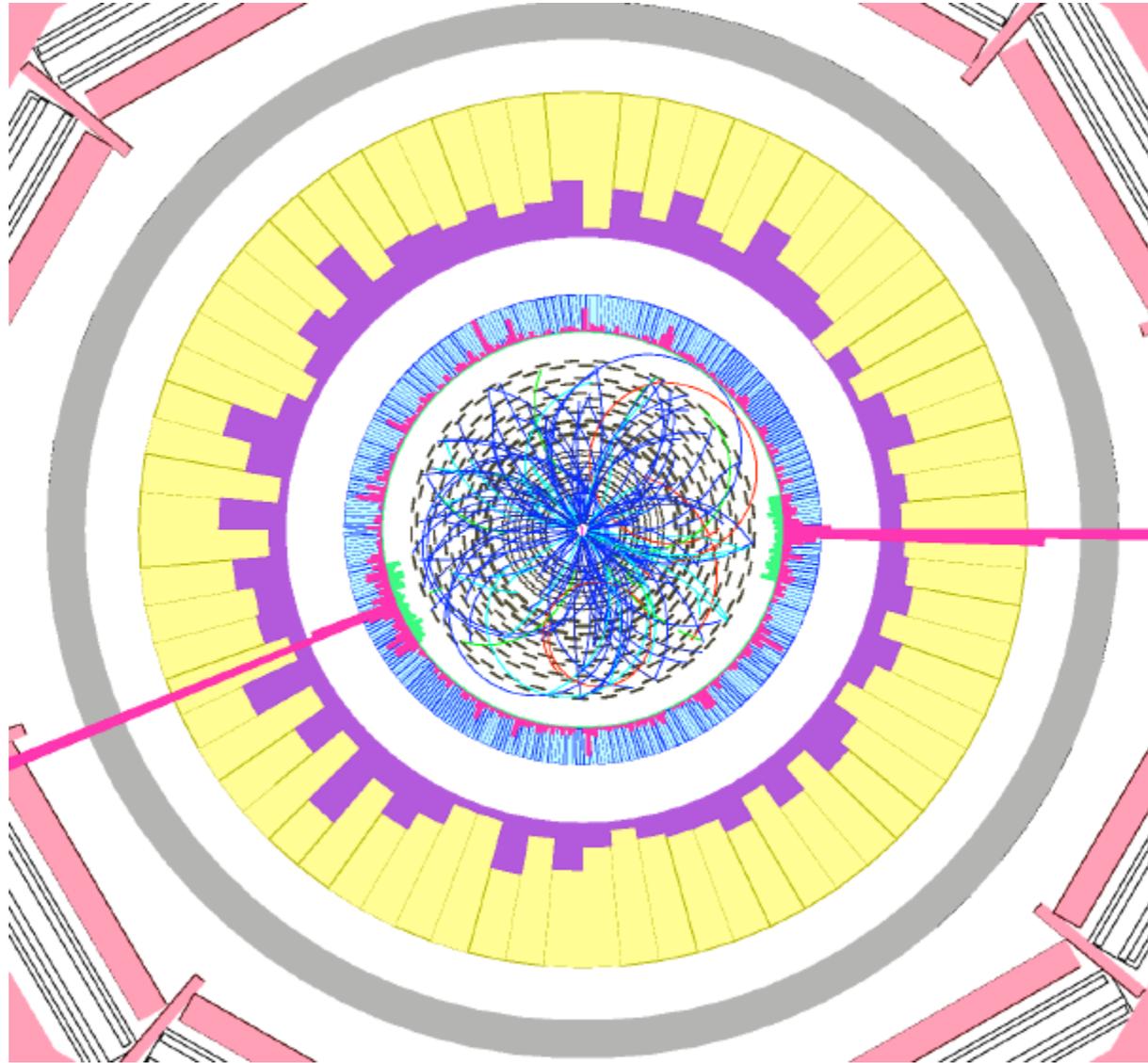
$$\sigma = \text{PDFs} \otimes 2 \rightarrow n \text{ process} \otimes \text{hadronization}$$



$$\sigma_{\text{QCD}} = \sum_{jk} \int dx_j dx_k f_j(x_j, \mu_F^2) f_k(x_k, \mu_F^2) \cdot \hat{\sigma}(x_j x_k S, \mu_F^2, \mu_R^2) \otimes \text{hadronization}$$

Complicated process – use MC techniques to calculate cross sections, phenomenological modes to describe hadronization process (quarks \rightarrow jets)

Example: simulated Higgs Decay in CMS



Can you see the Higgs?

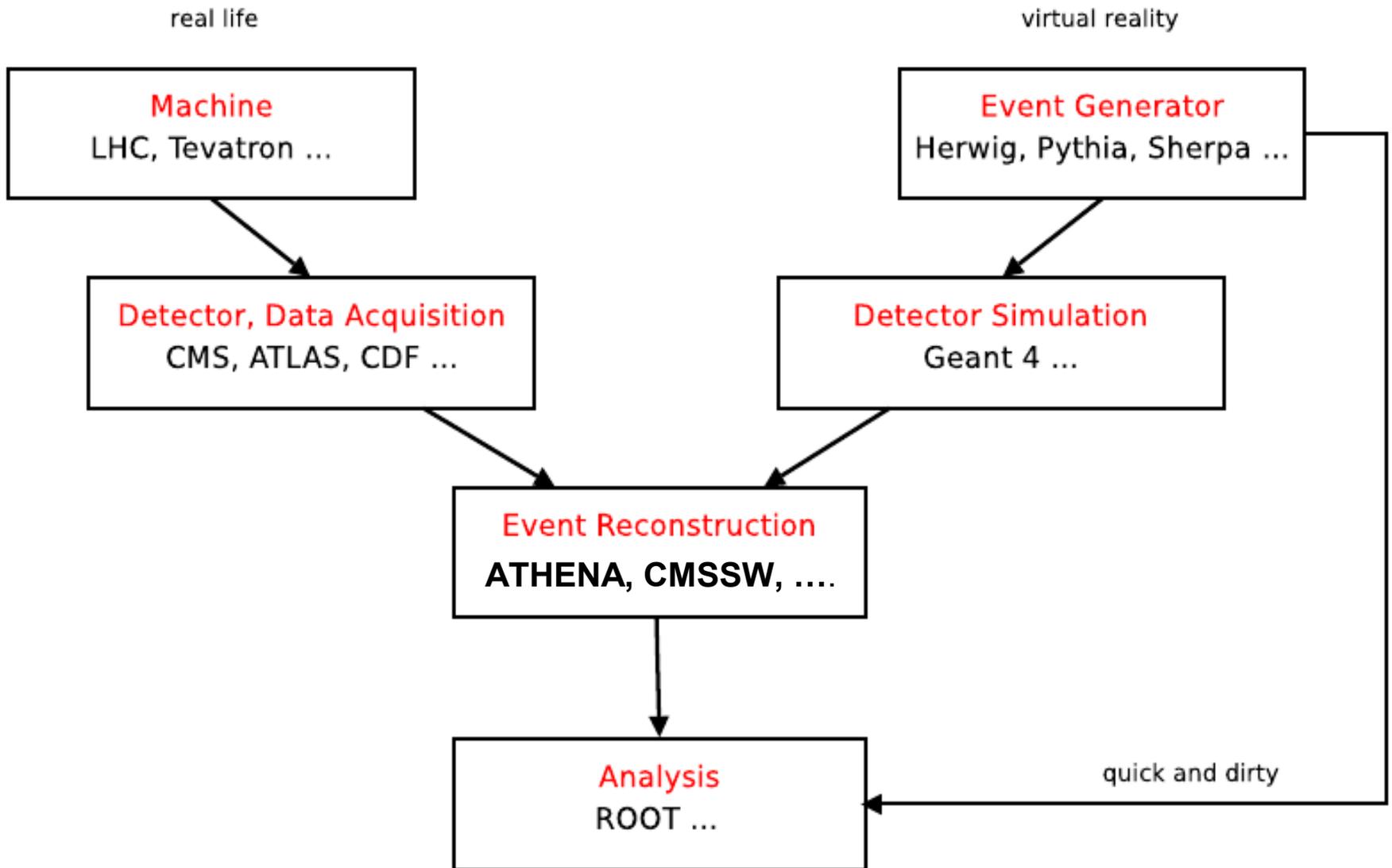
nice lecture, much more detailed than what can be shown here:

Monte Carlo School 2012, Helmholtz Alliance „PhyScis at the Terascale“
lecture by Stefan Giesecke, KIT

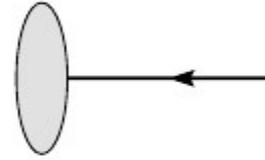
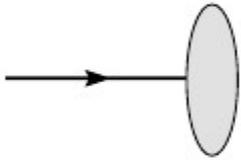
Technique in particle physics:

- Generate artificial events
reflecting all processes in the Lagrangian
using the Monte Carlo Technique
- obtain arbitrary distributions from simulated final state particles
- and compare with measurements

