

Use of Monte Carlo simulation at the LHC

Roberto Covarelli (*University / INFN of Torino*)

XIII meeting on B-Physics - Napoli, 23 May 2017

R. Covarelli

Outline

- ▶ MC tools widely used at the LHC
 - ▶ Matrix-element generators and parton showers
 - ▶ MC for B-physics
- ▶ Assessing theory uncertainties from MC
 - ▶ Matrix-element reweighting
 - ▶ Recent works on parton-shower uncertainties
 - ▶ PDFs and underlying-event tunes
- ▶ Technical aspects
 - ▶ General issues of large-sample generations
 - ▶ Issues specific to B-physics

Monte Carlo tools

Use of MC in pp experiments

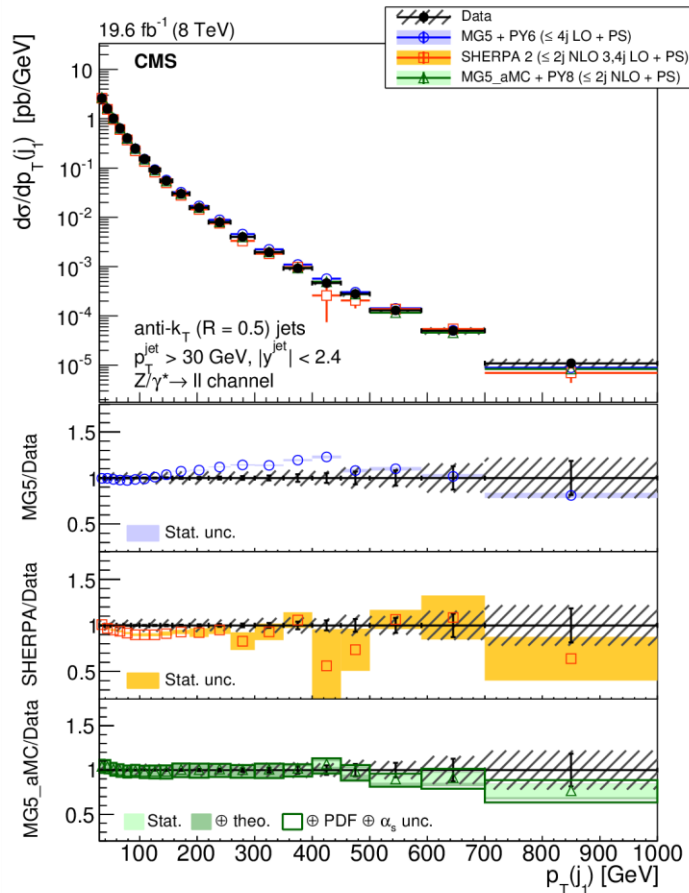
► 3 main ways of exploitation:

- 1) **Signal acceptance and efficiency**: associated uncertainties are typically **very small** (cancellation in ratios). **Severe mismodeling issues** can happen **only** in case of **large changes in kinematics**, e.g.:
 - LO vs. NLO $gg \rightarrow H$ production
 - Simulation vs. missing simulation of $g \rightarrow b \bar{b}$... etc.
- 2) **Background estimation** for rare signal searches. Should be mitigated where possible by use of data control samples (sidebands etc.) otherwise **MC uncertainties enter in full**
- 3) **Traning samples for multi-variate analyses** (MVA): propagation of MC uncertainties / mismodeling **not easy to assess**, especially in case of non-analytical MVA (e.g. BDT)

MC generator categorization

- ▶ Order in QCD:
 - ▶ LO or NLO (first NNLO generators appearing in recent years just for specific processes, e.g. $gg \rightarrow H$)
- ▶ Purpose:
 - ▶ Full-event generators (e.g. Pythia8, Herwig7, SHERPA) or just for some stages of MC generation:
 - ▶ Matrix-element, i.e. up to parton-level (e.g. MadGraph5_aMCatNLO)
 - ▶ Parton showers
 - ▶ Applications for particle decays (e.g. EvtGen, Tauola)
 - ▶ Use of staged generation requires special treatment of physics effects («matching») and software interfacing
- ▶ Automation:
 - ▶ Can only simulate specific processes (e.g. POWHEG, Pythia8) or a generic $pp \rightarrow X$ process can be defined by user (e.g. SHERPA, MadGraph5_aMCatNLO)

Matrix-element generators

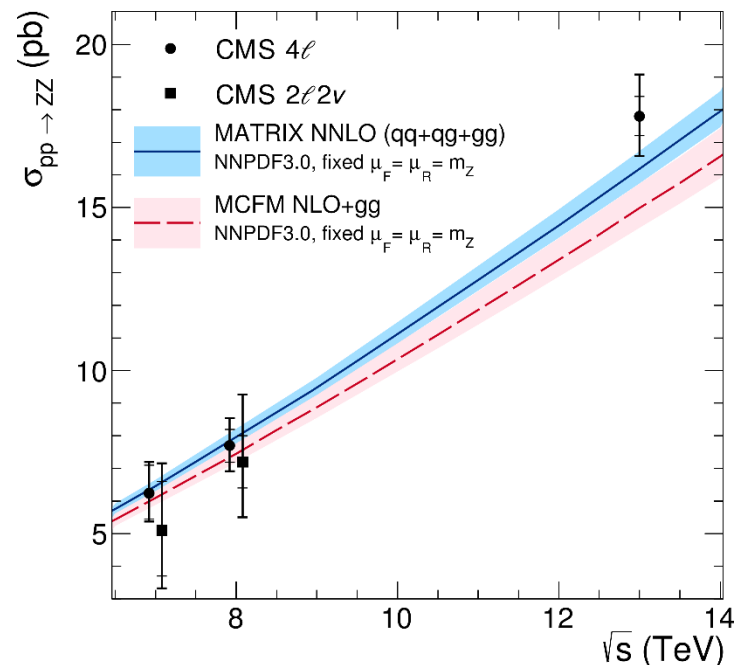


CMS-SMP-14-013, submitted to JHEP

- ▶ **NLO QCD event generation is now the standard** (SHERPA, POWHEG, MadGraph5_aMCatNLO)
- ▶ **Specific cases where LO is still used:**
 - ▶ Final states with large particle multiplicities (e.g. $V + 4\text{jets}$)
 - ▶ NP signals where NLO calculations are not available/implemented
 - ▶ Final states with particularly complicated kinematics (e.g. $pp \rightarrow G^*$ with full tensor structure)
 - ▶ ...etc.

