

Summary of LHC EWWG discussion on HEPData recommendations

Louie Corpe (UCL)

Talk to ALICE collaboration, 20 Jan 2020

Who am I?



- Post-doc at University College London, member of ATLAS collaboration
 - Previously at Imperial College, member of CMS between 2013-2017
 - Co-convenor for the ATLAS Generator Infrastructure and Tools subgroup of the Physics Modelling Group
- Giving this talk on behalf of the LHC Electroweak Working Group (Jets +Bosons):
 - We have been discussing how to propagate correlations and use them for tuning etc
 - Naturally lead to discussion on reviewing and agreeing conventions for HEPData uploads across experiments
 - Spoke to ATLAS/CMS/LHCb SM/Generators groups already : recommendations well received !

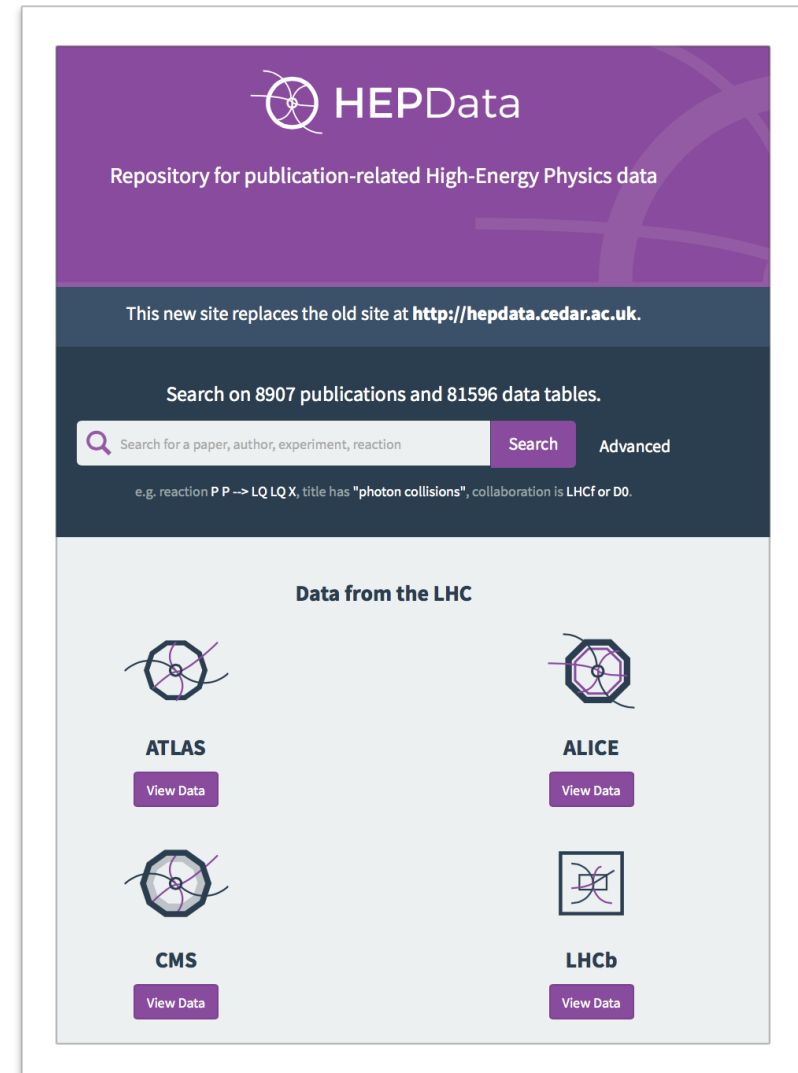
- Analysis preservation is increasingly important discussion in HEP
 - (in pp) new and ambitious goals of combinations, recasting, EFT fits etc...
 - many recommendations also apply to heavy ions too!
- Are we routinely storing enough information on HEPData to efficiently re-use the measurements we make at the LHC?

