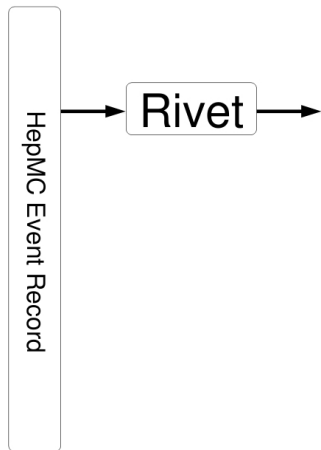


# NEW DEVELOPMENTS IN EVENT GENERATOR TUNING TECHNIQUES

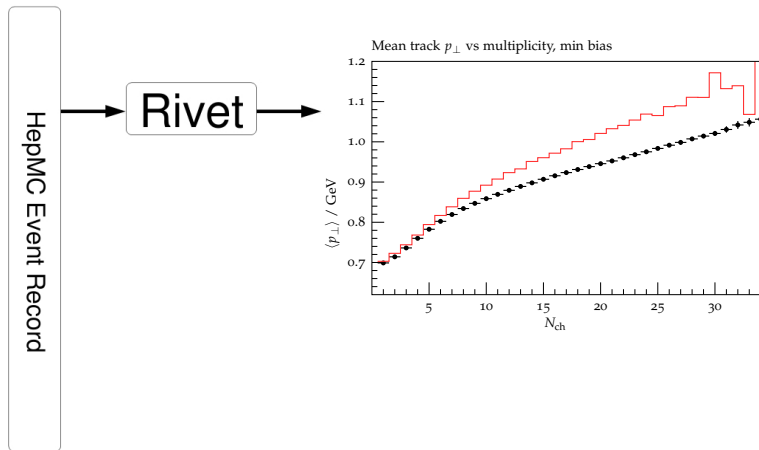
Holger Schulz, Heiko Lacker, Jan Eike von Seggern (HU Berlin),  
Andy Buckley (Edinburgh), Hendrik Hoeth (Durham)

ACAT Jaipur, 22-27 February 2010



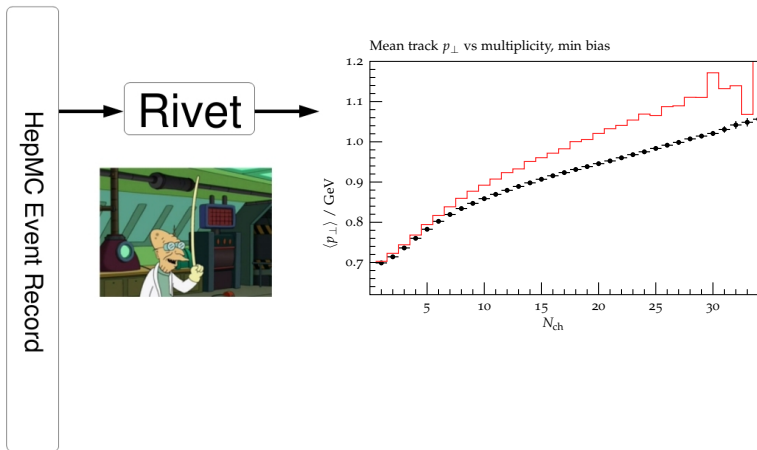


Rivet uses generator independent HepMC events as input



The output can be compared with data, but may be untuned

# INTRODUCTION

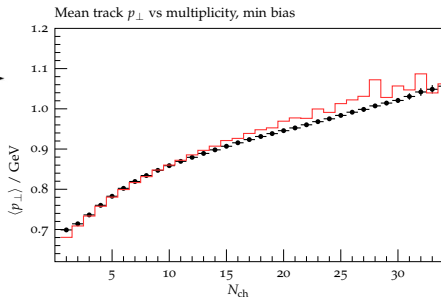


Luckily, there is Professor



HepMC Event Record

Rivet



Professor tunes generators based on Rivet output

- Version 1.2.0 to be released this month
- More than 60 validated analyses from experiments like LEP, TEVATRON, RHIC, SPS, HERA are now implemented
- Rivet is used by ATLAS, CMS, LHCb
- Output histograms are used by Professor for generator tuning
- HepData (<http://durpdg.dur.ac.uk>) now offers AIDA data export