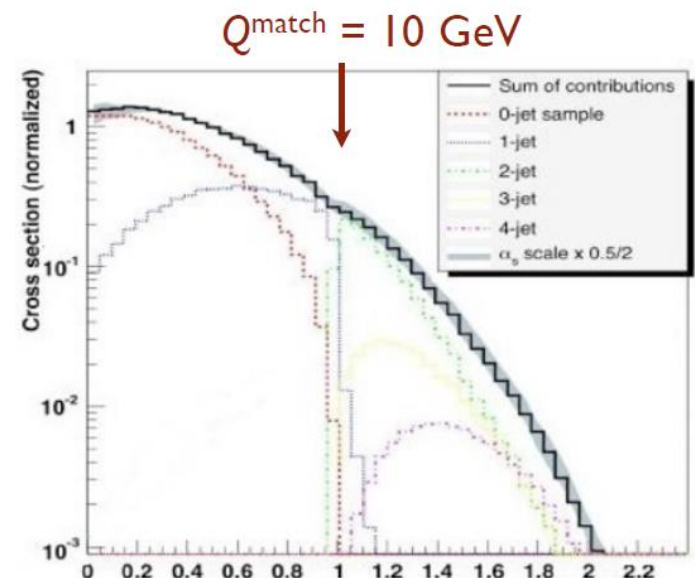
Higher-order QCD corrections

- ▶ NLO ensures:
 - ▶ Reasonable QCD scale uncertainties $\rightarrow o(10\%)$ for $V+jets$ / $t\bar{t}bar$
 - ▶ Quite accurate description of kinematics when matched to a parton shower (NLO+PS), except for some rare processes (e.g. those with a box diagram at the lowest order)
- ▶ In some cases where higher-order calculations are needed, experiments use k -factors ($k = \sigma_{NNLO}/\sigma_{NLO}$)
 - ▶ Integrated, i.e. scale cross-section by a constant number
 - ▶ Differential vs. specific process variables (event «reweighting» when filling MC distributions)



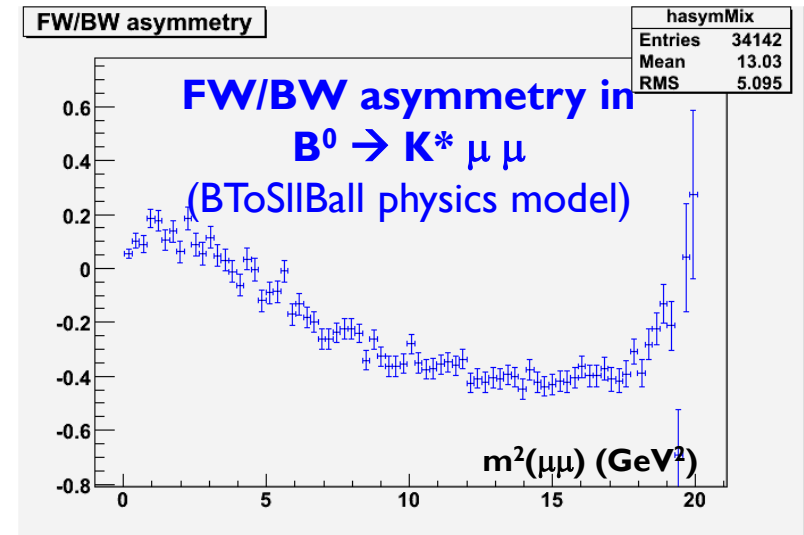
Parton showers

- ▶ Most matrix-element generators (basically all except SHERPA) need **matching with an external parton shower** to produce finalized pp events
- ▶ **Pythia8** and **Herwig++** (recently updated to Herwig7) are the two most widely used in LHC
- ▶ Matching is not a trivial task as it must deal with **jet double counting**:
 - ▶ An event with e.g. V+2jets can have two real g/q emissions from the ME, or one from ME and one from PS etc.
 - ▶ Theoretical recipes exist (**MLM** for LO, **FxFx** for aMCatNLO, **emission veto** for POWHEG etc.) and are included in PS tools



External decay programs

- ▶ PS tools provide their own particle decay utilities
 - ▶ For specific signals, external decay programs can be used
- ▶ EvtGen managing complex B (and hadron) decays
 - ▶ Created in BaBar, now maintained by the LHCb Warwick group
 - ▶ <http://evtgen.warwick.ac.uk/>
 - ▶ Spin amplitudes, CPV, automatic PHOTOS interface, etc.
- ▶ Interfacing to other tools is apparently trivial
 - ▶ Veto particle decays in other tools, if known to EvtGen
- ▶ In practice, many subtle issues
 - ▶ Coherence of PDG data in the two generators
 - ▶ Signal particles could be within an EvtGen decay tree
 - ▶ Fall-back to Pythia for high-multiplicity decays etc.



Use of MC in **B-physics** experiments

► 3 main ways of exploitation?