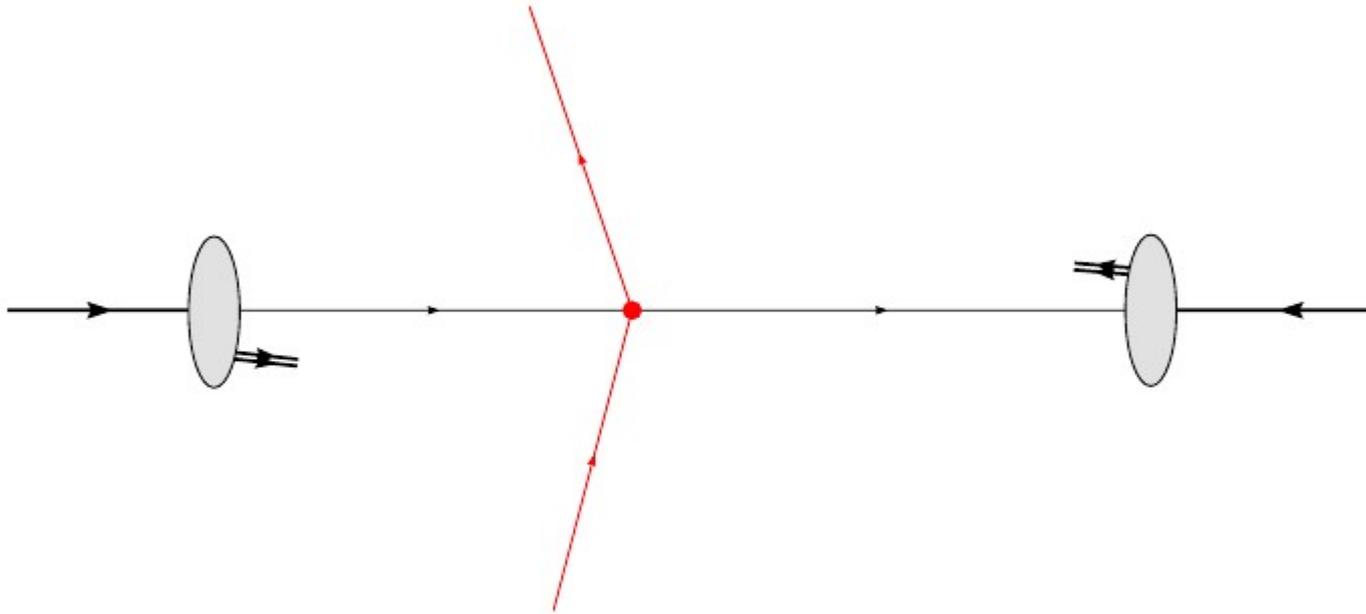
Steps of MC simulation



Example: pp collision



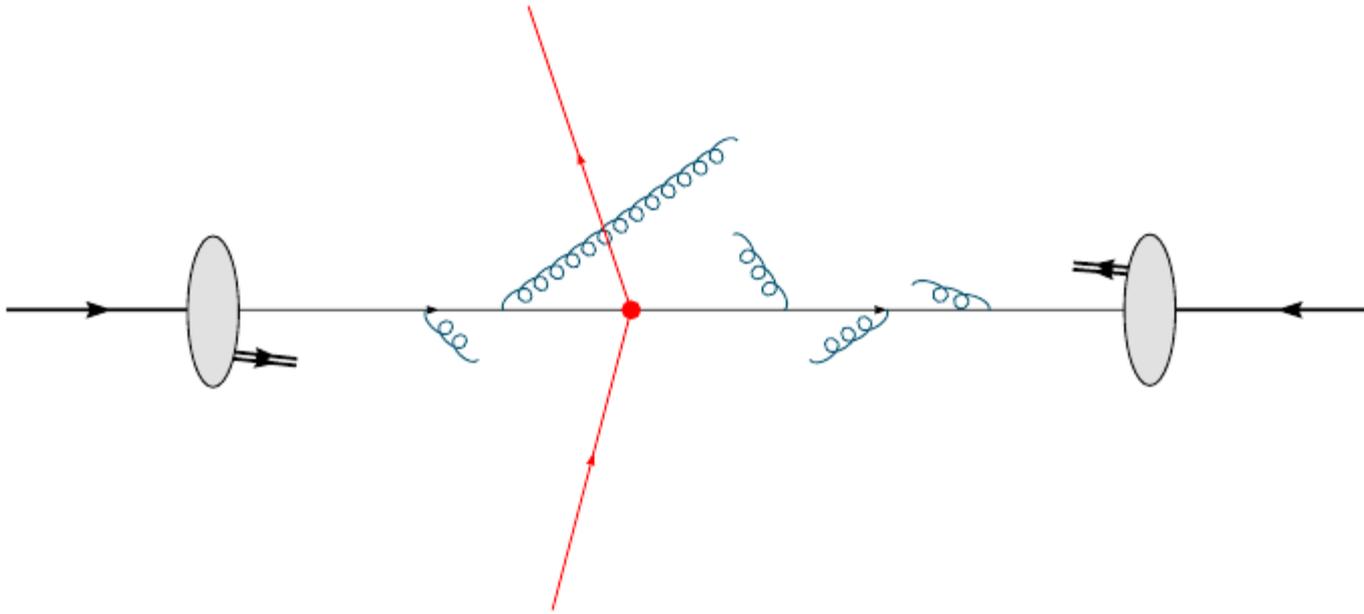
Example: pp collision



Stefan Gieseke · DESY MC school 2012

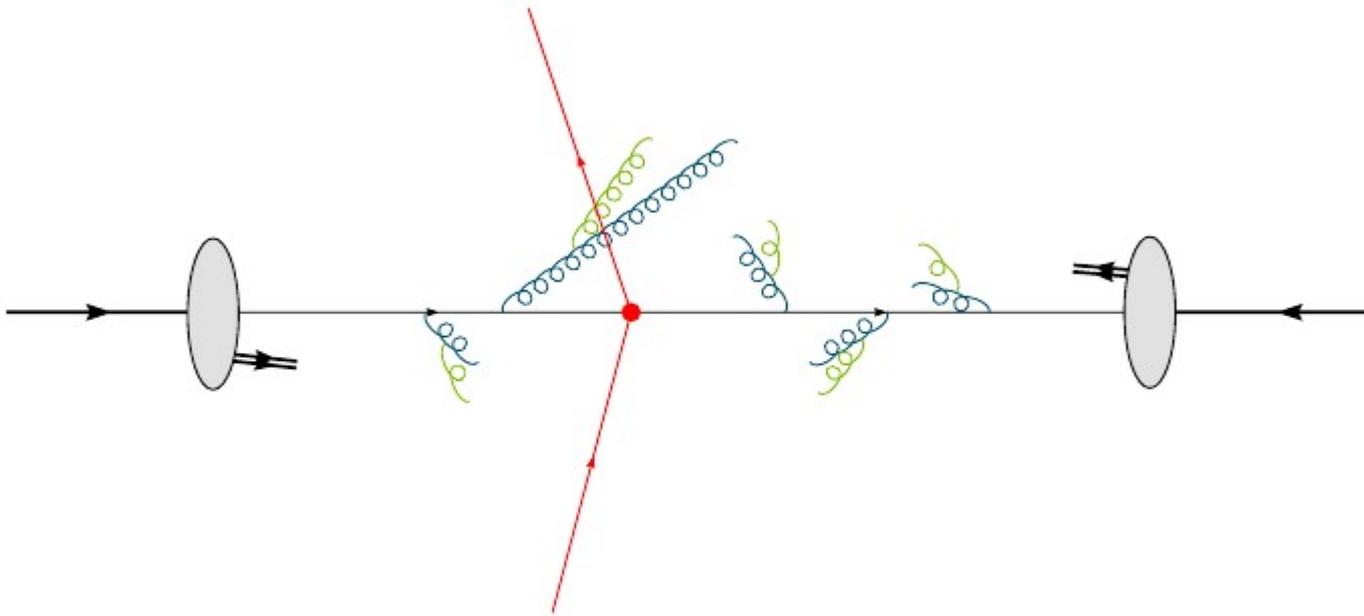
matrix element of hard process

Example: pp collision



parton shower

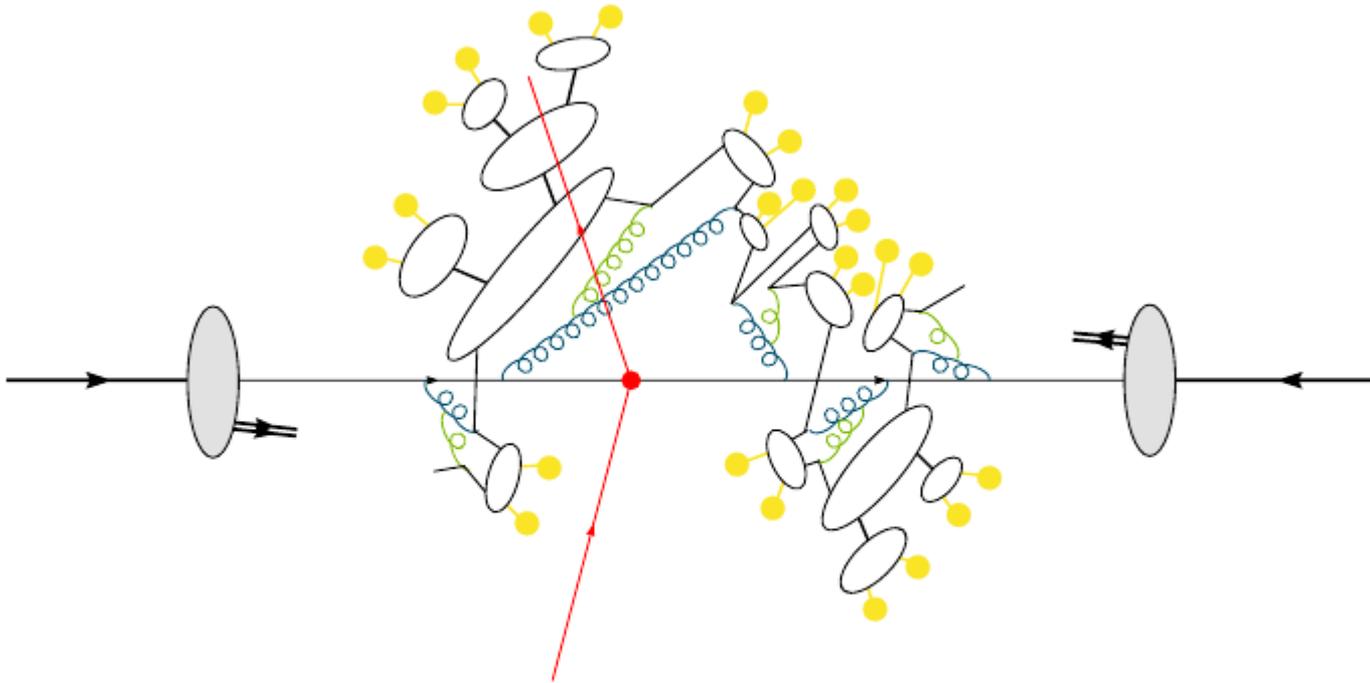
Example: pp collision



Stefan Gieseke · DESY MC school 2012

parton shower

Example: pp collision

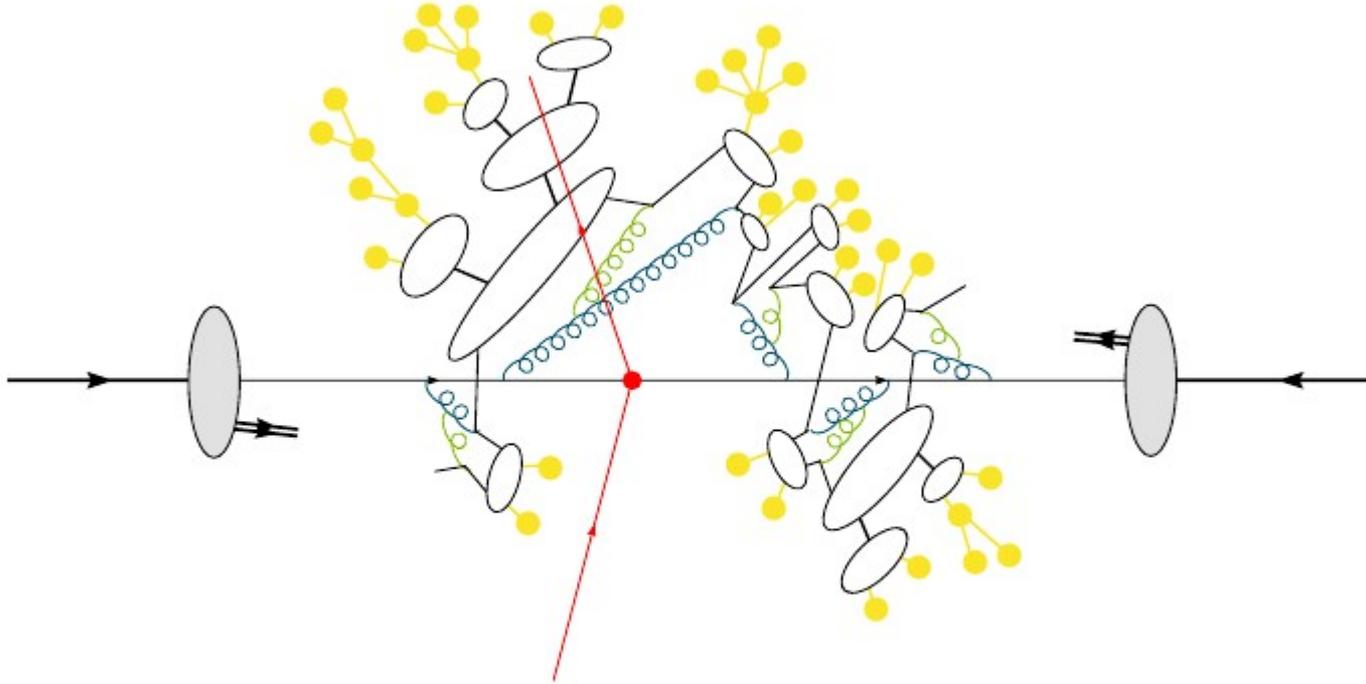


Stefan Gieseke · DESY MC school 2012

hadronization

phenomenological:
Lund string model
(Pythia)
or
cluster hadronisation
(Herwid(++))

Example: pp collision

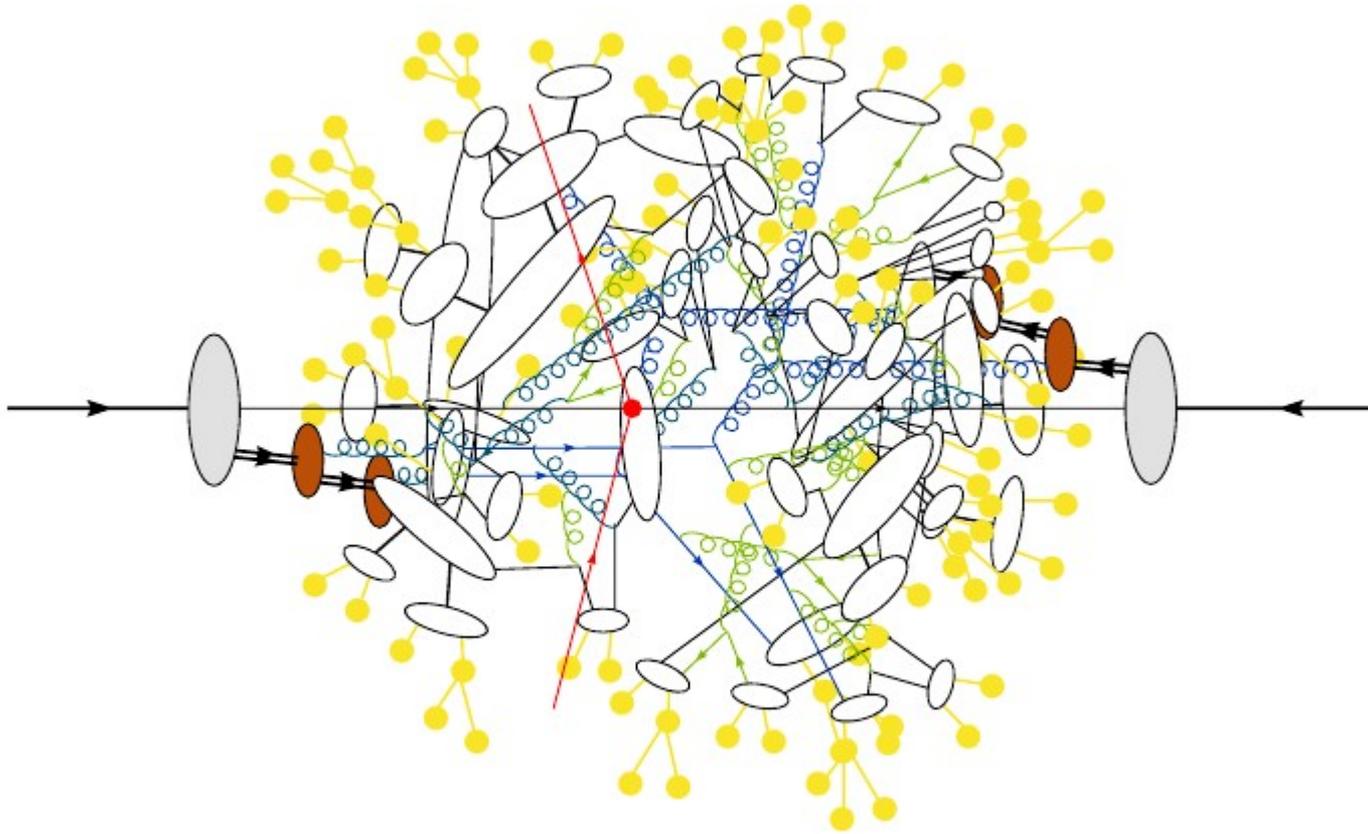


Stefan Gieseke · DESY MC school 2012

hadron decays

tedious -
relies on
measurements

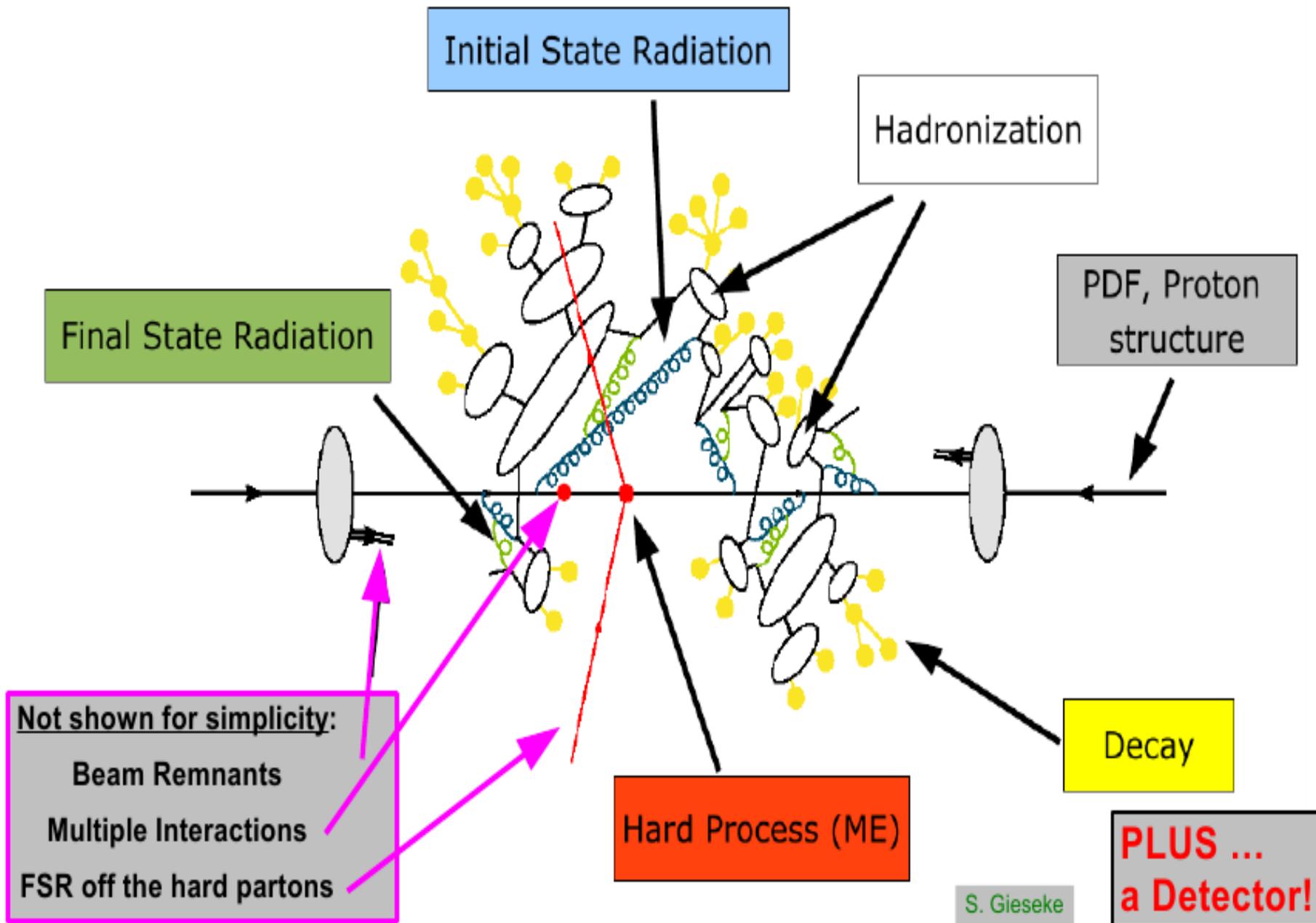
Example: pp collision



relies on models
& measurements
→ needs „tuning“

Multi-parton interactions and
underlying event

Summary: pp collision



Example: pp collision

last step:

- process stable particles through detector simulation to obtain „hits“ in detector cells;
- run reconstruction software to obtain „reconstructed objects“
- run selection procedures („Analysis“)
to obtain „identified reconstructed objects“

in total:

true properties of objects from hard process at parton level
are **folded** with

- parton distribution functions,
- hadronization effects,
- detector acceptance and efficiency,
- reconstruction efficiency and resolution,
- identification efficiency and purity

to obtain **reconstructed properties**

all steps involve multi-dimensional integrations;

Monte Carlo is the only choice !