--> Not always! **small policy shifts can boost impact of analyses**
- Prompted by discussion on correlations, **LHCEWWG**: [Dec 18](#), [Feb 19](#), [July 19](#)
- Attempt to **formalise recommendations** and **document them in [note](#) to be agreed between LHC experiments**
 - Give **recommendations on conventions to follow** depending on what level of re-interpretation is needed

- I'm not a heavy ions expert, so I apologise if some of the recommendations I'll talk about today are not relevant to you!
- The purpose of this talk is also to gather feedback:
 - If you have suggestions or comments about what is/is not applicable to heavy ions, I'll be very grateful!
 - That way we can make this document useful for the heavy ion community as well as p-p

What are the current recommendations?

In this section I'll highlight some of the pitfalls of the current ALICE HEPdata uploads..



The screenshot shows the HEPData website interface. At the top, the HEPData logo is displayed, followed by the text "Repository for publication-related High-Energy Physics data". Below this, a dark blue banner states: "This new site replaces the old site at <http://hepdata.cedar.ac.uk>".

The main content area features a search bar with the text "Search on 8907 publications and 81596 data tables." and a search input field containing "Search for a paper, author, experiment, reaction". To the right of the search bar are buttons for "Search" and "Advanced". Below the search bar, a small example text reads: "e.g. reaction P P -> LQ LQ X, title has 'photon collisions', collaboration is LHCf or D0."

The bottom section is titled "Data from the LHC" and contains four entries, each with a detector logo and a "View Data" button:

- ATLAS**: ATLAS logo, View Data button
- ALICE**: ALICE logo, View Data button
- CMS**: CMS logo, View Data button
- LHCb**: LHCb logo, View Data button

- Good practice to define fiducial volume/ region of measurement (eg **Rivet Routine**)
 - This was only present in one HEPData entry I could find...
 - But hopefully recent release of Rivet 3.0.1 (many Heavy-Ion developments!) will help improve the situation?

- Give results with uncertainties in each bin. Separate stat vs syst uncertainties at minimum. **Stat/syst not enough to model correlations if re-interpretation is to be trusted!**

π^0 and η production at 8 TeV pp [[link](#)] but there are very few others with Rivet Routines

π^0 and η meson production in proton-proton collisions at $\sqrt{s} = 8$ TeV

The ALICE collaboration

Acharya, Shreyasi , Adam, Jaroslav , Adamova, Dagmar , Adolfsson, Jonatan , Aggarwal, Madan Mohan , Aglieri Rinella, Gianluca , Agnello, Michelangelo , Agrawal, Neelima , Ahmed, Zubayer , Ahmad, Nazeer

Eur.Phys.J. C78 (2018) 263, 2018

<https://doi.org/10.17182/hepdata.79044.v2>

Journal INSPIRE Resources

Rivet Analysis

Measurement of D0, D+, D*+ and D+S production in Pb Pb [[link](#)]

RE	P PB --> D0(Q=PROMPT) X			
SQRT(S)/NUCLEON	5020.0 GEV			
YRAP(RF=CM)(D)	-0.96 TO 0.04			
PT(D) [GEV]	d ² σ/d p _T dy [μb/GeV]			
0.0 - 1.0	22300.0 ±2.23e+03	stat	+1.65e+03 -1.69e+03	syst ±1.0% syst, uncertainty on branching ratio ±3.7% syst, uncertainty on integrated luminosity
1.0 - 1.5	31500.0 ±2.72e+03	stat	+3.85e+03 -3.93e+03	syst ±1.0% syst, uncertainty on branching ratio ±3.7% syst, uncertainty on integrated luminosity
1.5 - 2.0	29600.0 ±1.36e+03	stat	+2.70e+03 -2.80e+03	syst ±1.0% syst, uncertainty on branching ratio ±3.7% syst, uncertainty on integrated luminosity
2.0 - 2.5	21700.0 ±7.06e+02	stat	+1.35e+03 -1.32e+03	syst ±1.0% syst, uncertainty on branching ratio ±3.7% syst, uncertainty on integrated luminosity
2.5 - 3.0	15400.0 ±4.02e+02	stat	+1.00e+03 -1.17e+03	syst ±1.0% syst, uncertainty on branching ratio ±3.7% syst, uncertainty on integrated luminosity

Sys uncertainty could be more granular?