- 1) **Signal acceptance and efficiency**: associated uncertainties are typically very small (cancellation in ratios). **Severe mismodeling issues** can happen **only** in case of **large changes in kinematics**

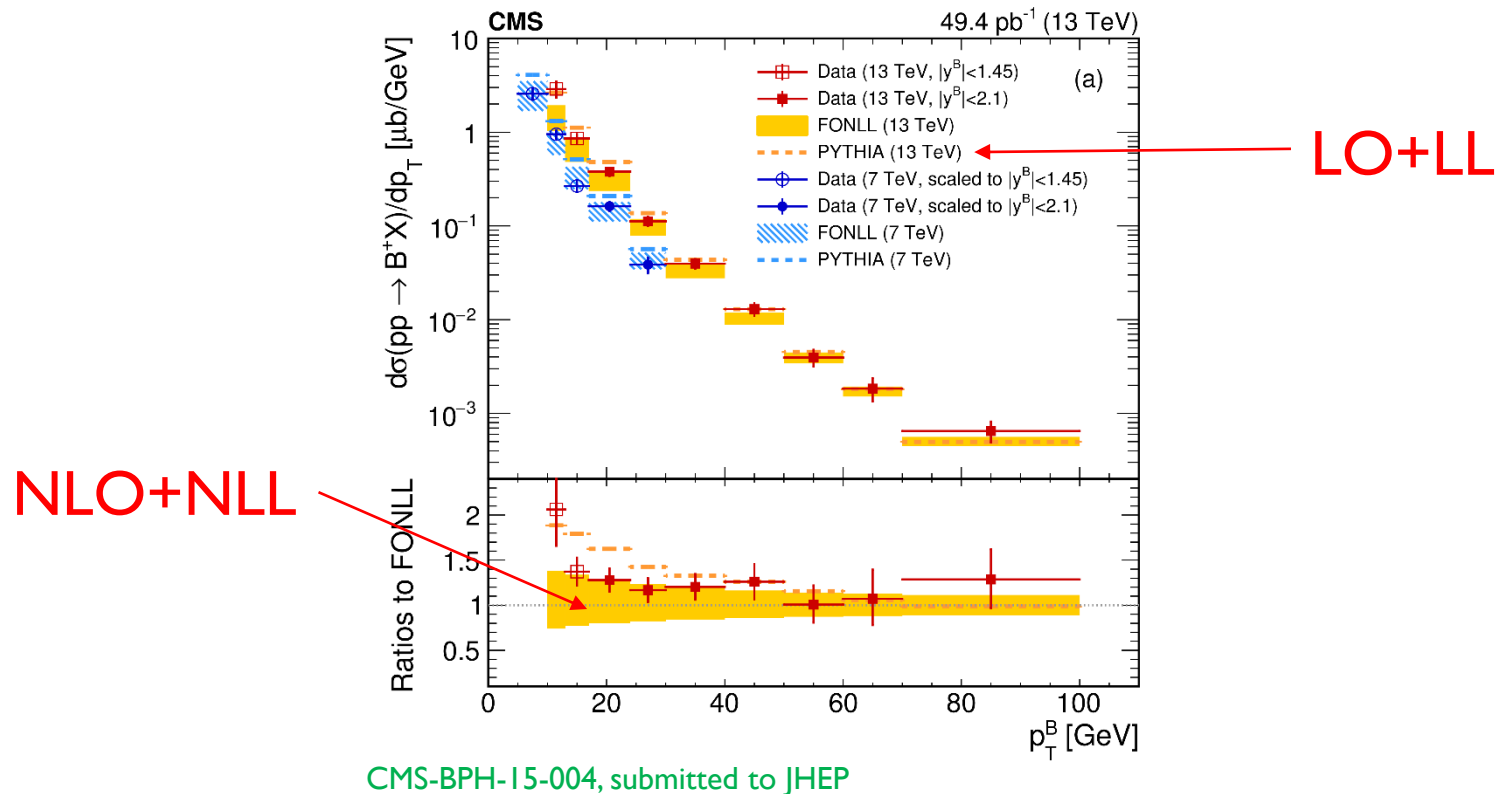
► LO vs. NLO: really important?

- ~~2) **Background estimation** for rare signal searches. Should be mitigated where possible by use of data control samples (sidebands etc.) otherwise **MC uncertainties enter in full**~~

ALMOST NEVER USED

- 3) **Traning samples for multi-variate analyses** (MVA): propagation of MC uncertainties / mismodeling **not easy to assess**, especially in case of non-analytical MVA (e.g. BDT)