Result of Simulation

Distribution(s) of

– signal events

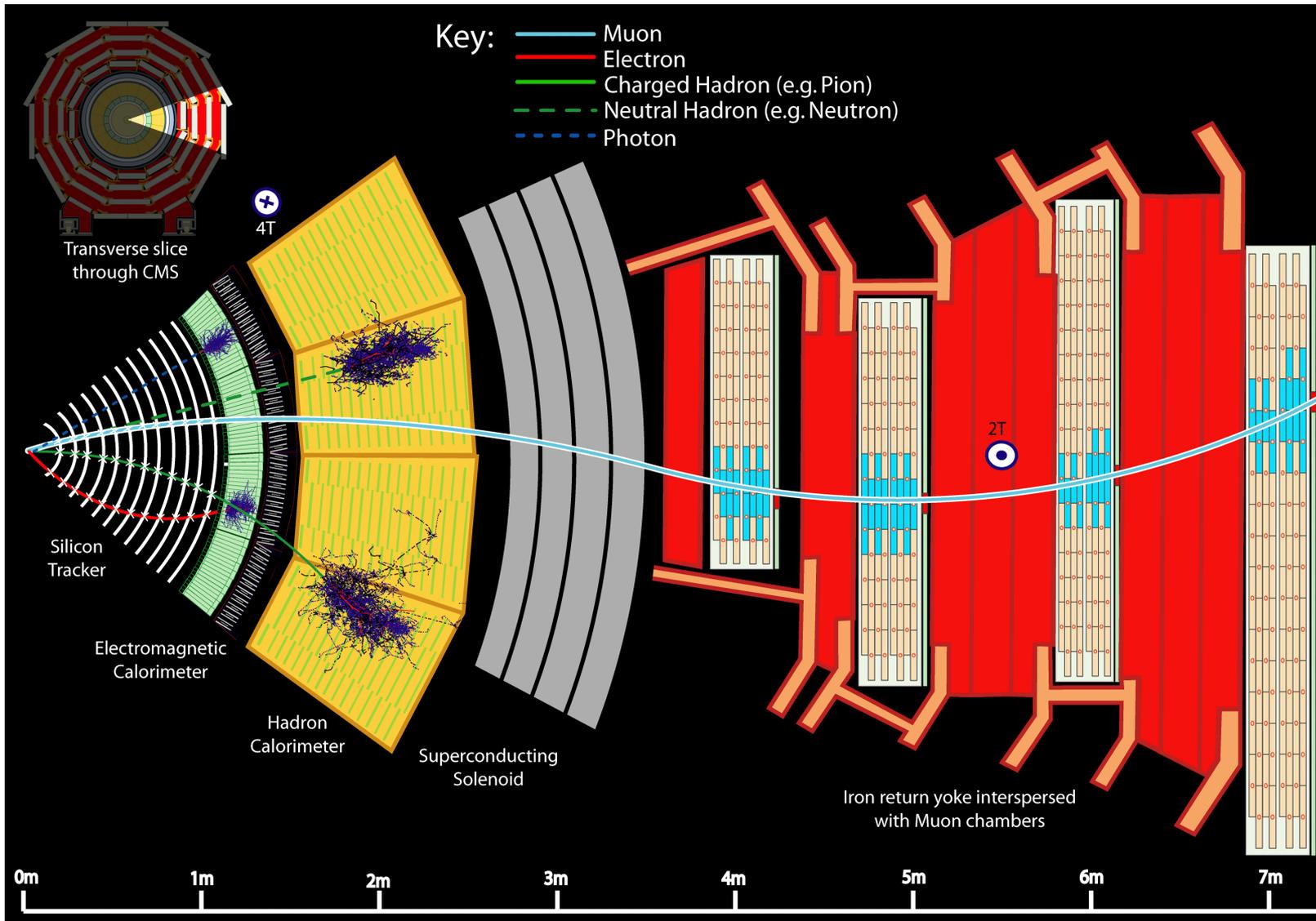
– background events

Used to formulate
„signal+background“ (S+B)
and
„background-only“ (B)
hypotheses for comparison
with data and
statistical inference

*Hint: in the real experiment, only very small numbers
are expected (see y-axis) – the question will be:
are they best described by the S+B or the B-only shape?*

The Real Experiment

Particle reconstruction



Detector registers only „stable particles“,
i.e. with life times long enough to traverse the detector

7 stable particles:
 γ , e , μ , p , n , π^\pm , K^\pm

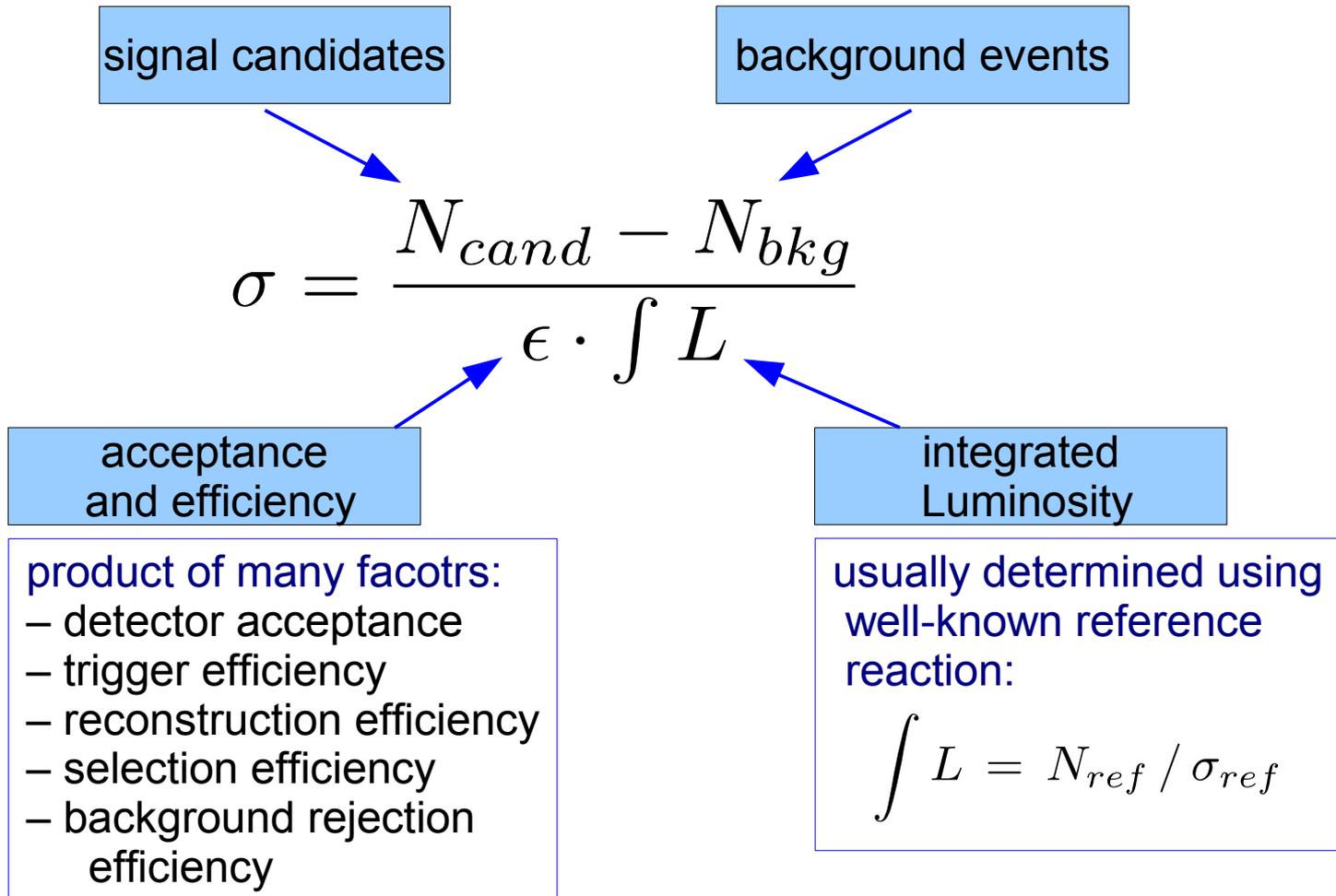
Steps of Event selection

- **hardware Trigger** and **on-line selection** identify „interesting“ events with particles in the sensitive area of the detector
(events not selected are lost)
→ detector acceptance and online-selection efficiency
- physics objects are **reconstructed** off-line
→ reconstruction efficiency
- **Analysis** procedure identifies physics processes and rejects backgrounds
→ selection efficiency and purity
- **statistical inference** to determine confidence intervals of interesting parameters (production cross sections, particle properties, model parameters, ...)

All steps are affected by systematic errors !

Cross section measurement

Master formula:



Cross Section measurement: errors

by error propagation →

$$\frac{\delta\sigma}{\sigma} = \sqrt{\frac{\delta N_{cand}^2 + \delta N_{bkg}^2}{(N_{cand} - N_{bkg})^2} + \left(\frac{\delta\epsilon}{\epsilon}\right)^2 + \left(\frac{\delta \int L}{\int L}\right)^2}$$

This is the error you want to minimize

- with signal as large as possible
- background as small as possible
- nonetheless, want large efficiency
- luminosity error small (typically beyond your control, also has a “theoretical” component)

Trigger

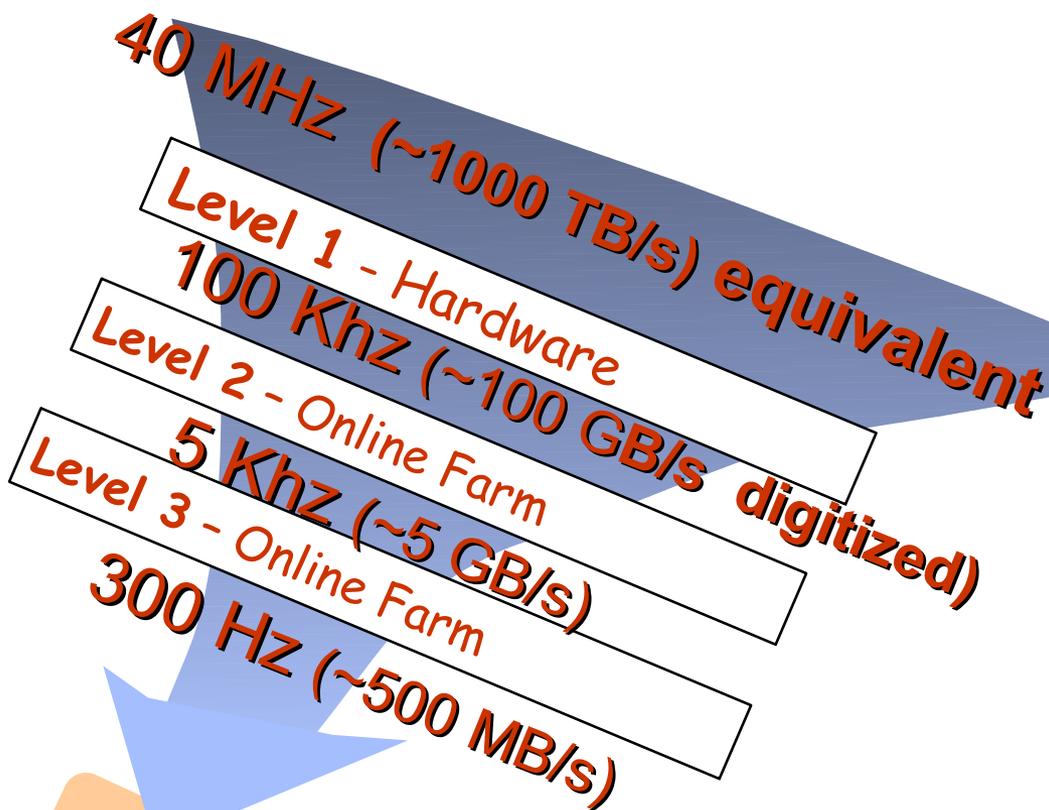
Online Data Reduction

- ~ 100 millionen detector cells
- LHC collision rate: 40 MHz
- 10-12 bit/cell

→ **~1000 Tbyte/s raw data**

Zero-Suppression & Trigger
reduce this to
„only“ some 100 Mbyte/s

i.e. 1  /sec



Large majority of events is not stored!

CMS Trigger & Data Acquisition

every 25 ns



40 MHz
COLLISION RATE

100 kHz
LEVEL-1 TRIGGER

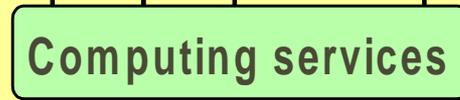
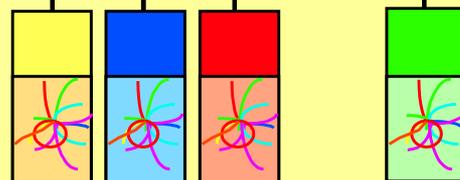
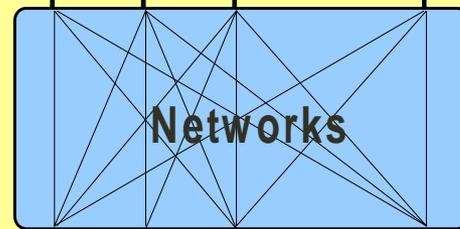
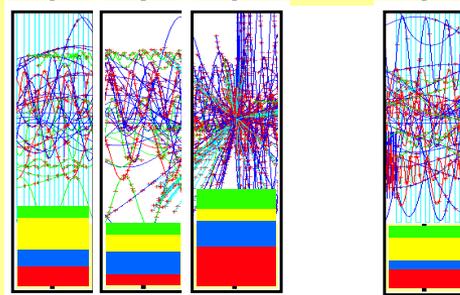
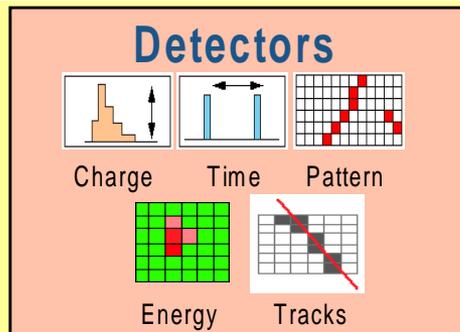
DAQ accepts
Level-1 Rate of 100kHz

1 Terabit/s
(50000 DATA CHANNELS)

500 Gigabit/s

HLT (High Level Trigger)
designed for O(100Hz)
- suppression factor ~1000
~2000 CPUs

Gigabit/s SERVICE LAN



16 Million channels
3 Gigacell buffers

1 Megabyte EVENT DATA

200 Gigabyte BUFFERS
500 Readout memories

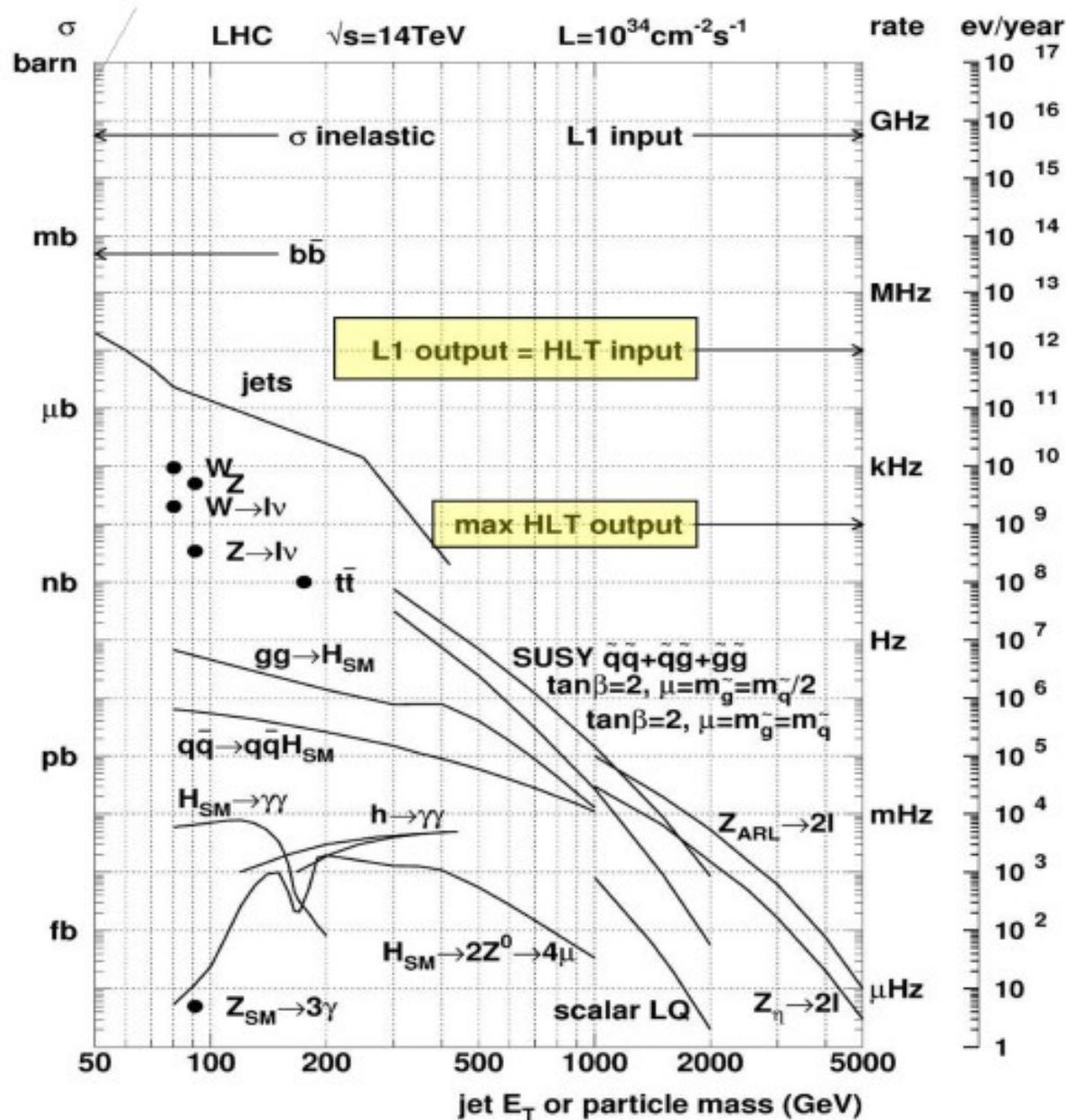
EVENT BUILDER. A large switching network (512+512 ports) with a total throughput of approximately 500 Gbit/s forms the interconnection between the sources (Readout Dual Port Memory) and the destinations (switch to Farm Interface). The Event Manager collects the status and request of event filters and distributes event building commands (read/clear) to RDPMs

5 TeraIPS

EVENT FILTER. It consists of a set of high performance commercial processors organized into many farms convenient for on-line and off-line applications. The farm architecture is such that a single CPU processes one event

Petabyte ARCHIVE

Trigger Rate vs. Cross section



Much of the
 “interesting physics”
 limited by maximum
 trigger rate !