- If strong correlations... 2 options:
 - a) explicit covariance or correlation matrix **OK only if measurement never intended to be combined with other measurements.**
I didn't see any examples of covariance matrices in the ALICE entries
 - b) give breakdown of signed(!) effect of each NP. **Can then rebuild covariance matrix if each uncertainty is defined as correlated/uncorrelated**
- Prefer to use b) since a) implicitly symmetrizes, and information to correlate with other measurements is insufficient.

Charged-particle production as a function of multiplicity and transverse sphericity
[link]

SQRT(S)	13 TEV		
ETARAP	-0.8 - 0.8		
RE	P P --> CHARGED X		
pT [GEV/C]	d²N/dEtadpT [C/GEV (X ['])]	d²N/dEtadpT [C/GEV (IX ['])]	d²N/dEtadpT [C/GEV (VIII ['])]
0.15 - 0.2	4.446e+00 ±3.247e-03 <i>stat</i> ±3.116e-01 <i>sys,total</i> ±2.916e-01 <i>sys,uncorrelated</i>	8.091e+00 ±5.746e-03 <i>stat</i> ±2.353e-01 <i>sys,total</i> ±1.618e-01 <i>sys,uncorrelated</i>	1.082e+01 ±9.454e-03 <i>stat</i> ±3.251e-01 <i>sys,total</i> ±2.164e-01 <i>sys,uncorrelated</i>

- Statistical correlations as correlation matrix. **Bootstrap Replicas (see backup) best for future combinations but need make TH*DBootstrap code public [Overkill in the ALICE case?]**

Currently still a draft, iterating
with LHCEWWG conveners

Proposed LHC-wide HEPData Recommendations

Public note to be agreed between expts

Defines 3 scenarios for levels of
information to provide on HEPData

Gives concrete recommendations for the
format of objects which are to be stored

Recommendations for preservation on analyses on
HEPData by the LHC experiments

LHC Electoweak Working group

24th October 2019

Contents

1	Introduction	2
2	Current workflow to preserve analysis results	2
3	Limitations of current recommendations for HEPData entries	3
4	Three scenarios for analysis preservation	4
4.1	Scenario A - Maximum Re-interpretability	5
4.2	Scenario B - Approximate Re-interpretability	6
4.3	Scenario C - Minimum Requirements for Analysis Preservation	7
4.4	Results which cannot be re-interpreted	7
5	Specific conventions for preserved objects and examples	8
5.1	Rivet routines	8
5.2	Bootstrap histograms	9
5.3	Statistical correlation matrices	9
5.4	Systematic covariance matrices	10
5.5	Uncertainty breakdowns	10
5.6	Post-fit impacts	10
6	Summary	10
	Appendix	11
A	How to store error breakdowns in HEPData entries	11
A.1	The YODA format, and improvements to store error breakdowns	12
A.2	A common library of tools to manipulate YODA files with covariance information	14
B	Pseudocode examples	14

3 Scenarios for re-interpretation

- Identify different levels of recommendations, depending on the analysis type and how re-interpretable it needs to be:

3.1 Scenario A - Maximum Re-interpretability

3.2 Scenario B - Approximate Re-interpretability

3.3 Scenario C - Minimum Requirements for Analysis Preservation

3.4 Results which cannot be re-interpreted

Best case - aims to provide maximal information for reinterpretations.
Should be gold standard for precision measurements

3 Scenarios for re-interpretation

- Identify different levels of recommendations, depending on the analysis type and how re-interpretable it needs to be:

3.1 Scenario A - Maximum Re-interpretability

3.2 Scenario B - Approximate Re-interpretability

3.3 Scenario C - Minimum Requirements for Analysis Preservation

3.4 Results which cannot be re-interpreted

Best case - aims to provide maximal information for reinterpretations. Should be gold standard for precision measurements

Closest to current situation. Plenty of information published. Not necessarily enough for strict combinations... but good enough for many analyses (especially searches)