HQ production: Pythia8 vs. FONLL

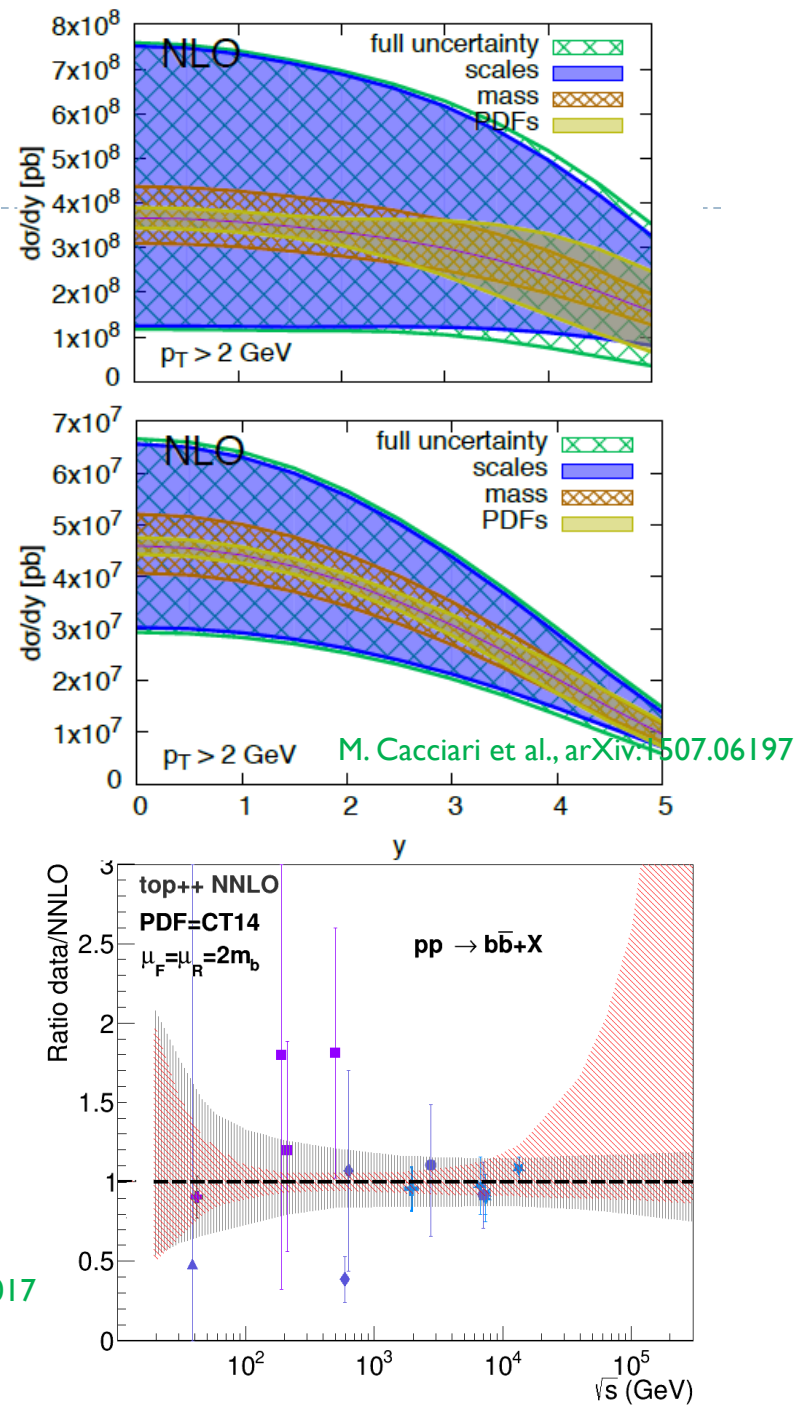


- ▶ In many practical cases, LO sufficient for B production kinematics
 - ▶ Generation speed is more important, see next slides

HQ production: beyond NLO?

- ▶ Scale uncertainties still not so small, especially for $c\bar{c}$
 - ▶ Dominant vs. other sources (gluon PDFs, b and c mass)
- ▶ Ratios at different \sqrt{s} remove such uncertainties
- ▶ First attempts to use **Top++ NNLO** computator (Mitov and Czakon, 2013) applied to other quarks
 - ▶ Encouraging, but NNLL seems not possible yet (essential for low p_T production)

D. D'Enterria, talk at Moriond 2017



Monte Carlo uncertainties

Reweighting

- ▶ Nowadays **matrix-element generators** can compute variations of certain calculation parameters **on an event-by-event basis** and write the result as a **per-event weight**
 - ▶ Most common
 - ▶ **QCD renormalization and factorization scales** (7 to 9 variations corresponding to $\times 0.5$, $\times 2$ in all possible combination)
 - ▶ **PDF variations** (can be 100's of them)
 - ▶ Written in **LHE files** in standard form

```
<weightgroup combine="envelope" name="scale_variation">
```

```
<weight id="1001"> muR=1 muF=1 </weight>
```

```
<weight id="1002"> muR=1 muF=2 </weight>
```

```
<weight id="1003"> muR=1 muF=0.5 </weight>
```

```
<weightgroup combine="hessian" name="PDF_variation">
```

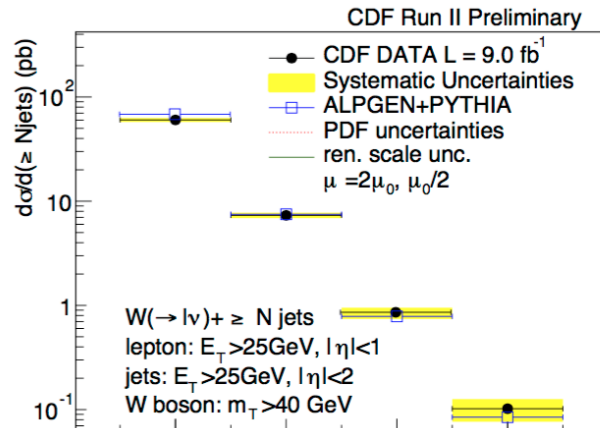
```
<weight id="2001"> PDF set = 260001 </weight>
```

```
<weight id="2002"> PDF set = 260002 </weight>
```