# TOOLS FOR SYSTEMATIC GENERATOR TUNING

DELPHI 1995, Hamacher et al.: bin-wise interpolation of MC generator response and  $\chi^2$  minimization

Professor 2008 ([arXiv:0907.2973](https://arxiv.org/abs/0907.2973), [arXiv:0906.0075](https://arxiv.org/abs/0906.0075), [arXiv:0902.4403](https://arxiv.org/abs/0902.4403)):

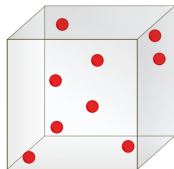
## “PROCEDURE FOR ESTIMATING SYSTEMATIC ERRORS”



- Pick up DELPHI idea
- Use flexible python interface
- Use quadratic **or** cubic interpolations
- Interpolation storing with python.cPickle
- Validation of results possible in many ways
- **NEW:** GUI tool, uncertainty estimates

# TUNING PROCEDURE IN PROFESSOR (1D, 1BIN)

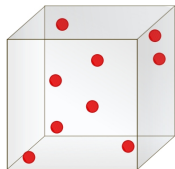
- 1 Random sampling:  $N$  parameter points in  $n$ -dimensional space
- 2 Run generator and fill histograms
- 3 For each bin: use  $N$  points to fit interpolation (2<sup>nd</sup> or 3<sup>rd</sup> order polynomial)
- 4 Construct overall (now trivial)  $\chi^2 = \sum_{bins} \frac{(interpolation - data)^2}{error^2}$
- 5 Numerically minimize using `pyMinuit`, `SciPy`





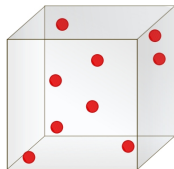
# TUNING PROCEDURE IN PROFESSOR (1D, 1BIN)

- 1 Random sampling:  $N$  parameter points in  $n$ -dimensional space
- 2 Run generator and fill histograms
- 3 For each bin: use  $N$  points to fit interpolation (2<sup>nd</sup> or 3<sup>rd</sup> order polynomial)
- 4 Construct overall (now trivial)  $\chi^2 = \sum_{bins} \frac{(interpolation - data)^2}{error^2}$
- 5 Numerically minimize using `pyMinuit`, `SciPy`



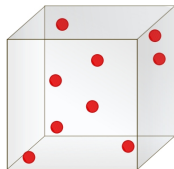
# TUNING PROCEDURE IN PROFESSOR (1D, 1BIN)

- 1 Random sampling:  $N$  parameter points in  $n$ -dimensional space
- 2 Run generator and fill histograms
- 3 For each bin: use  $N$  points to fit interpolation (2<sup>nd</sup> or 3<sup>rd</sup> order polynomial)
- 4 Construct overall (now trivial)  $\chi^2 = \sum_{bins} \frac{(interpolation - data)^2}{error^2}$
- 5 Numerically minimize using `pyMinuit`, `SciPy`



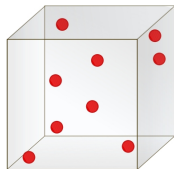
# TUNING PROCEDURE IN PROFESSOR (1D, 1BIN)

- 1 Random sampling:  $N$  parameter points in  $n$ -dimensional space
- 2 Run generator and fill histograms
- 3 For each bin: use  $N$  points to fit interpolation (2<sup>nd</sup> or 3<sup>rd</sup> order polynomial)
- 4 Construct overall (now trivial)  $\chi^2 = \sum_{bins} \frac{(interpolation - data)^2}{error^2}$
- 5 Numerically minimize using `pyMinuit`, `SciPy`