What is easy to trigger ?

**Trigger thresholds rise as luminosity goes up,
and are a topic of permanent debate !**

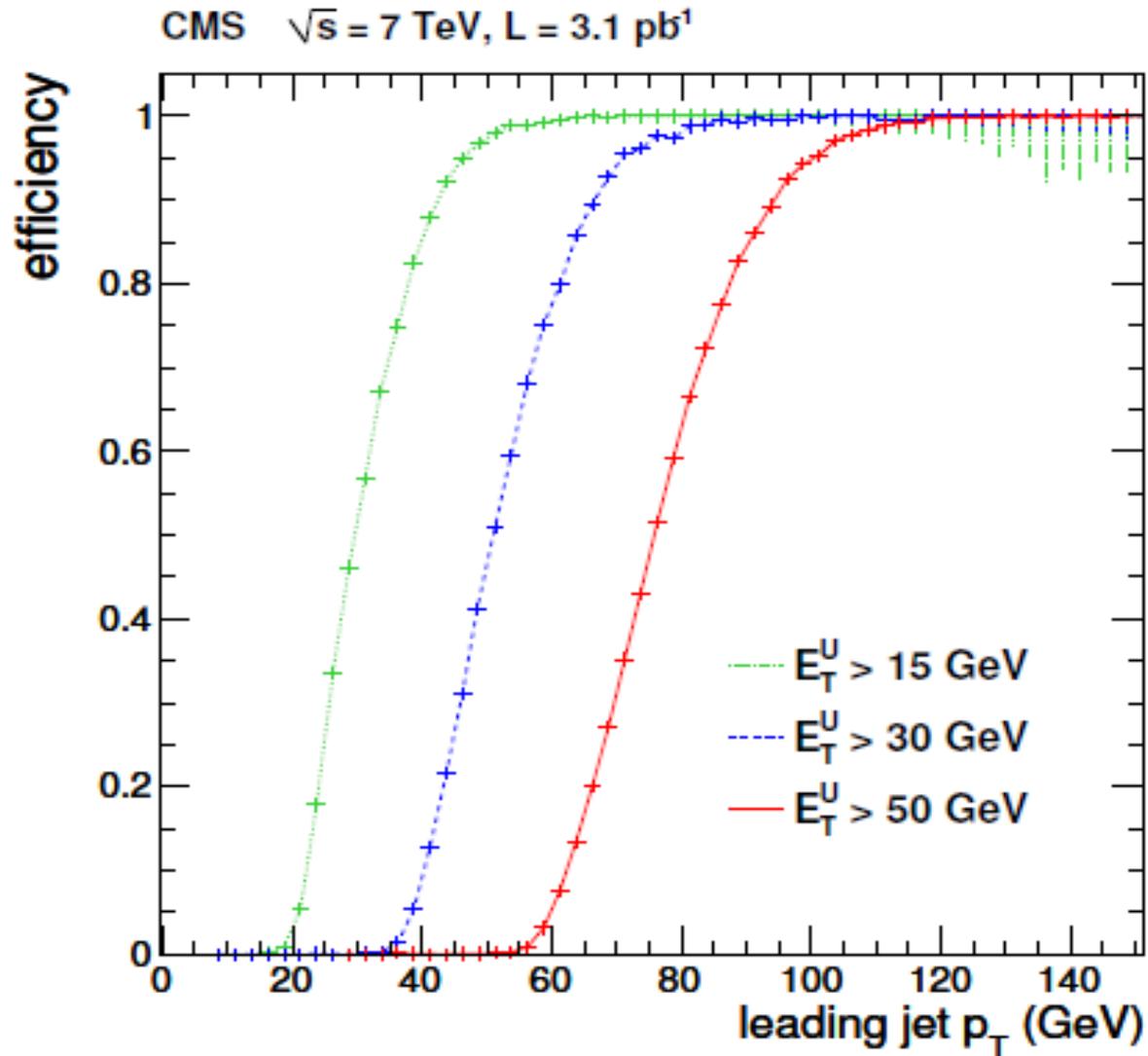
- isolated leptons with large transverse momentum $> \sim 20$ GeV
(from W, Z, top)
- di-lepton events with transverse momentum $> \sim 10$ GeV
- jets with very high transverse momentum (several 100 GeV)
- events with large missing energy (~ 100 GeV)
- isolated photons with transverse energy $> \sim 50$ GeV

lower-threshold triggers typically pre-scaled

Rest is difficult and probably not in recorded data !

for analysis, must know trigger efficiencies

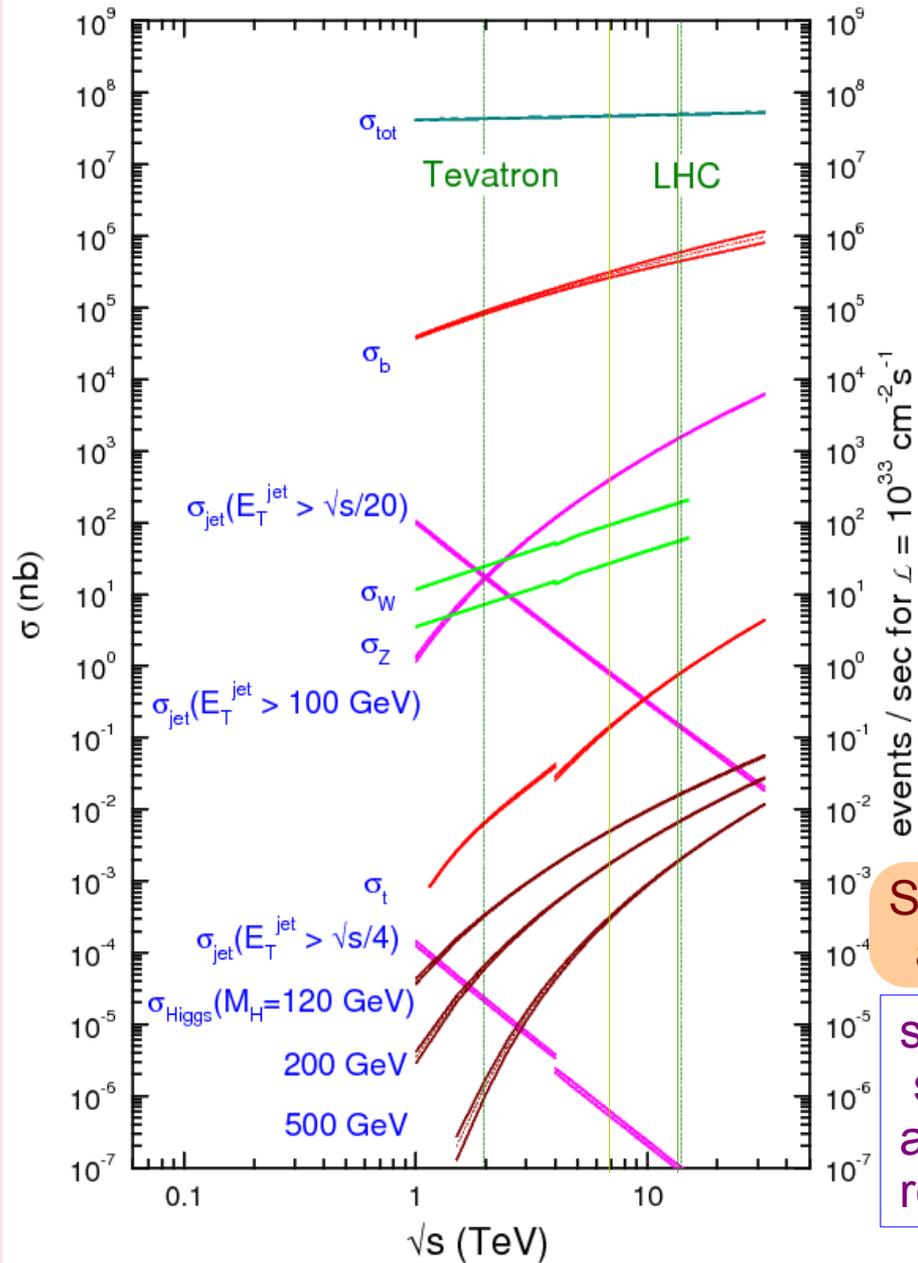
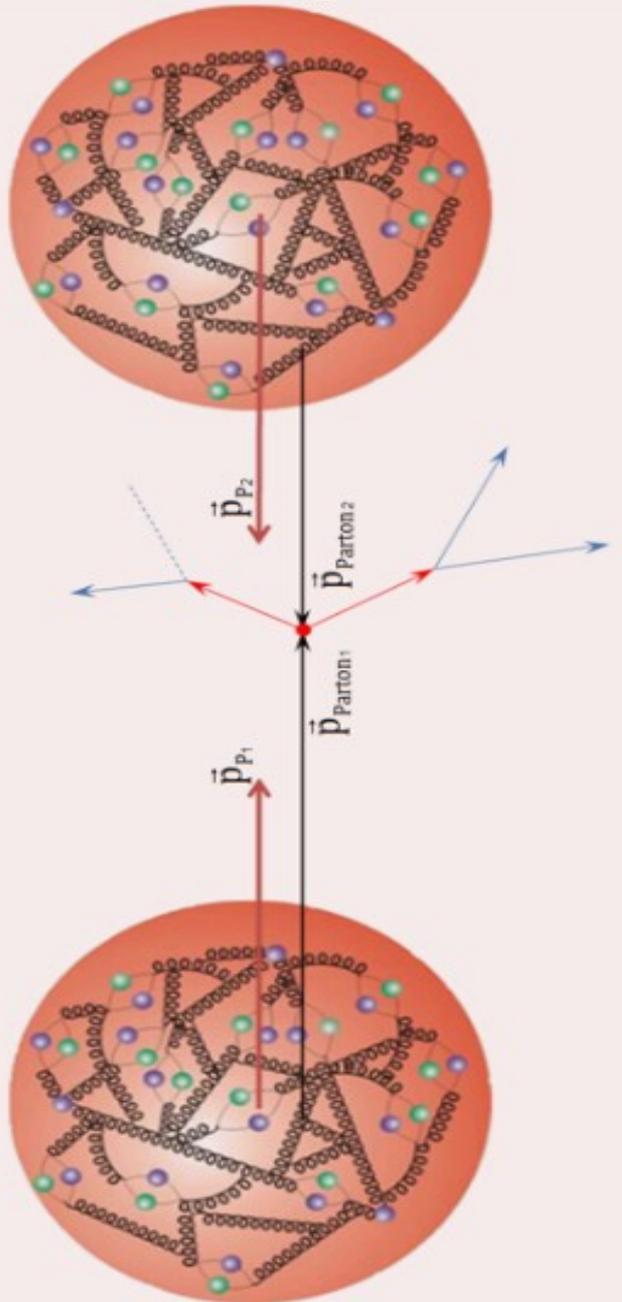
Example: trigger "turn-on" for jets



typical knee-shaped trigger efficiency curves (CMS, 2010), rising from 0 to 1

Data Analysis

Event Selection in the Analysis

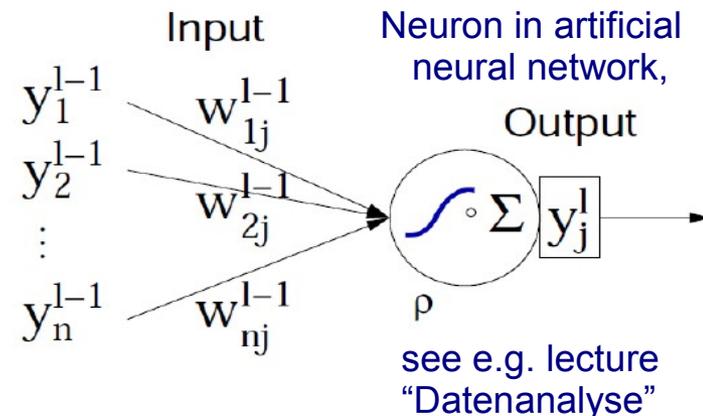


Some processes are very rare !

sophisticated signal selection and background rejection needed.

Analysis Steps

- recorded events are **reconstructed**: hits \rightarrow physical objects like electrons, muons, photons, hadrons, jets, missing energy ...
need to know reconstruction efficiency and resolution
- selection** of “interesting events” and objects for a particular analysis
affected by selection efficiencies for signal and background processes
- last step of analysis involves advanced algorithms for the optimal **separation of signal from background** and **extraction of parameters** of interest from the background-corrected signal distribution
(multivariate analysis, MVA, like discriminant methods, decorrelated likelihood, artificial neural networks, boosted decision trees)
understanding the systematics involved is required !



Finally, arrive at a result with statistical and systematic errors
evaluation of systematics requires much hard work

Much use of **simulated data** is made in this process
to evaluate known or suspected sources of uncertainties
and propagate them to the final results.

Reconstruction of Objects

1. **combine sub-detectors** to classify all stable objects, i.e. find electrons, muons, photons, hadrons.