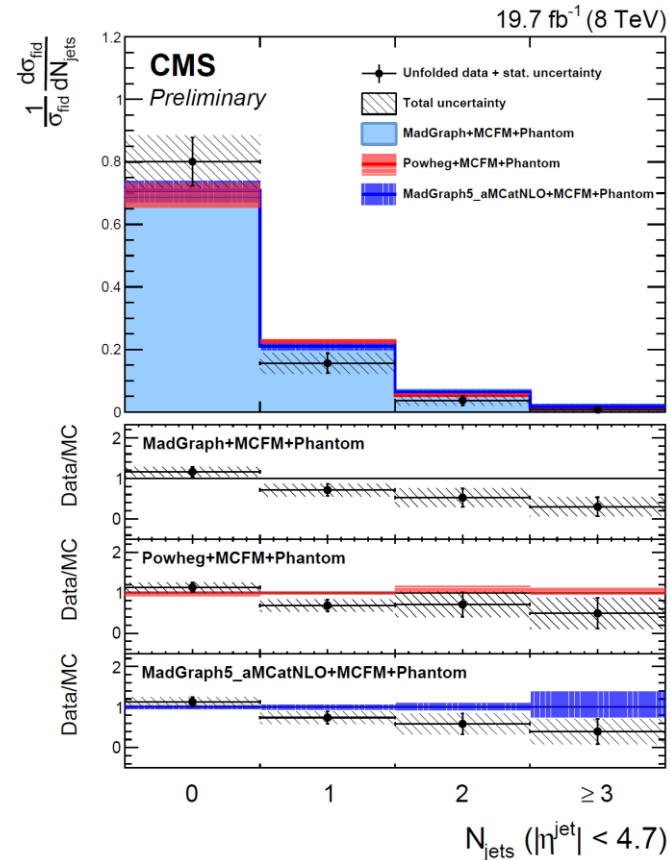
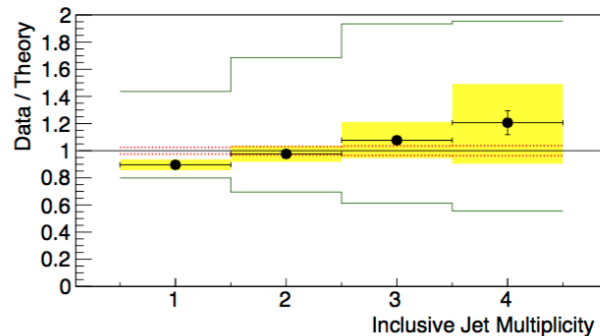
```
<weight id="2003"> PDF set = 260003 </weight>
```

- ▶ Crucial for **use in experiments!**
 - ▶ Would be impossible to produce 100's of samples, one per variation

Is this the full theory uncertainty?



CDF Public Note 11167



CMS-
PAS-
SMP-15-012

► LO vs. NLO cannot be responsible for such a difference

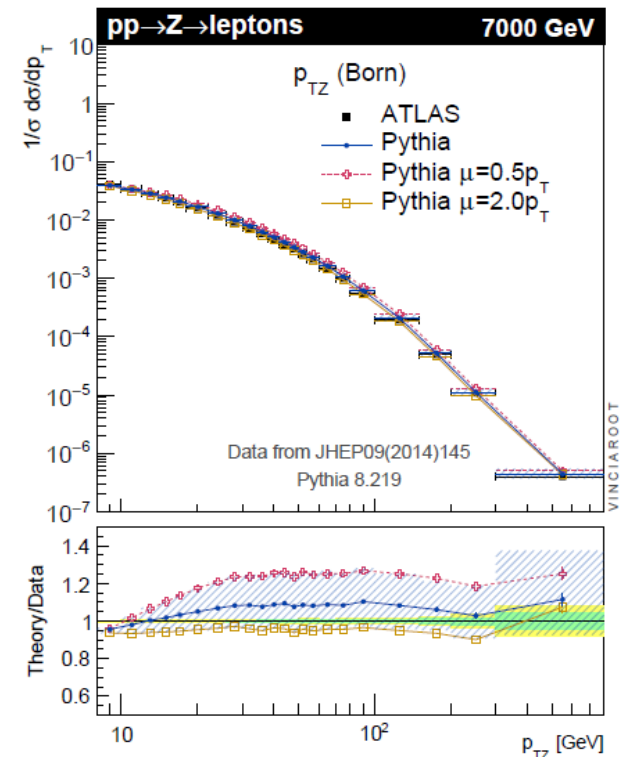
Parton shower uncertainties

- ▶ The matrix-element generator uncertainty is not the end of the story → **parton shower uncertainties can be large**
 - ▶ Especially important if MC used for **background estimation**
 - ▶ Even more if variables are heavily affected by showering/fragmentation/hadronization, e.g. **jet substructure**, **q/g discrimination** etc.
- ▶ Until 2016, mixed recipes used in experiments
 - ▶ Totally neglected
 - ▶ Used Pythia vs. Herwig difference in relevant variables
 - ▶ Used ME-consistent up and down μ_R variations
- ▶ In very recent years, substantial progress on the matter
 - ▶ **Improved PSs (DIRE,VINCIA)** EPJC 76 (2016) 589, arXiv:1705.00982
 - ▶ Three simultaneous works from authors of Pythia, Herwig and SHERPA clarifying the matter of **PS uncertainties**

arXiv:1605.01338, arXiv:1605.08352, arXiv:1606.08753

The Pythia8 example

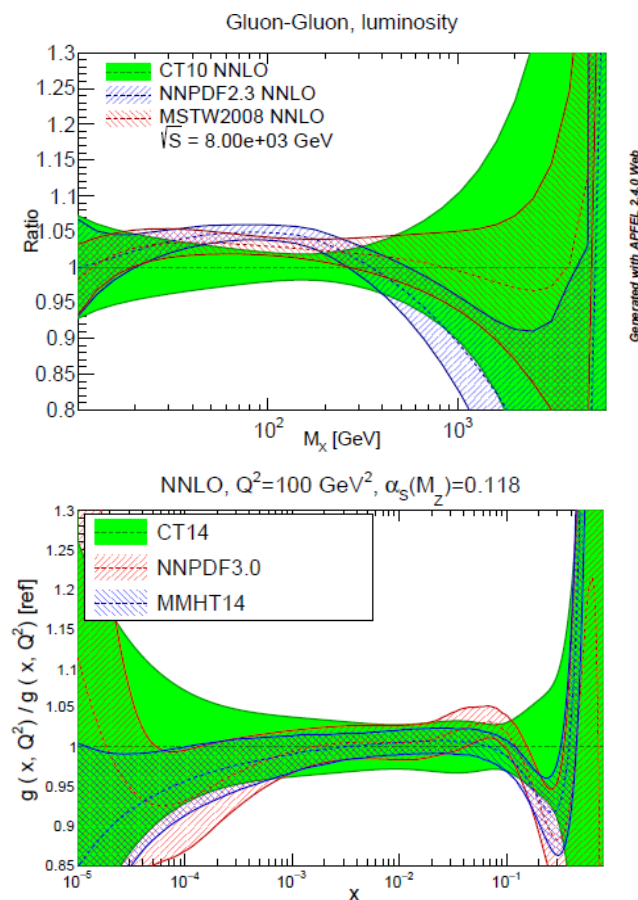
- ▶ More difficult to have **per-event weights in PSs** because of their literal «MC» nature (acceptance-rejection algorithms)
- ▶ **In Pythia8 now available** from a specific «recycling» of failed trials in the Sudakov veto algorithm
 - ▶ **Validation successful in Z+jets data**
- ▶ Also DGLAP splitting-kernel uncertainties studied
 - ▶ Become negligible when ME corrections are applied
 - ▶ Correlation w.r.t. ME uncertainties well understood?
- ▶ **Not yet used** in experiments