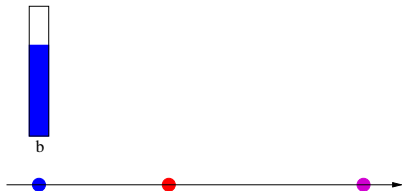
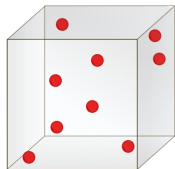
# TUNING PROCEDURE IN PROFESSOR (1D, 1BIN)

- 1 Random sampling:  $N$  parameter points in  $n$ -dimensional space
- 2 Run generator and fill histograms
- 3 For each bin: use  $N$  points to fit interpolation (2<sup>nd</sup> or 3<sup>rd</sup> order polynomial)
- 4 Construct overall (now trivial)  $\chi^2 = \sum_{bins} \frac{(interpolation - data)^2}{error^2}$
- 5 Numerically minimize using `pyMinuit`, `SciPy`



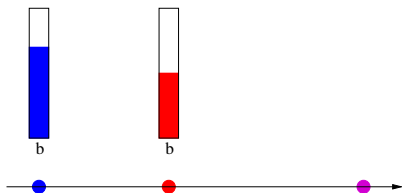
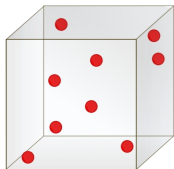
# TUNING PROCEDURE IN PROFESSOR (1D, 1BIN)

- 1 Random sampling:  $N$  parameter points in  $n$ -dimensional space
- 2 Run generator and fill histograms
- 3 For each bin: use  $N$  points to fit interpolation (2<sup>nd</sup> or 3<sup>rd</sup> order polynomial)
- 4 Construct overall (now trivial)  $\chi^2 = \sum_{bins} \frac{(interpolation - data)^2}{error^2}$
- 5 Numerically minimize using `pyMinuit`, `SciPy`



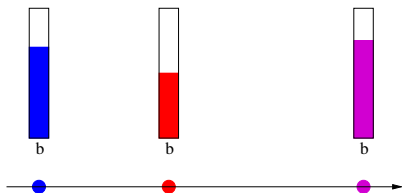
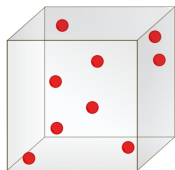
# TUNING PROCEDURE IN PROFESSOR (1D, 1BIN)

- 1 Random sampling:  $N$  parameter points in  $n$ -dimensional space
- 2 Run generator and fill histograms
- 3 For each bin: use  $N$  points to fit interpolation (2<sup>nd</sup> or 3<sup>rd</sup> order polynomial)
- 4 Construct overall (now trivial)  $\chi^2 = \sum_{bins} \frac{(interpolation - data)^2}{error^2}$
- 5 Numerically minimize using `pyMinuit`, `SciPy`



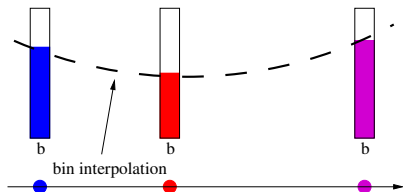
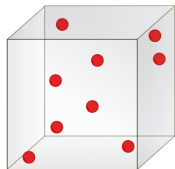
# TUNING PROCEDURE IN PROFESSOR (1D, 1BIN)

- 1 Random sampling:  $N$  parameter points in  $n$ -dimensional space
- 2 Run generator and fill histograms
- 3 For each bin: use  $N$  points to fit interpolation (2<sup>nd</sup> or 3<sup>rd</sup> order polynomial)
- 4 Construct overall (now trivial)  $\chi^2 = \sum_{bins} \frac{(interpolation - data)^2}{error^2}$
- 5 Numerically minimize using pyMinuit, SciPy