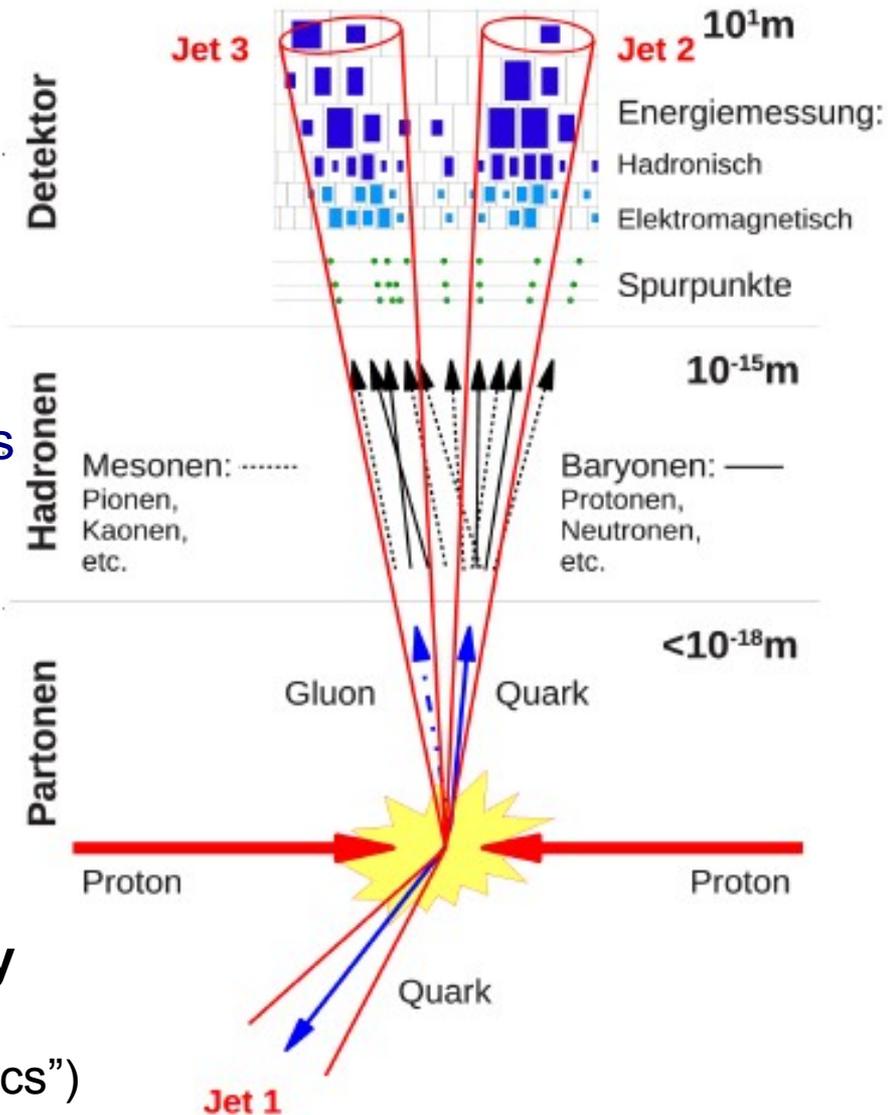
2. **cluster objects** into “jets”
relation between
measured final state objects
& hard partons

two types of algorithms:

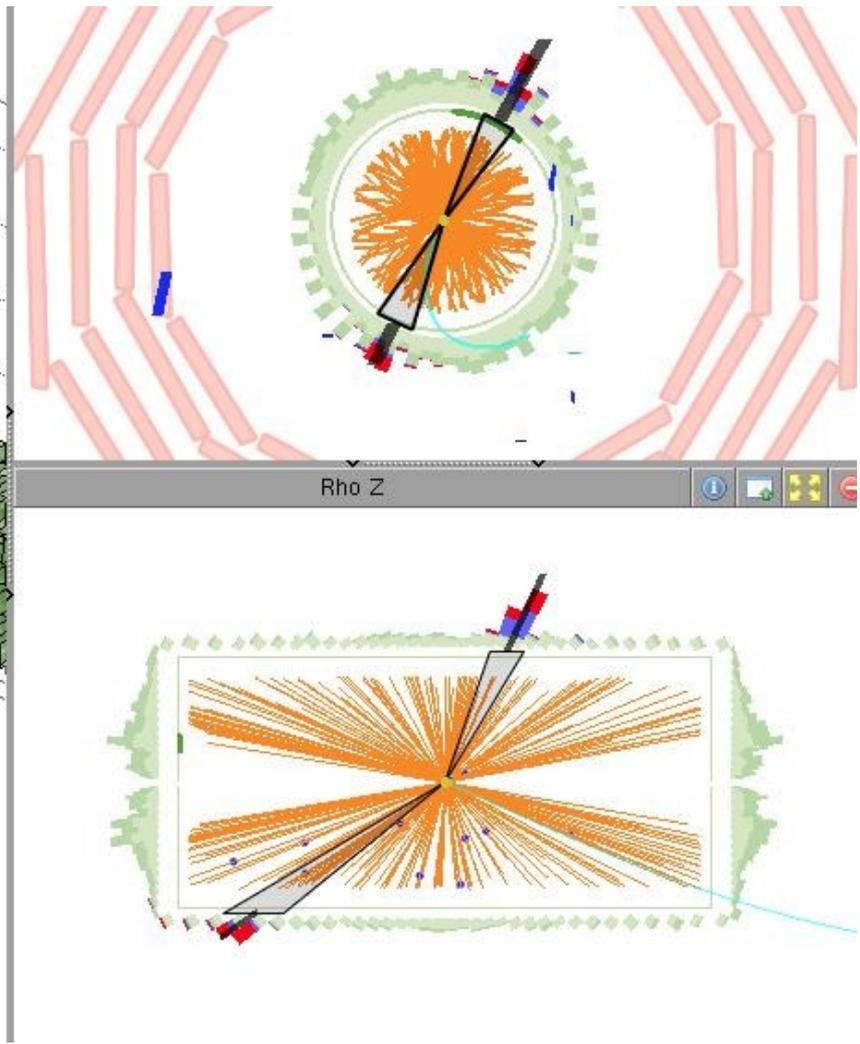
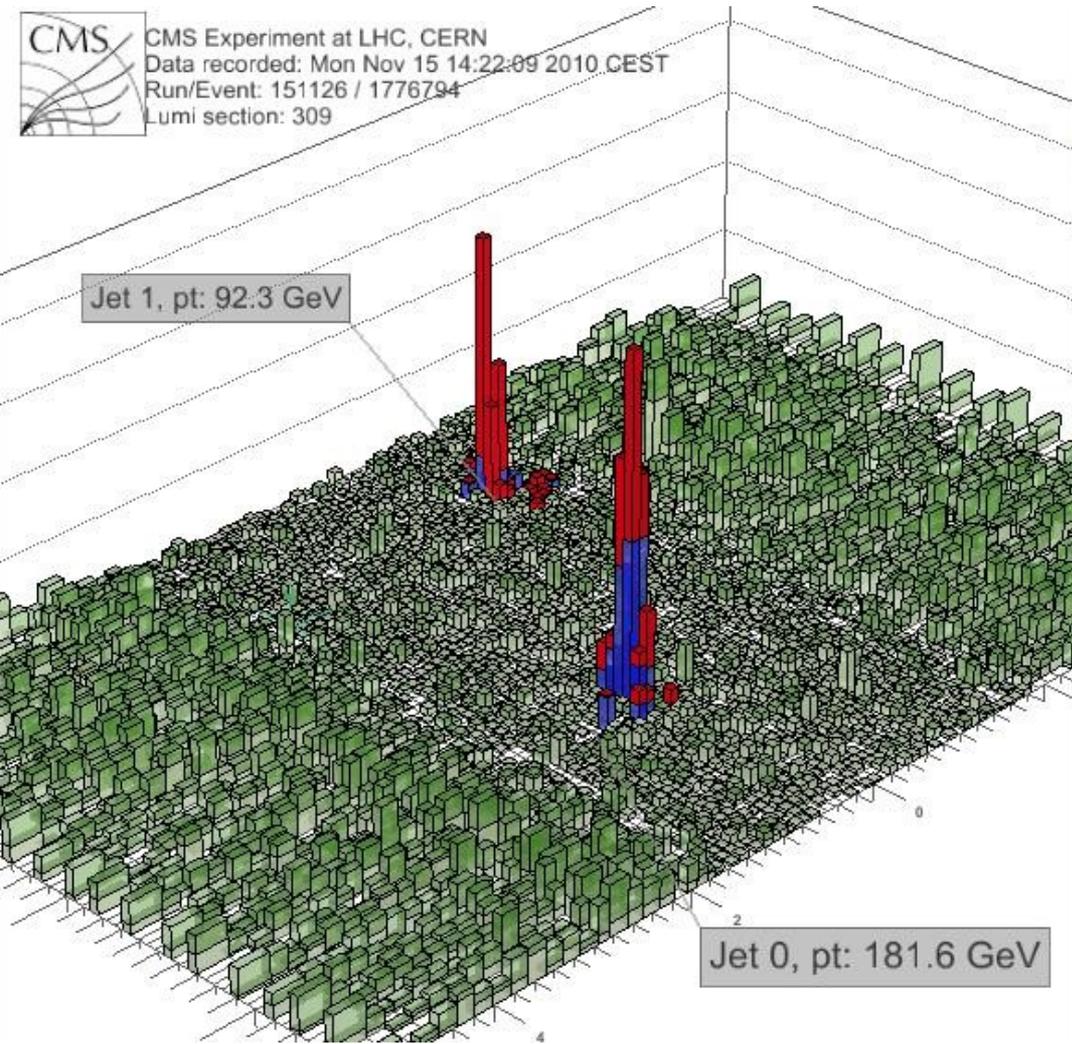
1. **“cone”**: geometrically assign objects to the leading object
2. **sequentially combine** closest pairs of objects – different measures of “distance” exist (kT, anti-kT) with some variation of resolution parameter, which determines “jet size”

CMS does this across detector components (“particle flow” analysis)

3. determine **missing transverse energy** carried away by undetectable particles (neutrinos, or particles signalling “new physics”)



Two-Jet Event in the CMS detector

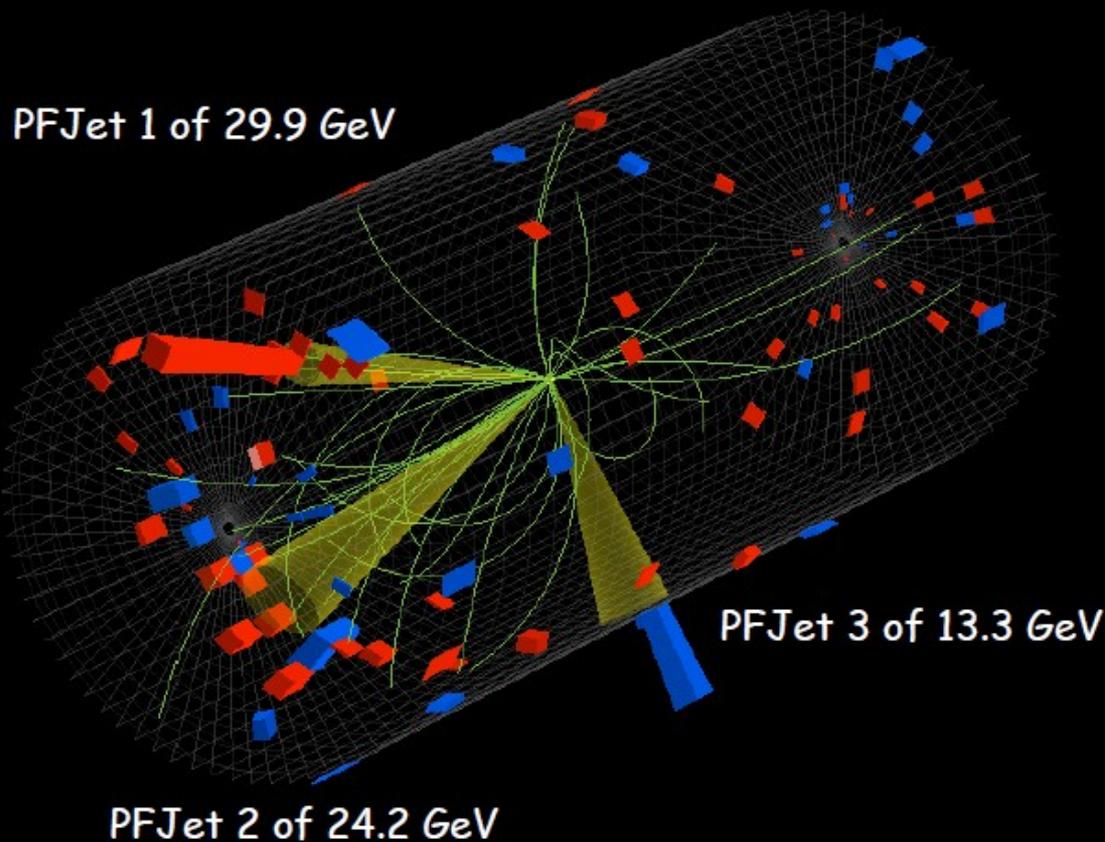


Three-jet event



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

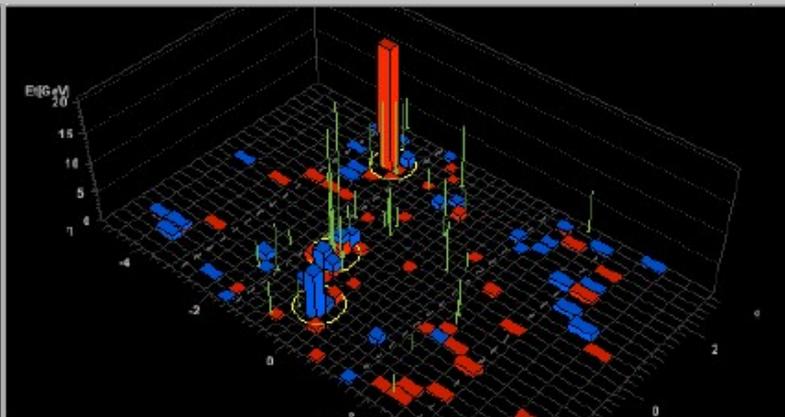
PFJet 1 of 29.9 GeV



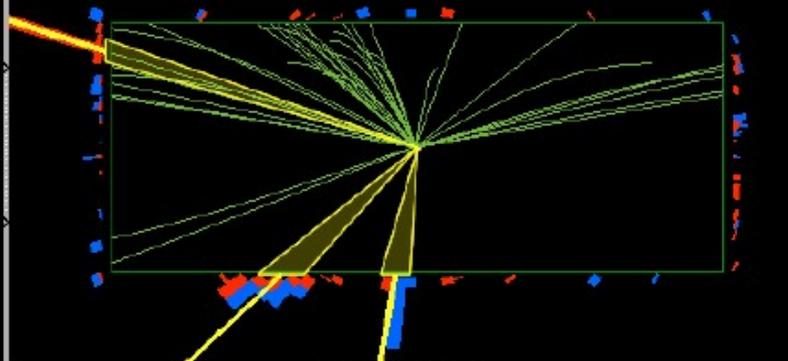
PFJet 3 of 13.3 GeV

PFJet 2 of 24.2 GeV

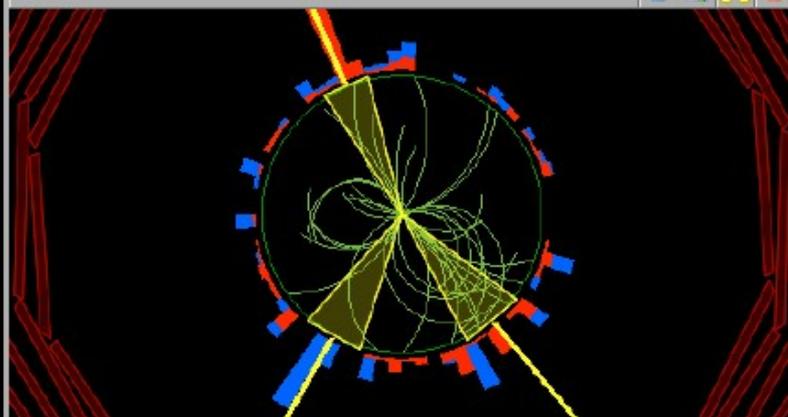
3 PFlow jets $p_T > 10$ GeV
 p_T cut on tracks displayed > 0.4 GeV