arXiv:1605.08352,
Phys. Rev. D 94, 074005 (2016)

PDFs

- ▶ Recent updates (in 2014 and 2017) of PDF fits in collaborations, including LHC data
 - ▶ Brought to significantly improved agreement on gluon PDFs
- ▶ Recipes used in practice in MC generation
 - ▶ Used NNPDF3 (become available first)
 - ▶ Store weights for at least:
 - ▶ All NNPDF error sets (100 at $\alpha_s = 0.118 + 2 \alpha_s$ variations of ± 0.001)
 - ▶ MMHT and CT
 - ▶ NNPDF can be replaced with PDF4LHC15 (preferably with the 100 error set, unless shown that using just the reduced set of 30 makes no difference)
- ▶ Of course for precision measurements (e.g. W/Z cross-sections) the recommendation is still to publish giving the result for various PDF sets

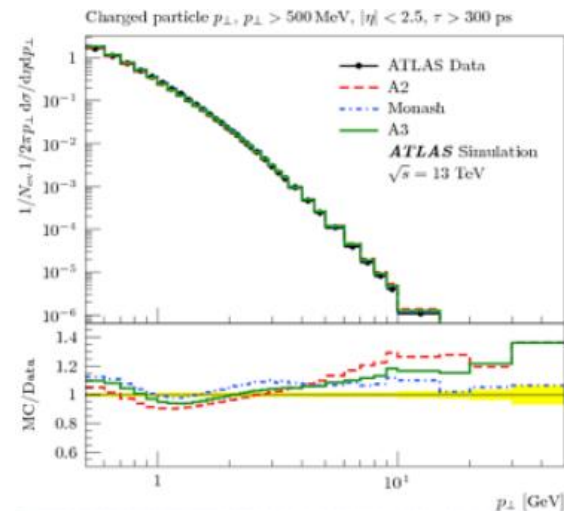
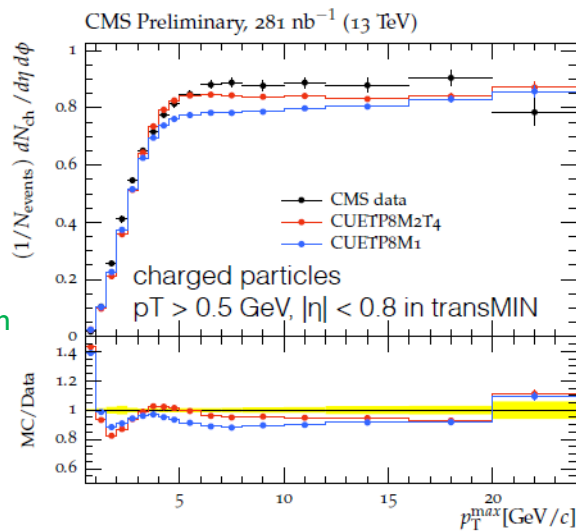


J. Phys. G: Nucl. Part. Phys. 43 023001 (2016)

Underlying event tunes

- ▶ Very significant work in both ATLAS and CMS to extract UE (DPS) tunes from Minimum Bias (4-jet, same-sign WW etc.) data

CMS-GEN-17-002,
Paper in preparation



ATL-PHYS-PUB-
2017-008

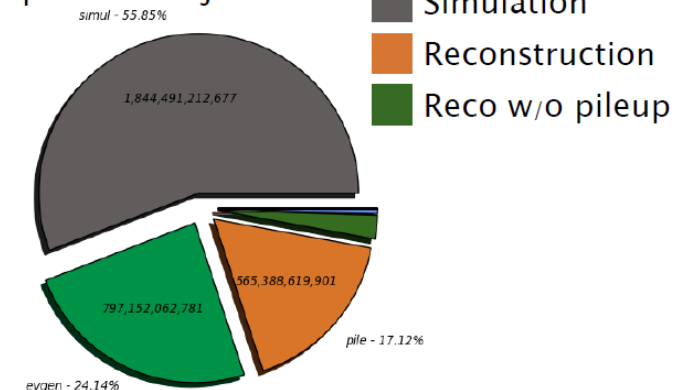
- ▶ Uncertainties estimated via «eigentunes»
- ▶ **Problem:** ATLAS/CMS produce MC samples with **its own** latest version of UE tunes → affects main-sample **and** PU simulation
 - ▶ In several cases, leads to impossibility of one-to-one comparisons

Technical aspects

MC production at LHC experiments

- ▶ Nowadays, MC «campaigns» (i.e. sets of samples targeting analysis of a specific LHC dataset) have typical sizes of **$\mathcal{O}(10\text{B})$ events**
 - ▶ By far, **the most CPU-expensive task in LHC computing**
 - ▶ Ideally, want to keep a MC/data ratio $\gg 1$
- ▶ **NLO generators changed the paradigm**
 - ▶ Before: Generation of the physical event has **negligible** CPU time/memory consumption w.r.t. other simulation steps (e.g. GEANT)
 - ▶ Now: Generation of the physical event can take a **substantial amount** of time/memory and can be impossible w/o some pre-computation steps