# TUNING PROCEDURE IN PROFESSOR (1D, 1Bin)

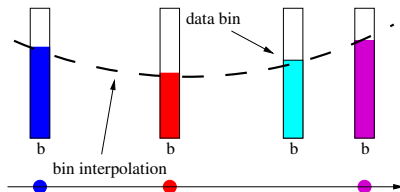
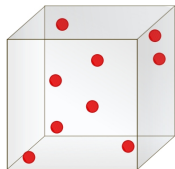
- 1 Random sampling:  $N$  parameter points in  $n$ -dimensional space
- 2 Run generator and fill histograms
- 3 For each bin: use  $N$  points to fit interpolation (2<sup>nd</sup> or 3<sup>rd</sup> order polynomial)
- 4 Construct overall (now trivial)  $\chi^2 = \sum_{bins} \frac{(interpolation - data)^2}{error^2}$
- 5 Numerically minimize using pyMinuit, SciPy





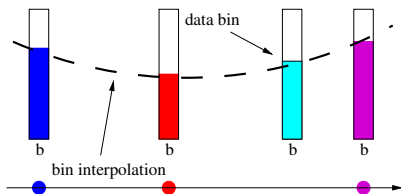
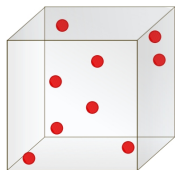
# TUNING PROCEDURE IN PROFESSOR (1D, 1Bin)

- 1 Random sampling:  $N$  parameter points in  $n$ -dimensional space
- 2 Run generator and fill histograms
- 3 For each bin: use  $N$  points to fit interpolation (2<sup>nd</sup> or 3<sup>rd</sup> order polynomial)
- 4 Construct overall (now trivial)  $\chi^2 = \sum_{bins} \frac{(interpolation - data)^2}{error^2}$
- 5 Numerically minimize using pyMinuit, SciPy



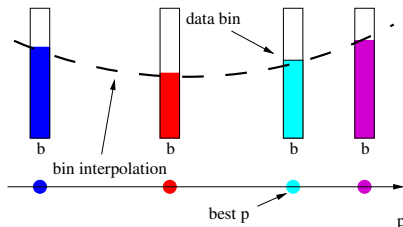
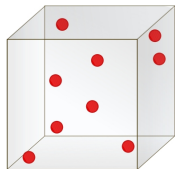
# TUNING PROCEDURE IN PROFESSOR (1D, 1Bin)

- 1 Random sampling:  $N$  parameter points in  $n$ -dimensional space
- 2 Run generator and fill histograms
- 3 For each bin: use  $N$  points to fit interpolation (2<sup>nd</sup> or 3<sup>rd</sup> order polynomial)
- 4 Construct overall (now trivial)  $\chi^2 = \sum_{bins} \frac{(interpolation - data)^2}{error^2}$
- 5 Numerically **minimize** using pyMinuit, SciPy



# TUNING PROCEDURE IN PROFESSOR (1D, 1Bin)

- 1 Random sampling:  $N$  parameter points in  $n$ -dimensional space
- 2 Run generator and fill histograms
- 3 For each bin: use  $N$  points to fit interpolation (2<sup>nd</sup> or 3<sup>rd</sup> order polynomial)
- 4 Construct overall (now trivial)  $\chi^2 = \sum_{bins} \frac{(interpolation - data)^2}{error^2}$
- 5 Numerically **minimize** using pyMinuit, SciPy



# OVERSAMPLING AND INTERPOLATION QUALITY

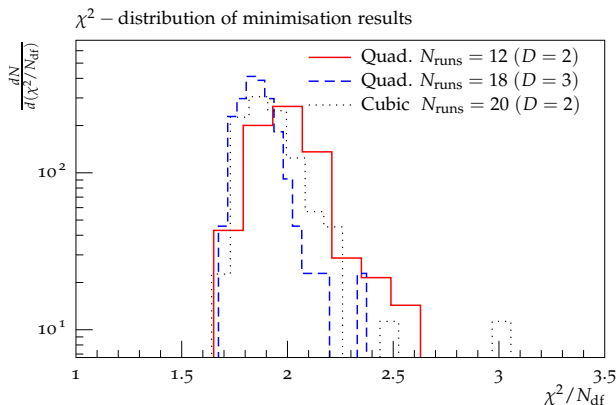
- In  $n$  dimensions: Singular Value Decomposition (SVD) requires  $N_{\min}^{(n)}$  generator runs:

$$N_{\min}^{(n)} = 1 + n + n(n+1)/2 + \underbrace{(n+1)(n+2)/6}_{\text{cubic only}}$$