Rho Z



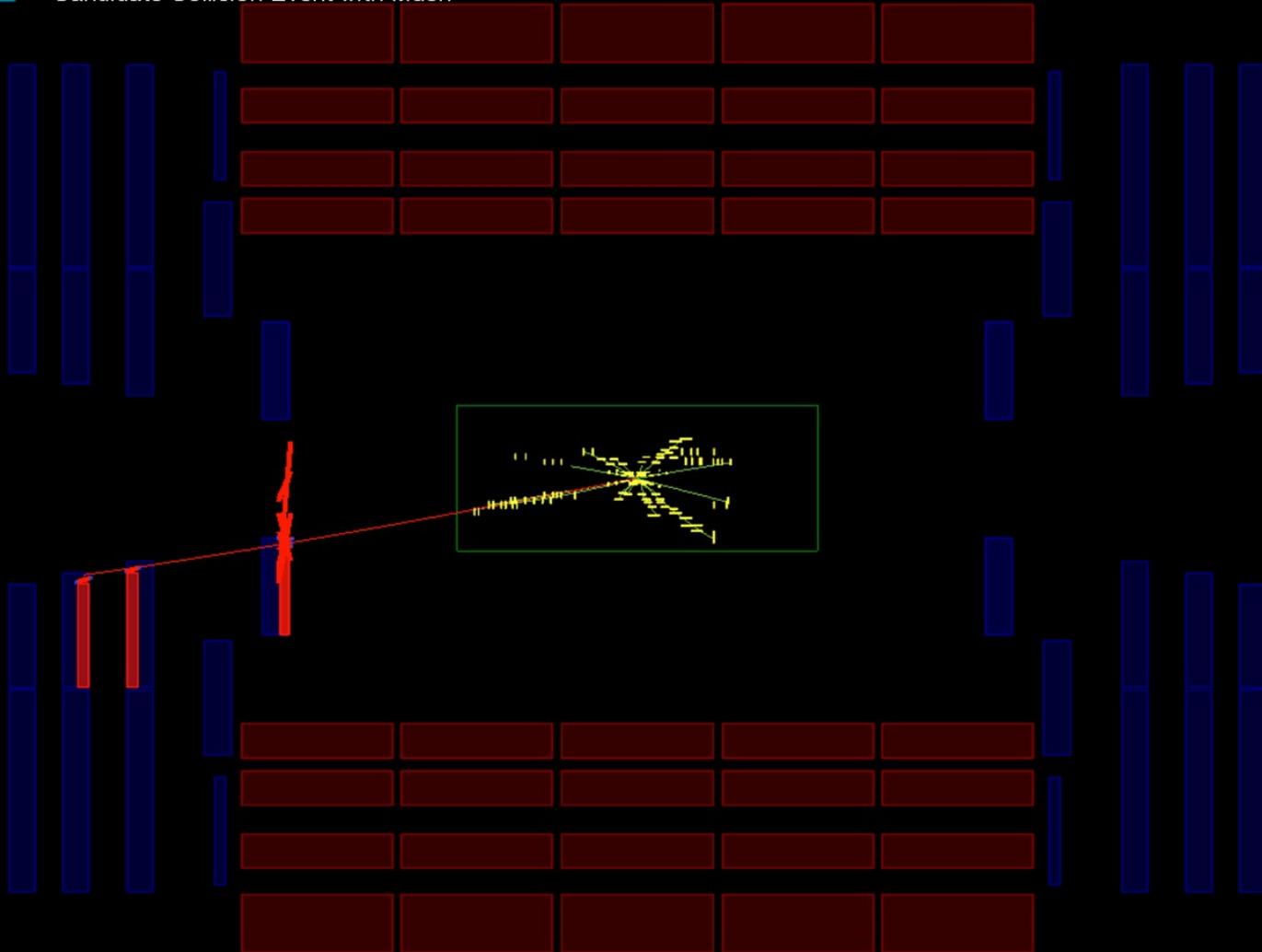
Rho Phi



event with 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



2 electrons in CMS

Summary View

- Add Collection
- ECal
 - HCal
 - Jets
 - Tracks
 - Muons
 - Electrons
 - Vertices
 - DT-segments
 - CSC-segments
 - Photons
 - MET
 - pfMet

Views

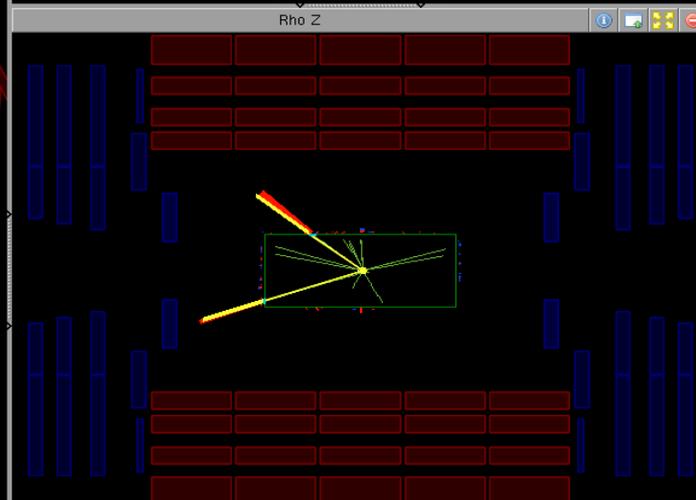
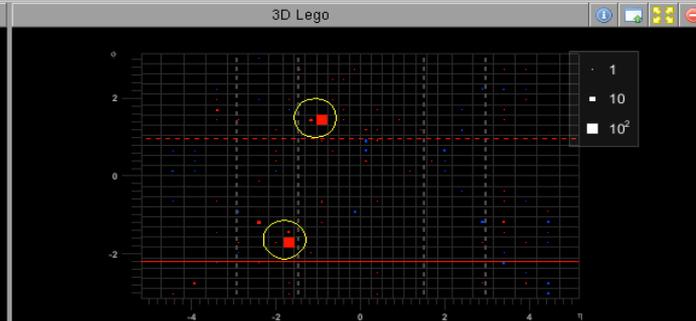
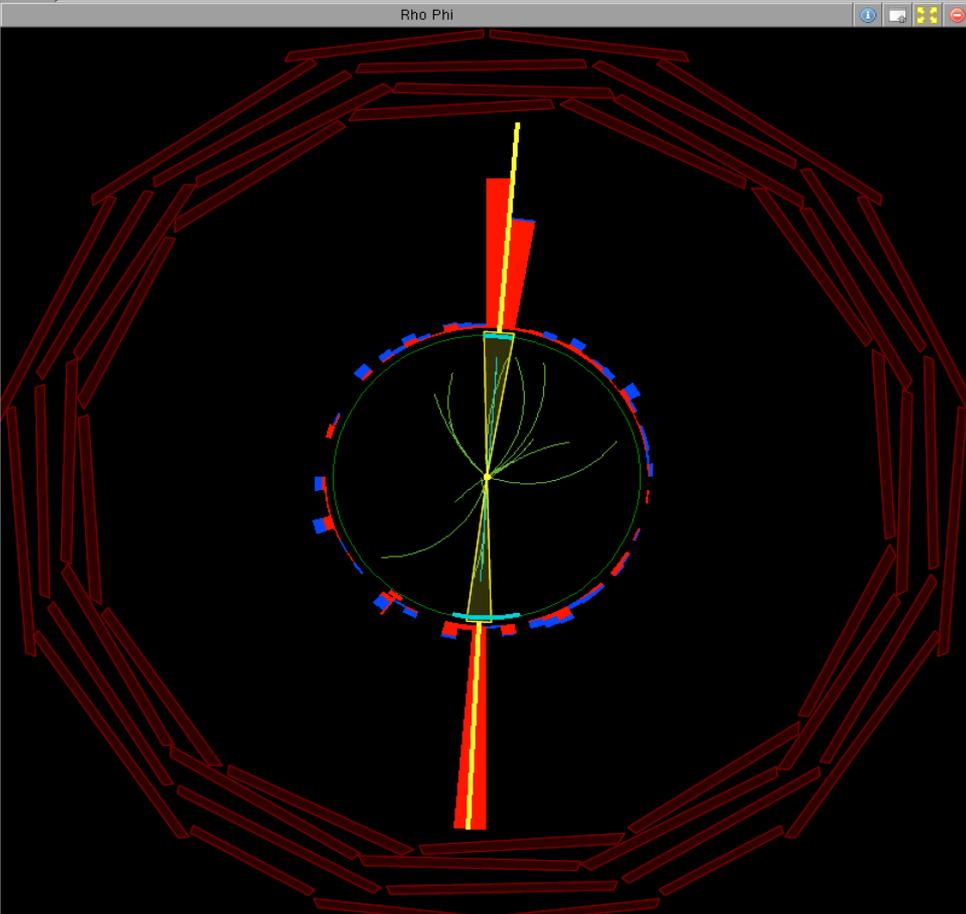


Table				
Collection	phi	sumEt	mEISig	
pfMet	4.1	-2.189	187.8	0.302

Table								
Collection	Pt	eta	phi	ECAL	HCAL	emf	size_eta	size_phi
Jets	43.1	-1.070	1.492	69.5	0.7	0.988	0.014	0.053
Jets	41.1	-1.802	-1.621	127.9	0.0	1.000	0.011	0.035
Jets	5.0	0.073	0.890	2.9	2.3	0.557	0.102	0.163
Jets	2.6	2.003	0.450	1.2	8.7	0.120	0.072	0.173
Jets	2.1	3.024	-1.223	5.7	15.9	0.264	0.160	0.067

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

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

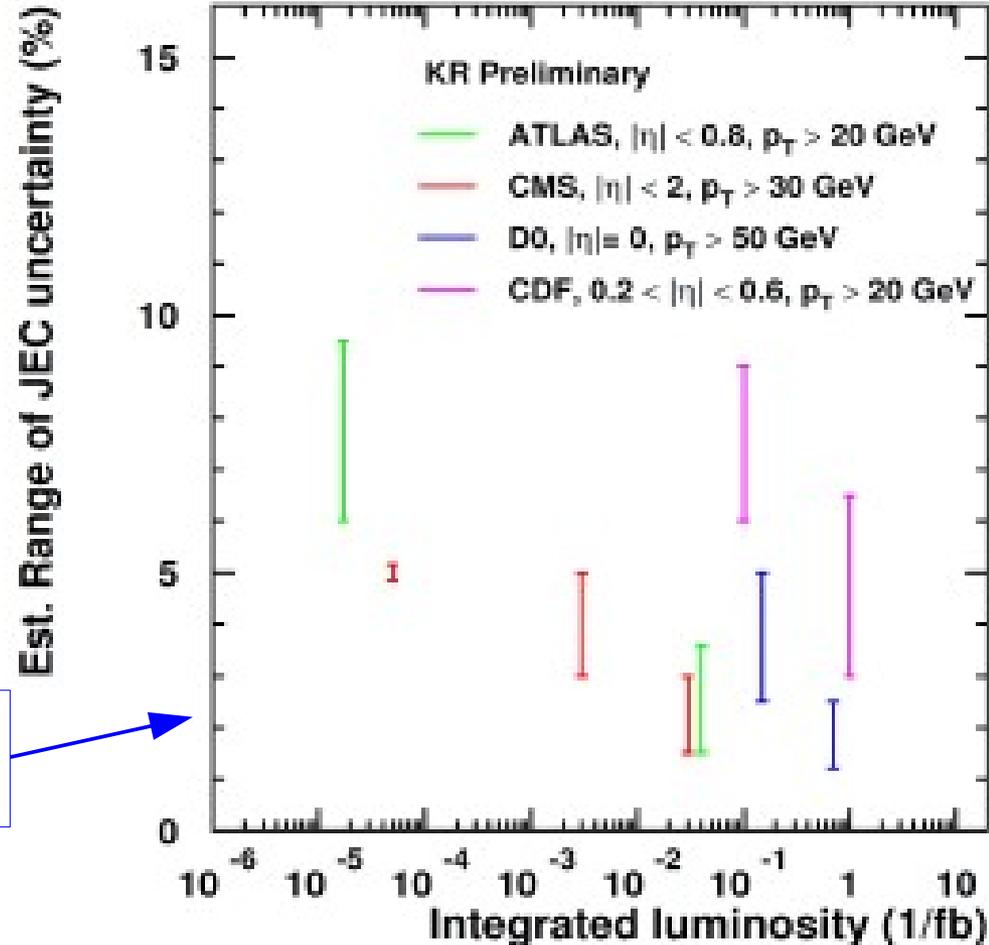
Calibration

Jets and missing transverse energy must be calibrated

relies on special topologies:

- **di-jet events** to equalize detector response
- **Z or γ balanced by a jet** to determine absolute scale
- events with **genuine missing energy** ($Z \rightarrow \nu\bar{\nu}$, W, Top)

Precision of Jet energy calibration reaches level of a few % !

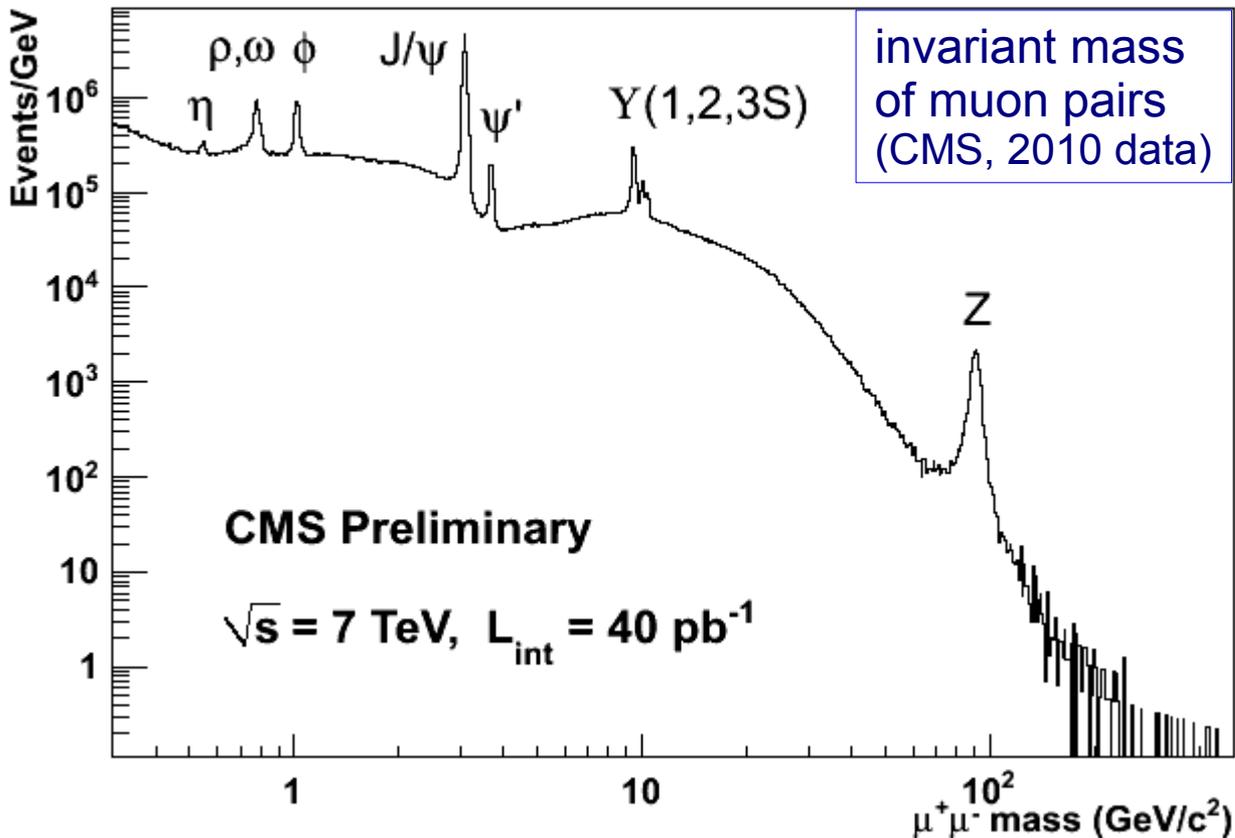


More complicated observables

Calculate **derived quantities** from objects:

- invariant mass of groups of objects
- missing energy or missing transverse momentum
- scalar sum of jet energies
- event shape variables (for QCD analyses)
- all kinds of “classifiers” for event classification

60 years of particle physics in only one year:



Determination of efficiencies

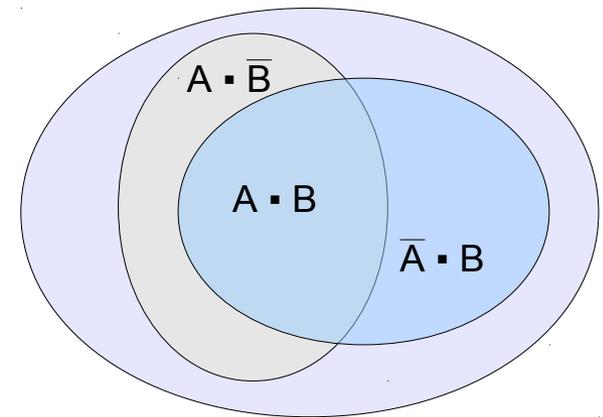
two options:

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

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

(statistical errors governed by Binomial distribution)