CPU consumption for good MC production jobs



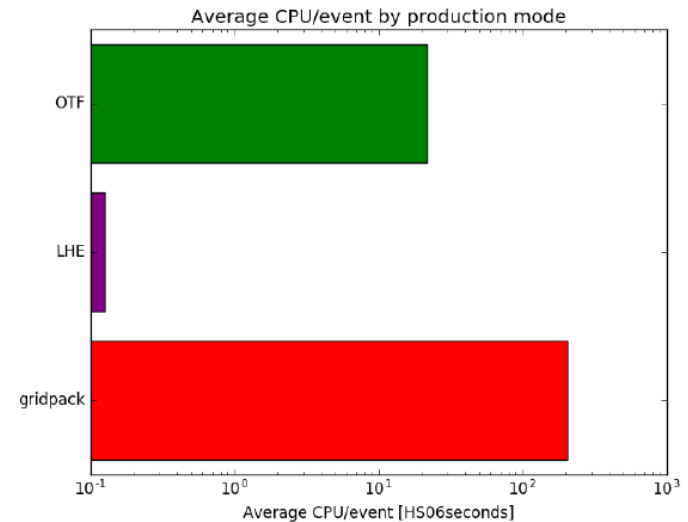
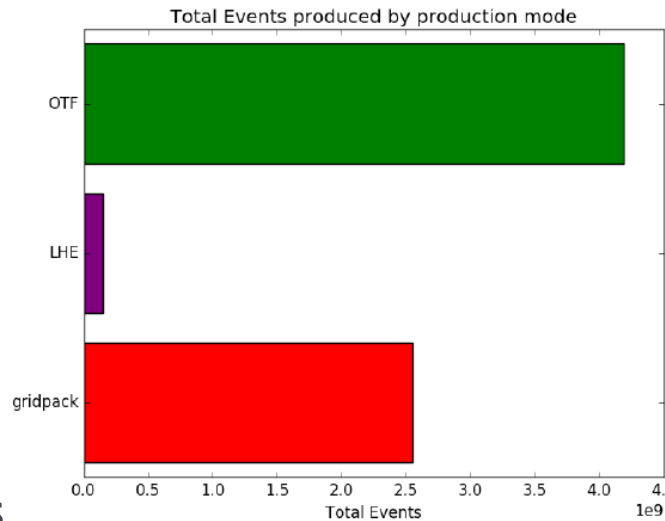
Josh Mc Fayden
@ATLAS-CMS MC workshop 2016

The «gridpack/SHERpack» idea

- ▶ When event generation is split in several jobs (a must for large generations) there can be some **pre-computation steps that are valid once for all**, e.g.:
 - ▶ **Calculation of diagrams and amplitudes** for automated ME generators
 - ▶ **Optimization of «grids»** for MC integration
- ▶ Idea: **run these steps in advance** and store all needed files in a large archive («gridpack») retrieved in each job
 - ▶ Code available from authors for **MadGraph / SHERPA**
 - ▶ **Written by experimentalists for POWHEG**
- ▶ **Tested and working**
 - ▶ **Issue:** the «gridpack» creation step does not fit LHC GRID-based computing models
 - ▶ Must be performed by **a single user** using batch queues → can take a long time for NLO processes with high-multiplicity final states

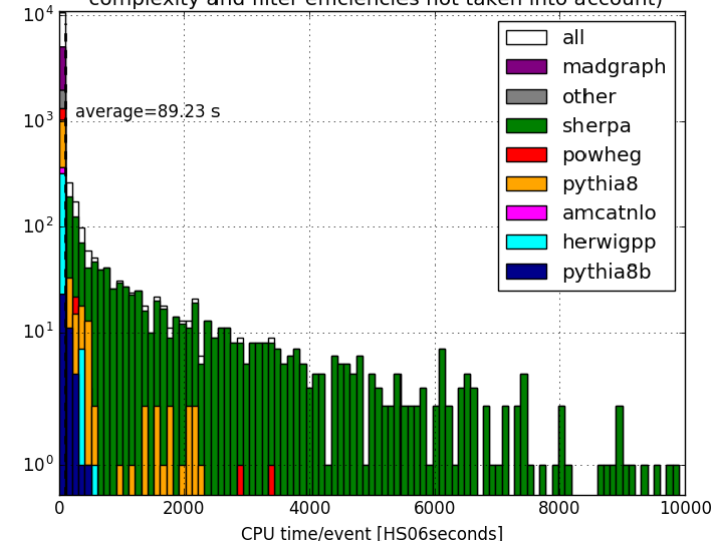
The «gridpack/SHERpack» idea (cont'd)

- ▶ RunII MC gridpack use in **ATLAS**
 - ▶ More imbalance vs. gridpacks in **CMS**



- ▶ Time/event split by MC generator in **ATLAS**
 - ▶ Better picture in **CMS** because of less frequent SHERPA use

CPU time/event for 2015 MC event generation at $\sqrt{s}=13$ TeV
(All physics processes included, correlations with process complexity and filter efficiencies not taken into account)