- SVD allows for oversampling.
- Degree of oversampling:  $D = N_{\text{runs}}/N_{\min}^{(n)}$
- What is a sensible  $D$ ?
- $\rightarrow$  use  $\mathcal{O}(1000)$  different interpolations with different  $N_{\text{runs}}$
- Perform minimisations, investigate g.o.f. measures
- Examples shown are from a two-dimensional Tuning of Jimmy

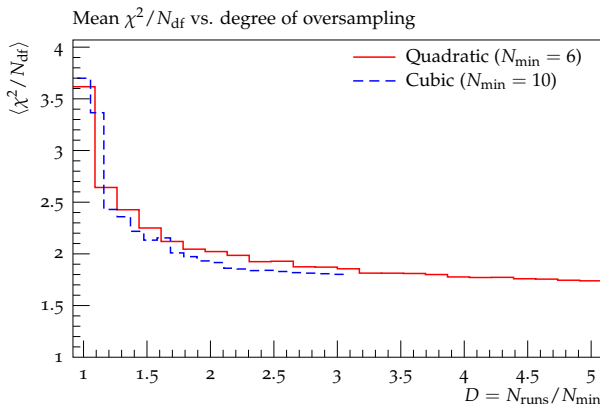
# DISTRIBUTION OF GOODNESS OF FIT VALUES

- Spread of results decreases with increasing  $D$ , polynomial degree
- Observe lower  $\chi^2 / N_{\text{df}}$ -boundary



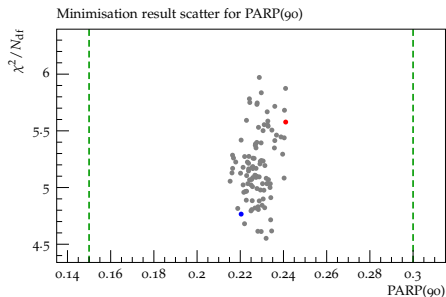
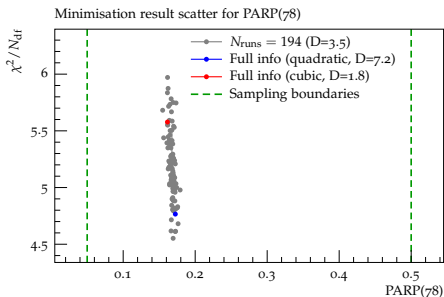
# GOODNESS OF FIT VS. DEGREE OF OVERSAMPLING

- Oversampling is necessary,  $D > 2 \dots 3$  seems sensible
- However, g.o.f. improves slowly for  $D > 4$ , almost saturates



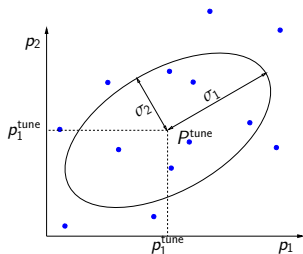
# COMBINATION TUNING UNCERTAINTY

- “Kebap plots”: investigate minimisation results parameterwise
- Used as rough estimate on how well a parameter is constrained
- Goal: translate spread into tuning uncertainty, *combination* error



# STATISTICAL TUNING UNCERTAINTY

- Exploit covariance matrix returned by minimiser
- Sample points from the valley of the minimum
- Translate these points into *statistical* tuning uncertainty

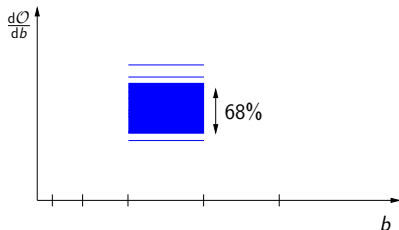
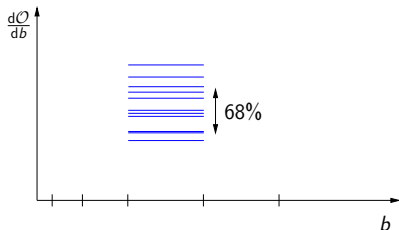