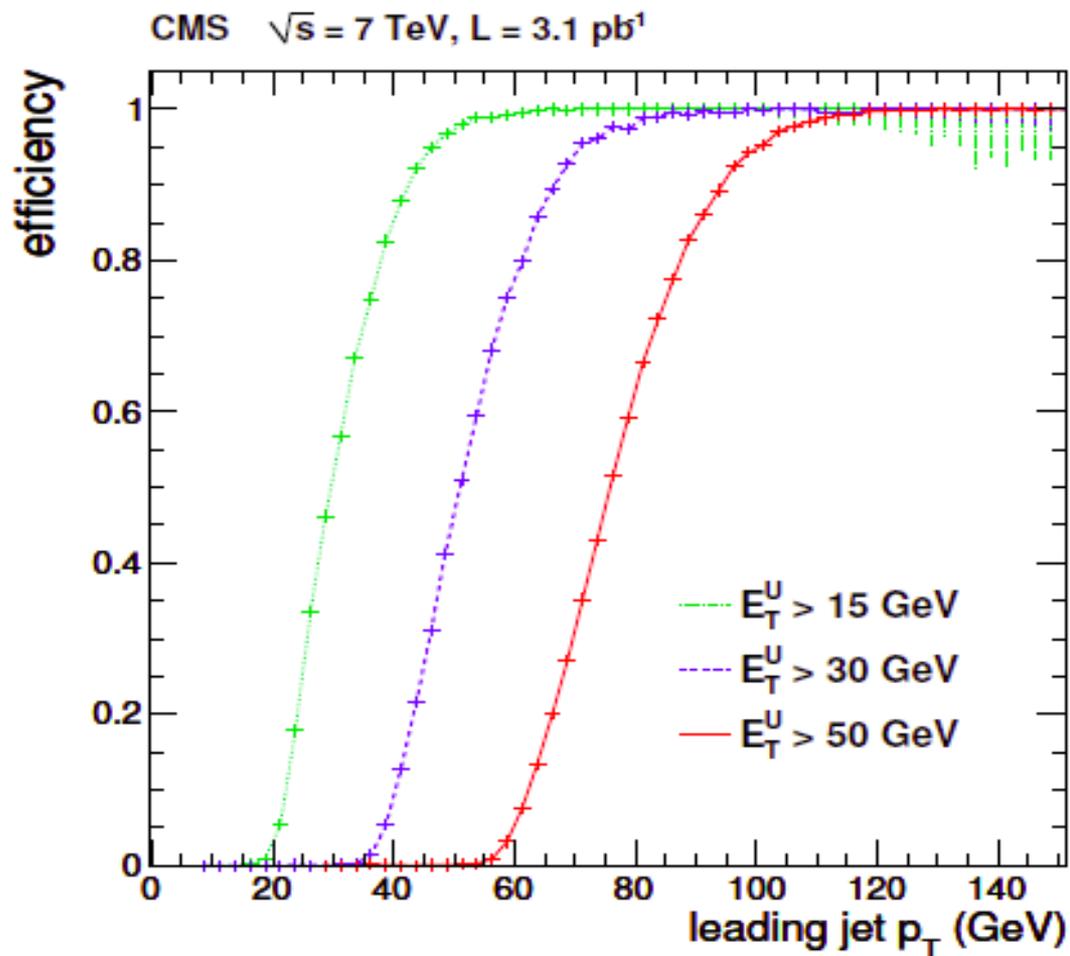


Example: 1 tight muon and one loose muon with tight selection on Z mass (“tag”) allows to measure the selection or trigger efficiency of second muon (“probe”)

Example: Trigger efficiencies

Typical “turn-on” curves of trigger efficiencies

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

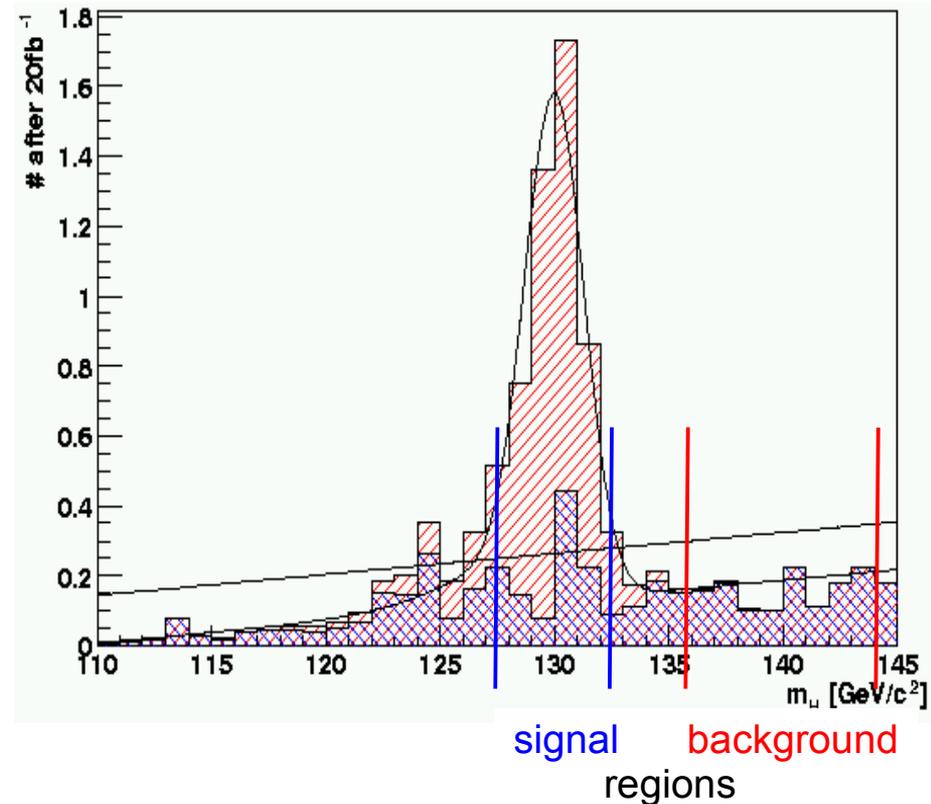


Remarks:

- 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

Determination of background

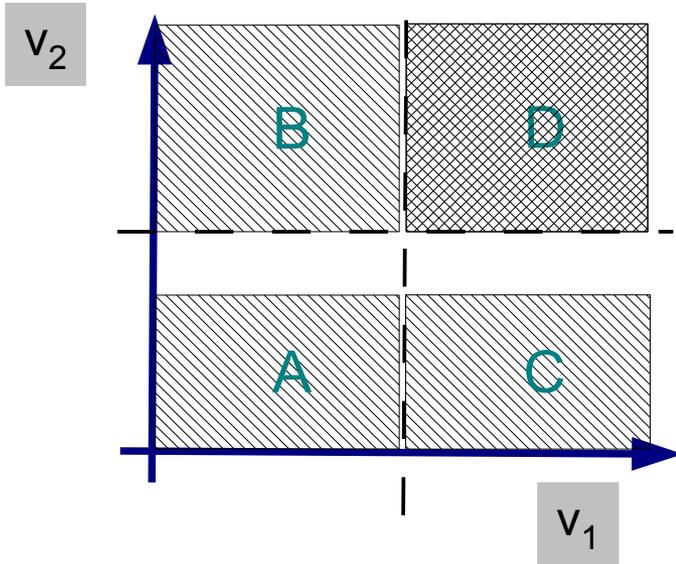
- take from **MC** (same comments as above)
- extrapolation from “side band” assuming “simple” signal shape or taking signal shape from MC



- if a **second, independent variable** can be found, background extrapolation from data becomes possible → **ABCD method**

Determination of background

– 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*

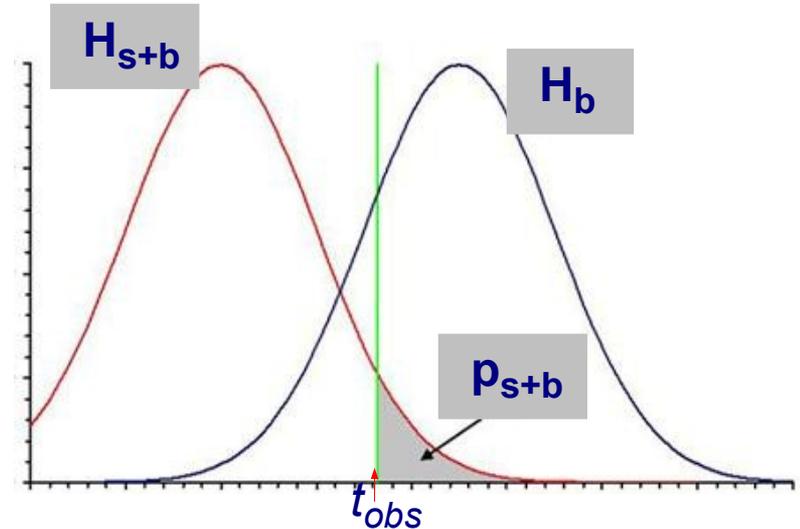
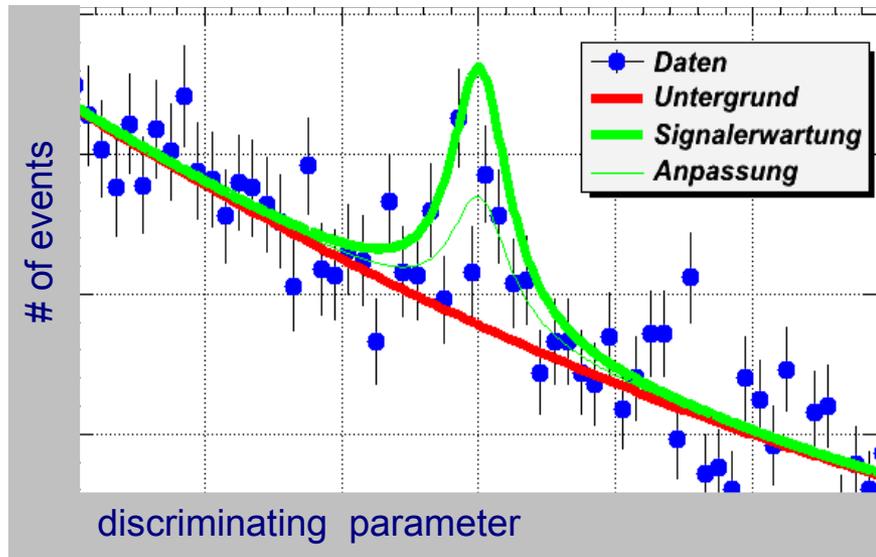
Example: invariant mass of two unlike-sign particles,
combinatorial background from sample with like-sign particles.

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

Statistical analysis

The Problem: an excess of observed events can have two sources:

1. signal in addition to expectation
2. a statistical upward fluctuation or insufficient understanding of background distribution (systematic error)



from the statistical view point, a **Hypothesis Test** H_{S+B} vs. H_B

Definition of a suitable teststatistic t as a function of the data: t_{obs}
calculation of probabilities

$p_{S+B} = \text{Prob}(t > t_{obs} | H_{S+B})$
resp. $p_B = \text{Prob}(t < t_{obs} | H_B)$
„p-values“ w.r.t. of S+B resp. B-only hypothesis

To postulate the observation of a **new signal**, background fluctuations must be excluded with very high probability !

Ein Beispiel: Würfel



Dodekaeder,

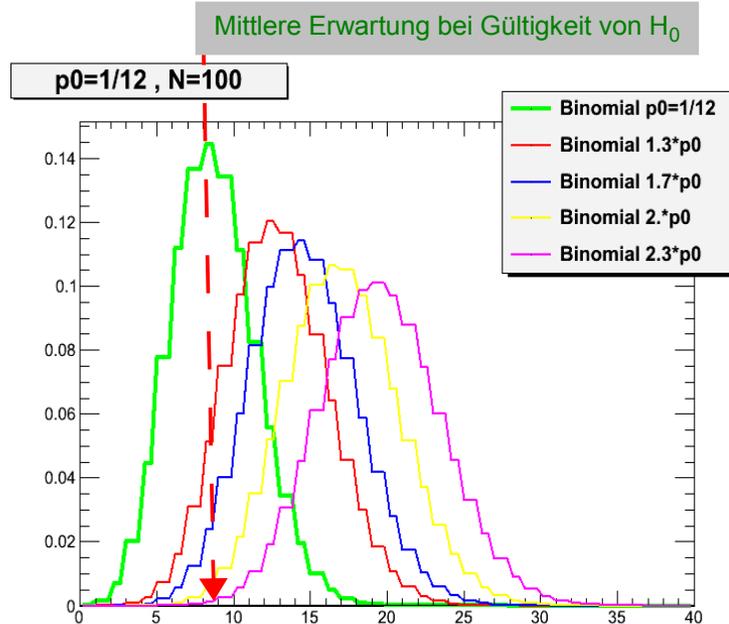
Frage: funktioniert so etwas als „Würfel“ im Spiel ?

Statistischer Test mit einer Zahl (z.B. 12), sollte mit Wahrscheinlichkeit $p=1/12$ auftreten

Nullhypothese, $H_0: p_{12} = 1/12$

Alternative, $H_a: p_{12} = f \times 1/12$

Teststatistik: Zahl der geworfenen 12en in $N=100$ Versuchen, n_{12} , folgt Binomialverteilung $B(n_{12}; N, 1/12)$



verschiedene Alternativen, $f=1.3$ bis $f=2.3$

Was können wir über f sagen ?

- Festlegen eines Konfidenzniveaus, z.B. 5%
- Bestimme f_{95} , so dass Fläche links von n_{12} unter der Verteilung 5% ist
- $\rightarrow f < f_{95}$ mit 95% Konfidenzniveau „beobachtete Grenze“

Legen wir statt der Beobachtung n_{12} die rote Linie (Erwartungswert $\langle n_{12} \rangle$ unter H_0) zu Grunde, erhalten wir die „erwartete Grenze“

Bedeutung für Higgs-Suche:

$\langle n_{12} \rangle$: erwartete Standard-Prozesse

f : „Signalstärke“ für Higgs-Beitrag

A complication:

expected signal and background processes and event selection are affected by additional uncertainties

→ nuisance parameters

need to treat these properly in the interpretation of the parameter of interest (e.g. the production rate of a new process)

Two main methods exist:

- *profile likelihood*
- *marginalization*
of the Bayesian posterior probability density

Nuisance parameters with profile likelihood

Minimization of $-2 \ln \mathcal{L}(\mathbf{x} | \theta, \lambda_i)$ w.r.t. all λ_i

for different values of the parameter of interest θ (e.g. # of signal events)

1. profile likelihood

provides all relevant quantities:

Best estimator: $\hat{\theta}$

at minimum of likelihood

Square of signal significance:

$$-2(\ln \mathcal{L}(N_s = 0) - \ln \mathcal{L}(N_s = \hat{\theta})) \\ = 2 \ln \frac{\mathcal{L}(N_s = \hat{\theta})}{\mathcal{L}(N_s = 0)} = 2 \ln Q = s^2$$