3 Scenarios for re-interpretation

- Identify different levels of recommendations, depending on the analysis type and how re-interpretable it needs to be:

3.1 Scenario A - Maximum Re-interpretability

3.2 Scenario B - Approximate Re-interpretability

3.3 Scenario C - Minimum Requirements for Analysis Preservation

3.4 Results which cannot be re-interpreted

Best case - aims to provide maximal information for reinterpretations. Should be gold standard for precision measurements

Closest to current situation. Plenty of information published. Not necessarily enough for strict combinations... but good enough for many analyses (especially searches)

Bare minimum for a search to be re-interpretable

- **Minimum amount of info** for result to be re-used meaningfully.
e.g if only rough estimate of MC/data agreement or sensitivity to new models needed
- **Analysis logic preservation**: Ideally, Rivet routine... if not...
 - detailed description of the region of interest
 - per-object efficiency tables
 - explicit definitions of each variable used in the selection,
 - cutflows of the effect of each selection on well-defined signals
- **Statistical correlations**: omitted if negligible bin migrations.
Stat error per bin still needed (assumed uncorrelated between bins)
- **Systematic correlations**: uncert breakdown or explicit cov matrices
- **Background**: SM bkg prediction of MC generators, w/ breakdown of theory uncertainty if possible [*N/A for Heavy Ions yet?*]

- **For standard measurements or searches** to be re-interpreted approximately. E.g generator tuning , and recasting of searches
- **Analysis logic preservation**: Rivet analysis must be provided **at the same time as the preprint !**
 - If results only at detector level, Rivet analysis should still be provided, with adequate smearing and efficiency tables
- **Statistical correlations**: **correlation matrices**. *Can't infer corrs between analyses, but OK if re-interpreting result in isolation*
- **Systematic correlations**: **uncertainty breakdown**, = effect of each NP on each bin -> cov matrix + correlate w/ other measurements
 - OR, cov matrix for each distribution: e.g. for simplified likelihoods
- **Background**: include **SM prediction from latest MC generators** w/ breakdown of theory uncertainty if possible [*N/A for Heavy Ions yet?*]

- **For precision analyses:** for future combinations, measurements of SM parameters, PDF fitting... Enough info for exact combination
- **Analysis logic preservation:** Rivet analysis must be provided **at the same time as the preprint !**
 - If results only at detector level, Rivet analysis should still be provided, with adequate smearing and efficiency tables
- **Stat correlations:** Bootstrap Replicas attached to HEPData entry
- **Syst correlations:** uncertainty breakdown, = effect of each NP on each bin -> cov matrix + correlate w/ other measurements
- **Background:** include **SM prediction from latest MC generators w/** breakdown of theory uncertainty if possible [*N/A for Heavy Ions yet?*]
- If likelihood fit used: **post-fit values of the NPs in each bin**

- LHC EW WG is reviewing recommendations for HEPData. Recommendations document in preparation...
- Recent developments are excellent opportunity to review status and see what we can do better
 - > **maximise impact** of our measurements
 - > **Agree conventions across experiments!**
- In particular:
 - we should be more diligent about **preserving analysis logic/fiducial volume in e.g. Rivet routines**
 - we should be careful to **give full uncertainty breakdown** instead of just stat vs syst
 - we should **provide SM generator predictions** in the HEPData entry if possible [*N/A for Heavy Ions yet?*]

- LHC EW WG is reviewing recommendations for HERData

Comments and feedback are very very welcome !

Backup

- Bootstrap histograms: like regular histograms, but each time histogram is filled, a certain number (normally ~ 1000) of replicas are also filled, where the filled weight is varied according to a random weight drawn from a Poisson distribution.
- The random number seed is set uniquely by the run and event number of the event in question. \rightarrow statistical uncertainty and correlations can be correctly evaluated when combining results between different analyses, from the same or different collaborations.
- The replicas can be attached to the additional material of the HEPData entry. This has been done for example here:
<https://www.hepdata.net/record/ins1604271>
<https://www.hepdata.net/record/ins1604271>
- Root extension for these classes exists in ATLAS: we are working to make this code public!