**Ultimately:** combine *statistical* and *combination* error to a combined error



# CONFIDENCE BELT CONSTRUCTION

- 1 Use points sampled from ellipsis or different minimisation results
- 2 Use parameterisation to get bin-content predictions
- 3 For each bin  $b$  and each observable  $\mathcal{O}$ : determine central 68, 95 pct.



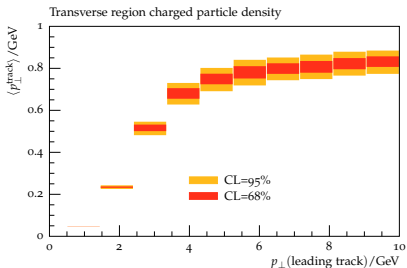
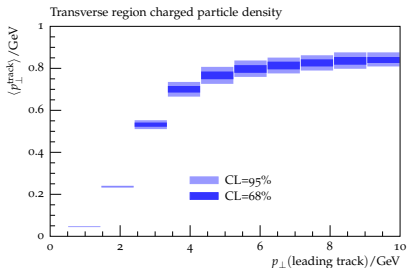
In the following: study behaviour of tuning uncertainty estimates when (pseudo-) data is added

# CONFIDENCE BELT - WITHOUT PSEUDODATA

These plots show the tuning uncertainties of Jimmy if tuned with Professor to Tevatron data for an observable at 7 TeV

statistical uncertainties

combination uncertainties

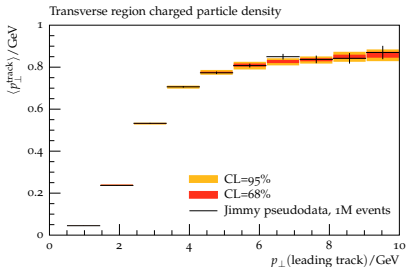
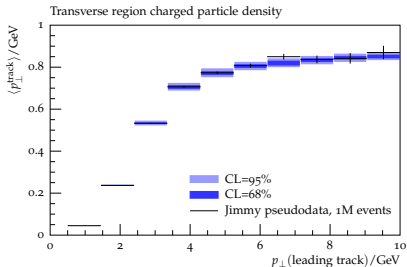


# CONFIDENCE BELT - ADDING PSEUDODATA

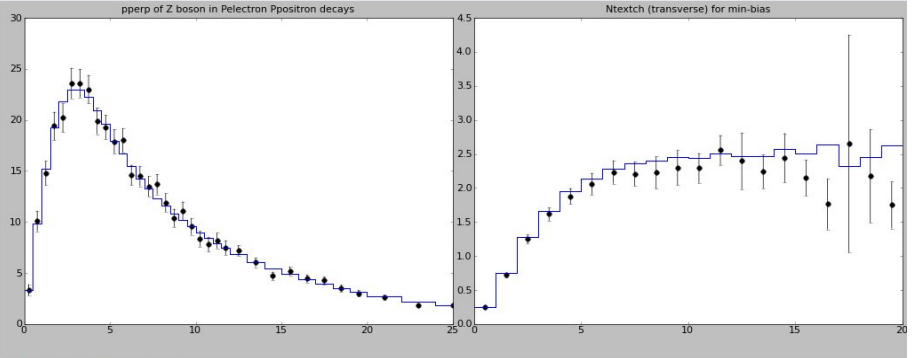
And this is what happens if we add 1 million events of 7 TeV pseudo-data to the Tevatron data in the tuning

statistical uncertainties

combination uncertainties



- Event generation usually very time consuming
- Hard to get a feeling how parameters influence observables
- Luckily, we have the generator parameterisation handy
- Prediction for bin content can be approximately calculated
- Why not put this into a nice GUI (wxPython), add a vertical slider for each parameter and play a bit?
- → real-time response to shifts in parameter space