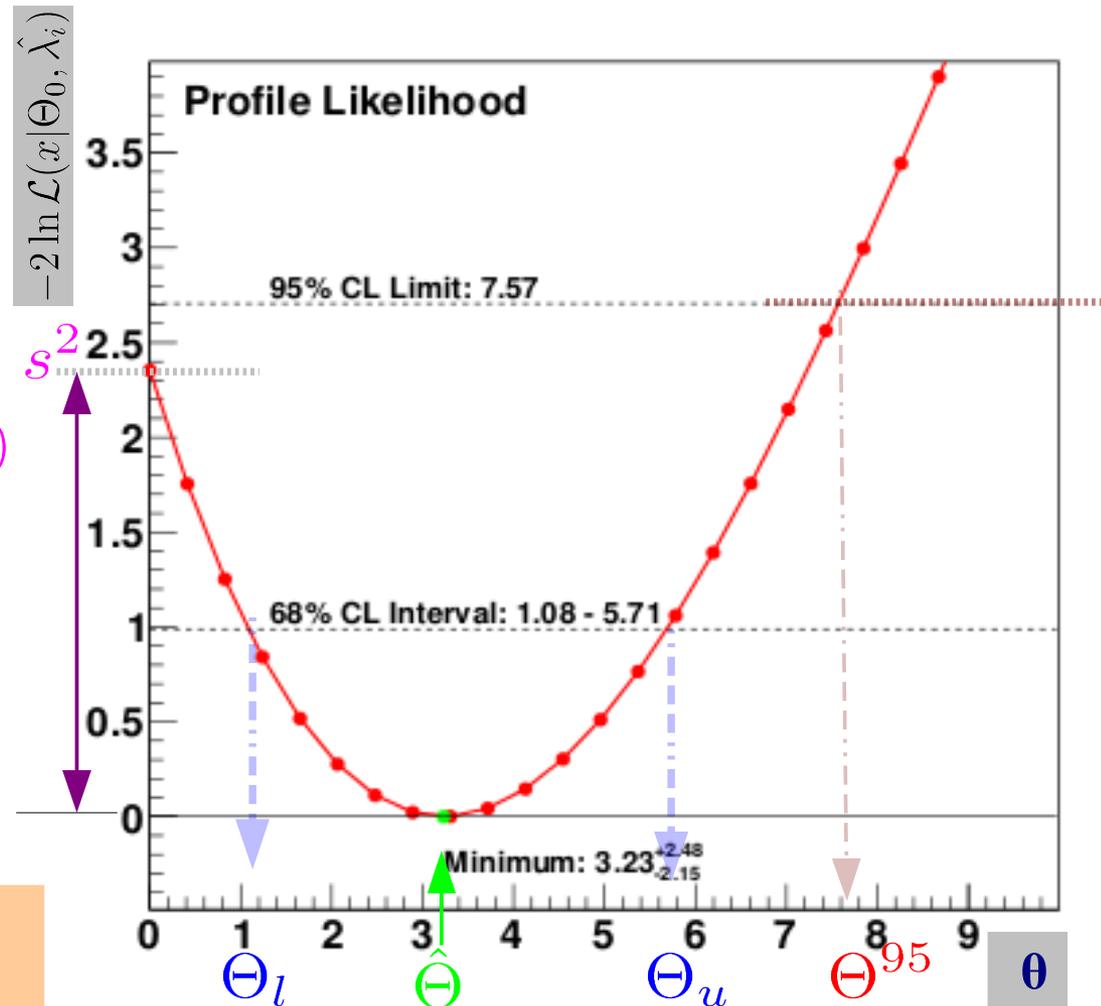
68% confidence interval:

$$-2\Delta \ln \mathcal{L}(\theta_{l,u}) = 1$$

95% limit:

$$-2\Delta \ln \mathcal{L}(\theta^{95}) = 1,645^2 = 2,71$$

Problem, if minimum near or close to physical boundary
– unphysical for number of signal events



Nuisance parameters via marginalisation

Integrate out the nuisance parameters in the Bayesian posterior probability density:

$$P(\theta) = \int \dots \int P(\theta, \lambda_i) d\lambda_1 \dots d\lambda_n$$

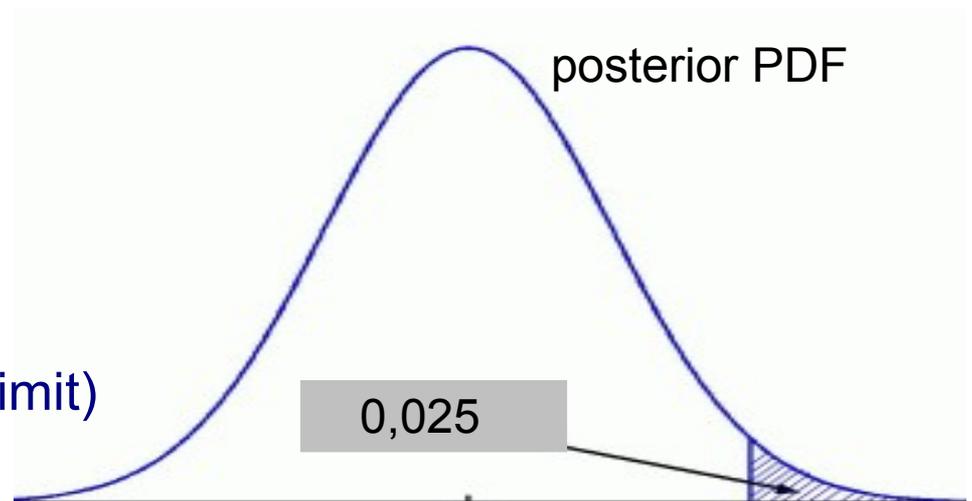
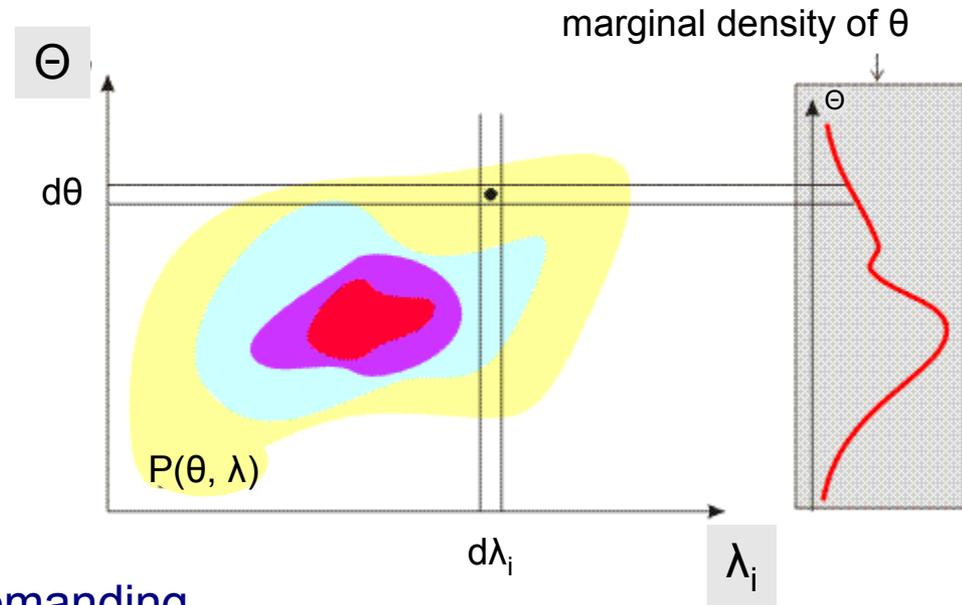
Multi-dimensional integral may be very demanding,
Markov-Chain Monte Carlo often used

confidence intervals, in Bayesian statistics often called

“credibility intervals”,

by determination of the appropriate quantiles of the marginalized posterior PDF

(e.g. one-sided upper limit)



Teststatistic for LHC Higgs search

for Higgs search: use **likelihood ratio** as teststatistic

profile likelihood w.r.t. the signal strength μ ($\mu=0$: no signal, $\mu=1$: nominal signal)
normalized to the global maximum of likelihood

**teststatistic
for limits**

signal strength

Set of nuisance parameters (= systematic uncertainties)

$$q_\mu = -2 \ln \frac{\mathcal{L}(\text{data} | \mu, \hat{\Theta}_\mu)}{\mathcal{L}(\text{data} | \hat{\mu}, \hat{\Theta})}, \quad 0 \leq \hat{\mu} \leq \mu$$

condition to ensure
 $\mu \geq 0$ and a
one-sided limit

best-fit values of all parameters

data distribution

$\hat{\cdot}$:= value of parameter that maximizes likelihood

$$\mathcal{L}(\text{data} | \mu, \Theta) = \prod \text{Poisson}(N_i | \mu \cdot s_i(\{\Theta\}) + b_i(\{\Theta\})) \cdot p(\{\tilde{\Theta}\} | \{\Theta\})$$

determination of the **distribution of q_μ** , $f(q_\mu | \mu)$, for background ($\mu=0$) resp. signal hypothesis ($\mu \neq 0$), via pseudo-experiments or asymptotic formulae in the limit of large data sets

Statistical Analysis

next:

- determination of p-values:

$$p_{\mu} = \text{Prob}(q_{\mu} > q_{\mu}^{\text{obs}} | \mu \cdot s + b)$$

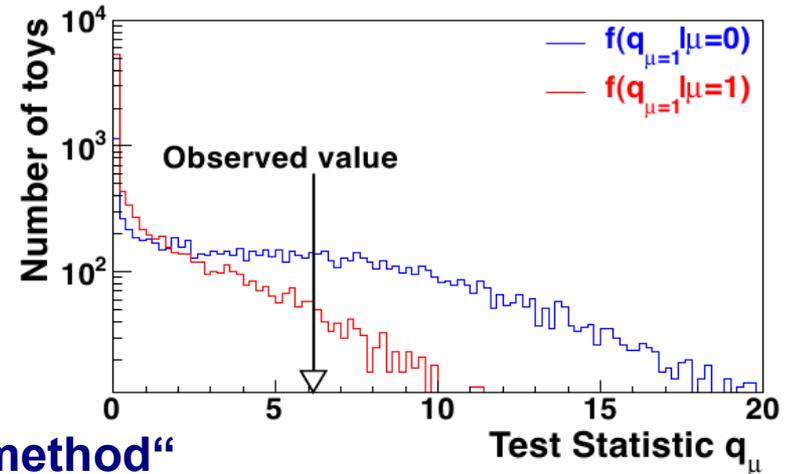
$$p_0 = \text{Prob}(q_{\mu} > q_{\mu}^{\text{obs}} | b)$$

- Calculation of confidence level with „CL_s method“

$$\text{CL}_s = \frac{p_{\mu}}{p_0}$$

robust against downward fluctuations of background

CL_s quantifies agreement with signal hypothesis



- for $\mu=1$, $\text{CL}_s = \alpha$ a Higgs Boson is excluded with confidence level $(1 - \alpha)$
convention: $\alpha=0.05$, exclusion 95% CL.
- usually:** specify value of μ that is excluded at 95% CL

perform pseudo experiments to determine **expected limit**,

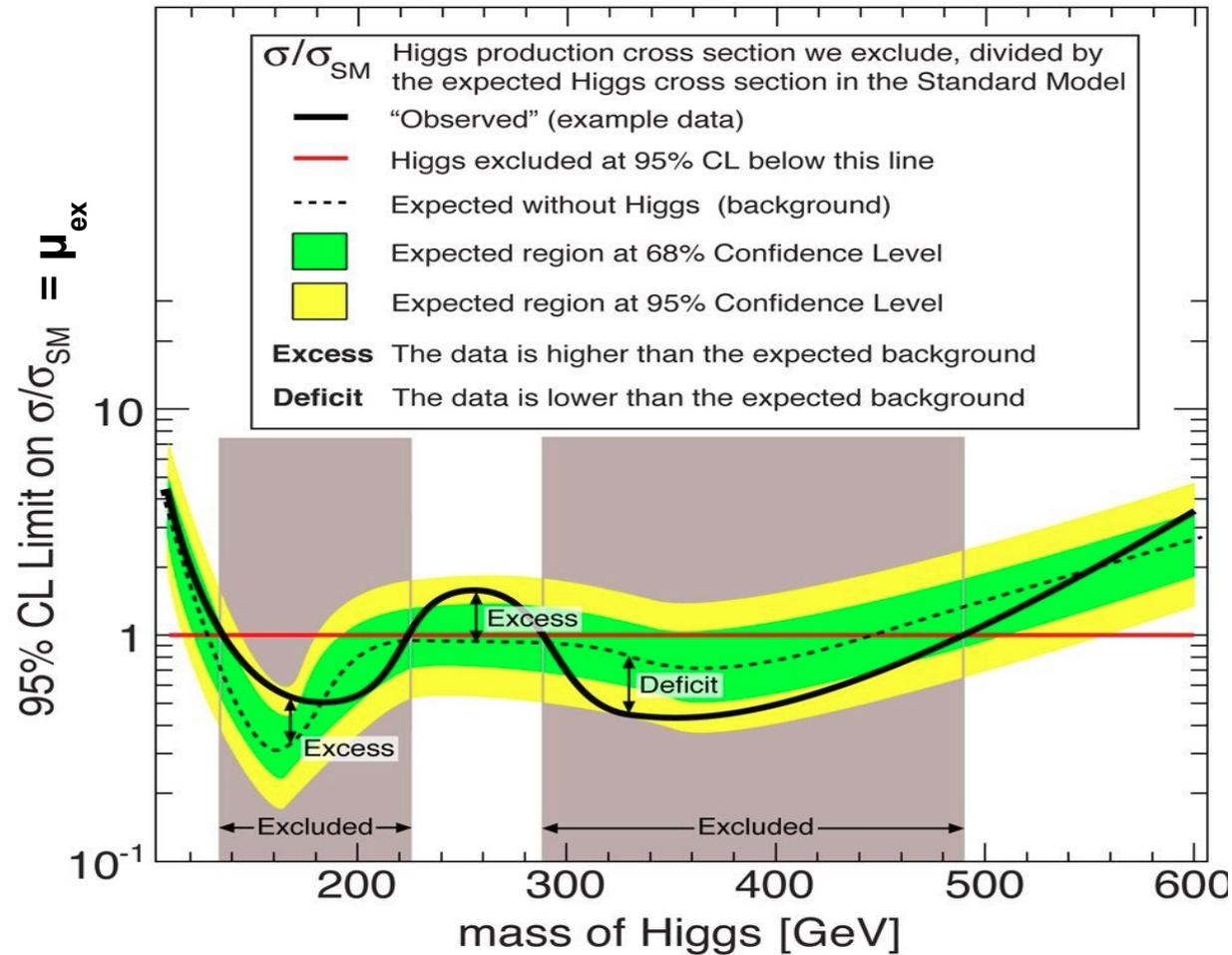
- e. the median of the distribution of obtained limits (dashed line),
and the regions for 68% („1 σ “, green band) and 95% („2 σ “, yellow band)

Statistical Analysis

- repeat all of the above for different values of Higgs mass



- Plot of the signal strength μ_{ex} excluded at **95% CL**
 - for data (black line)
 - for the median and 68%- and 95%-regions of the limit distribution obtained from the pseudo experiments (dashed line and green resp. yellow bands)



Higgs boson excluded

for $\mu_{\text{ex}} < 1$ (below red line, grey bands)

Statistical Analysis: Significance of a Discovery

If a signal cannot be excluded, what is the “significance” of a possible discovery?

Concept of “local significance”:

assume that N events have been observed over an expected background N_b

→ number of signal events: $N_s = N - N_b$

compare N_s with statistical fluctuations of background (N_b):

in Gaussian limit ($N_b > \sim 50$)

$$S = \frac{N_s}{\sqrt{N_b}}$$

– measure of a signal excess in terms of “number of sigmas” (“z-value”)

“The observed signal is S times larger than the standard deviation of the expected background fluctuations”

Statistical Analysis: Significance of a Discovery

typically, N_s and N_b are small numbers if a signal is just being discovered;

→ use Poisson statistics $P_0(N; N_b) = \frac{1}{N!} N_b^N e^{-N_b}$ background-only hypothesis

$$P_1(N; N_s + N_b) = \frac{1}{N!} (N_s + N_b)^N e^{-(N_s + N_b)}$$

and take logarithm of likelihood ratio:

s+b hypothesis

$$2 \ln (\mathcal{L}_1 - \mathcal{L}_0) = 2 \ln Q = 2 \left(N \ln \left(1 + \frac{N_s}{N_b} \right) - N_s \right)$$

and assume $N_s = N - N_b$

“in the asymptotic limit”

$2 \ln Q$ can directly be interpreted as the z-value of the observation!
(as it is the twice the difference in log-likelihood for $N_s = N - N_b$ and $N_s = 0$)

works also if observation is made in many bins,

$$N \rightarrow \sum N_i : 2 \ln Q = \sum 2 \ln Q_i$$

full treatment: define suitable test statistic ($q=q_0$ instead of $q=q_\mu$ for LHC),
determine distribution of q under signal and background hypotheses
and calculate p-values.

Statistical Analysis: Significance of a Discovery

If a signal cannot be excluded, what is the “significance” of a possible discovery?

Concept of “local significance”:

assume that N events have been observed over an expected background N_b

→ number of signal events: $N_s = N - N_b$

compare N_s with statistical fluctuations of background (N_b):

in Gaussian limit ($N_b > \sim 50$)

$$S = \frac{N_s}{\sqrt{N_b}}$$

– measure of a signal excess in terms of “number of sigmas” (“z-value”)

“The observed signal is S times larger than the standard deviation of the expected background fluctuations”