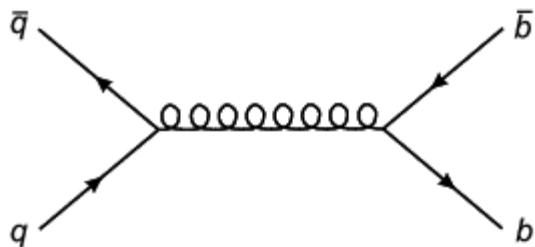


Josh Mc Fayden
@ATLAS-CMS MC workshop 2016

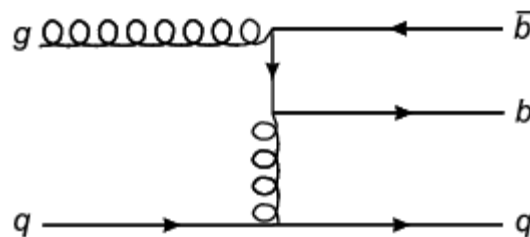
Specific issues for B-physics

- ▶ As we discussed in the previous slides, MC for B-physics is typically **LO+LL** (mostly **Pythia**)
 - ▶ DPS activated inside Pythia to estimate contribution to double Qqbar production
 - ▶ Interface to an external decay program (usually **EvtGen**) is a must for basically the whole set of samples
- ▶ **Generation of generic QCD is necessary**, because the pure $pp \rightarrow Q \bar{Q}$ generation misses **2 \rightarrow 3 processes (large σ !)** occurring only in subsequent stages of the PS

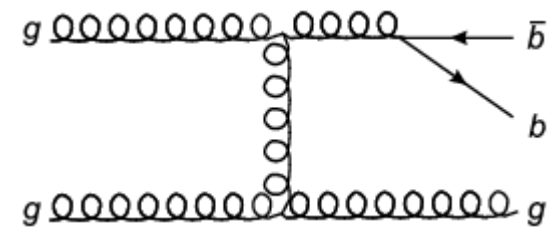
Tree-level



Flavor excitation



Gluon splitting

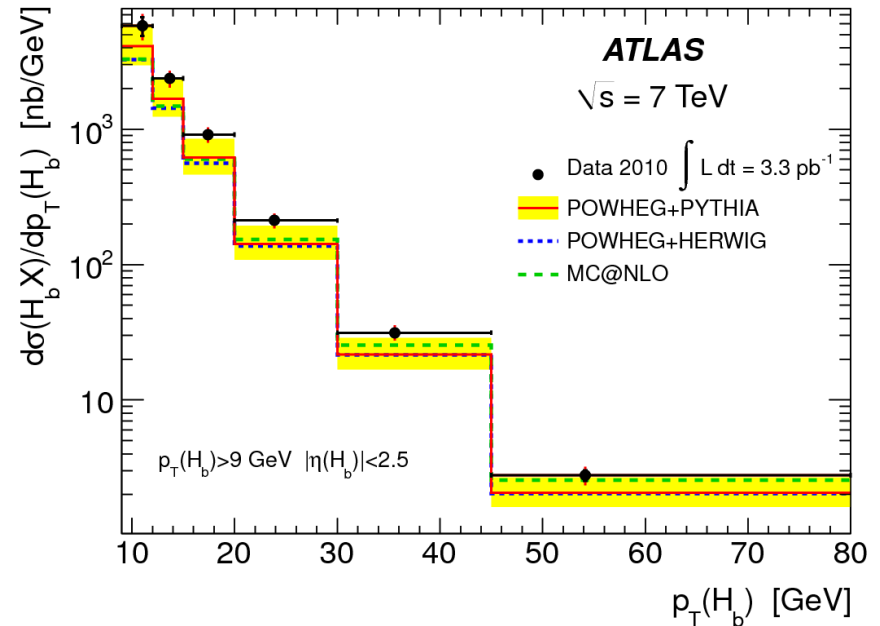


The generation efficiency problem

- ▶ Example: generate $B_s \rightarrow \mu\mu$ in CMS
 1. Start from generic QCD and find a $b\bar{b}$ pair: $\mathcal{O}(10^{-3})$
 2. b or \bar{b} fragmentation into B_s : $\mathcal{O}(10^{-1})$
 3. Decay to $\mu\mu$, can be forced in EvtGen: 100%
 4. Muons inside CMS trigger acceptance: $\mathcal{O}(10^{-2})$
- ▶ Generation efficiency for a detectable event is $\mathcal{O}(10^{-6})$
 - ▶ Timing of a Pythia8 QCD event is around 10 ms \rightarrow 10,000 s to obtain a single event, unacceptable for large productions
- ▶ Things are even worse for background processes e.g. $QCD \rightarrow \mu\mu$, since acceptance is much smaller
- ▶ Various workarounds found by experiments
 - ▶ Approximate cuts on b -parton kinematics (e.g. not at extreme $|\eta|$)
 - ▶ Re-hadronization of events, much faster than producing a new event \rightarrow Much easier in Pythia8 w.r.t. 6 («UserHooks»)
- ▶ Still not a completely settled issue

Moving to NLO?

- ▶ NLO generators for $pp \rightarrow Q \bar{Q}$ available for years
- ▶ And generally used for unfolded data/theory comparisons



- ▶ Could **not be enough in the description of gluon splitting** that can happen at a later stage of the PS
- ▶ In precision measurements where $X+b$ ($b\bar{b}$) is a major background (e.g. $VHbb$, $t\bar{t}Hbb$), the recipe is to **merge a NLO $X+b$ ($b\bar{b}$) sample** with a **second sample** where LHE-level b 's are vetoed and events are filtered on **the presence of a b in the PS**
- ▶ Does not solve the efficiency problem

Conclusions

- ▶ MC-generator evolution showed an **impressive rate** in recent years
 - ▶ Not just in the **available tools** but also in **input parameters** (e.g. updated PDFs)
- ▶ **Not always followed by timely application** in experiments
 - ▶ Generator-integration or CPU timing issues
 - ▶ Limited manpower
 - ▶ «Campaign» structure: full set of consistent MC events produced in one shot and then kept for long periods (considering the re-Reco option, can be a few years)
 - ▶ Complex applications of theory recipes difficult to incorporate into sample production workflows

Back up