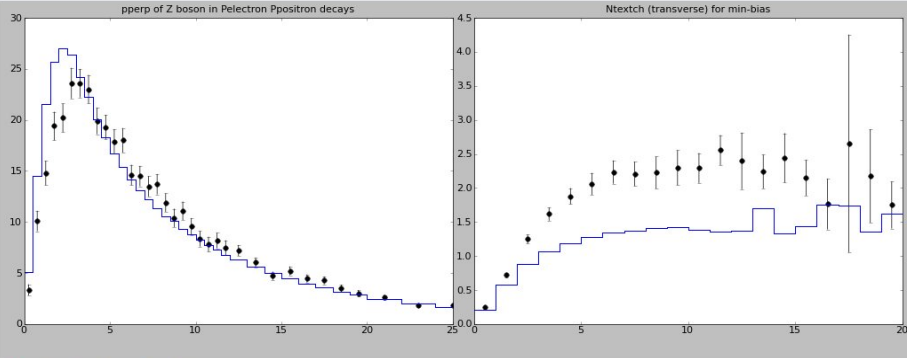


Obs 1: /CDF\_2000\_S4155203/d01-x01-y01  logy Obs 2: /CDF\_2001\_S4751469/d03-x01-y02  logy  Show g.o.f.

PARP(64) scale of alpha_S	<input type="text"/>	1.2981	Limits 1:	Limits 2:
PARP(71) ME-FSR shower merging	<input type="text"/>	1.9987	XMin	<input type="text" value="None"/>
PARP(78) Colour reconnection	<input type="text"/>	0.1699	XMax	<input type="text" value="25"/>
PARP(79) Beam remnant x-enhancement	<input type="text"/>	1.1772	YMin	<input type="text" value="None"/>
PARP(82) UE cut-off	<input type="text"/>	1.8499	YMax	<input type="text" value="None"/>
PARP(83) Hadronic matter distribution	<input type="text"/>	1.7992		
PARP(90) UE energy evolution exponent	<input type="text"/>	0.2199		
PARP(91) Width of primordial kt	<input type="text"/>	1.9997		
PARP(93) Upper cut-off of primordial kt	<input type="text"/>	6.9984		

Buttons: Set params, Reset limits 1, Reset limits 2

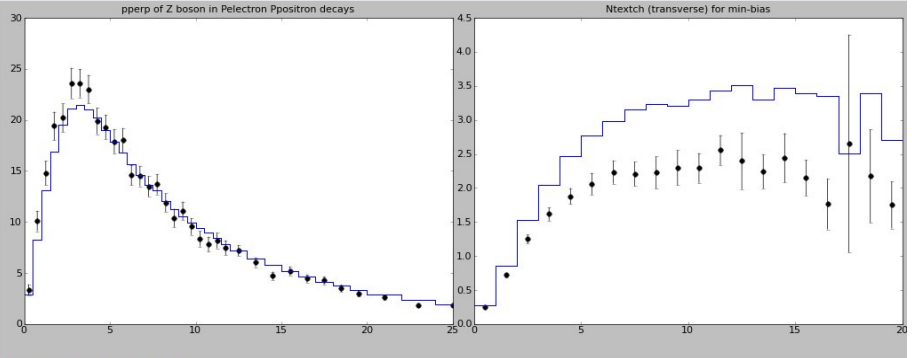




Obs 1: /CDF\_2000\_S4155203/d01-x01-y01  logy Obs 2: /CDF\_2001\_S4751469/d03-x01-y02  logy  Show g.o.f.

PARP(64) scale of alpha_S	<input type="text"/>	3.7611	Set params	Limits 1:	Limits 2:		
PARP(71) ME-FSR shower merging	<input type="text"/>	3.3664		XMin	None	XMin	None
PARP(78) Colour reconnection	<input type="text"/>	0.1699	Reset limits 1	XMax	25	XMax	20
PARP(79) Beam remnant x-enhancement	<input type="text"/>	1.1772		YMin	None	YMin	None
PARP(82) UE cut-off	<input type="text"/>	2.3230	Reset limits 2	YMax	None	YMax	None
PARP(83) Hadronic matter distribution	<input type="text"/>	1.4979					
PARP(90) UE energy evolution exponent	<input type="text"/>	0.2199					
PARP(91) Width of primordial kt	<input type="text"/>	1.6985					
PARP(93) Upper cut-off of primordial kt	<input type="text"/>	9.3402					





Obs 1: /CDF\_2000\_S4155203/d01-x01-y01  logy Obs 2: /CDF\_2001\_S4751469/d03-x01-y02  logy  Show g.o.f.

PARP(64) scale of alpha_S	<input type="text"/>	0.6665	Limits 1:	Limits 2:
PARP(71) ME-FSR shower merging	<input type="text"/>	1.3692		
PARP(78) Colour reconnection	<input type="text"/>	0.1235	XMax: <input type="text" value="25"/>	XMax: <input type="text" value="20"/>
PARP(79) Beam remnant x-enhancement	<input type="text"/>	2.7371	YMin: <input type="text" value="None"/>	YMin: <input type="text" value="None"/>
PARP(82) UE cut-off	<input type="text"/>	1.6301	YMax: <input type="text" value="None"/>	YMax: <input type="text" value="None"/>
PARP(83) Hadronic matter distribution	<input type="text"/>	1.9291		
PARP(90) UE energy evolution exponent	<input type="text"/>	0.2995		
PARP(91) Width of primordial kt	<input type="text"/>	2.2521		
PARP(93) Upper cut-off of primordial kt	<input type="text"/>	2.8783		

Buttons: Set params, Reset limits 1, Reset limits 2



- We studied how the interpolation benefits from oversampling
- $N_{\text{runs}} / N_{\text{min}}^{(n)} > 2 \dots 3$  is advisable
- Working on quantification of tuning uncertainties (to be submitted to Les Houches proceedings)
- We have designed a GUI tool to explore behaviour of observables interactively

---

Thank you!

---





Backup

2nd order polynomial includes lowest-order correlations between parameters

$$MC_b(\vec{p}) \approx f^{(b)}(\vec{p}) = \alpha_0^{(b)} + \sum_i \beta_i^{(b)} p_i + \sum_{i \leq j} \gamma_{ij}^{(b)} p_i p_j$$

Now use N generator runs, i.e. N different parameter sets x,y:

$$\underbrace{\begin{pmatrix} v_1 \\ v_2 \\ \vdots \\ v_N \end{pmatrix}}_{\vec{v} \text{ (N values, i.e. N bin contents)}} = \underbrace{\begin{pmatrix} 1 & x_1 & y_1 & x_1^2 & x_1 y_1 & y_1^2 \\ 1 & x_2 & y_2 & x_2^2 & x_2 y_2 & y_2^2 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & x_N & y_N & x_N^2 & x_N y_N & y_N^2 \end{pmatrix}}_{\tilde{\mathbf{P}} \text{ (N sampled parameter sets)}} \underbrace{\begin{pmatrix} \alpha_0 \\ \beta_x \\ \beta_y \\ \gamma_{xx} \\ \gamma_{xy} \\ \gamma_{yy} \end{pmatrix}}_{\vec{c} \text{ (coeffs)}}$$

Therefore:  $\vec{c}_b = \tilde{\mathcal{I}}[\tilde{\mathbf{P}}]\vec{v}$  where  $\tilde{\mathcal{I}}$  is the pseudoinverse operator.

$$\vec{c}_b = \tilde{\mathcal{I}}[\tilde{\mathbf{P}}]\vec{v}$$

- Use Singular Value Decomposition (SVD), a general diagonalisation for all normal matrices  $M: M = U\Sigma V^*$
- Method available in SciPy.linalg
- Minimal number of runs = number of coefficients in  $\vec{c}_b$ :

$$N_{\min}^{(n)} = 1 + n + n(n+1)/2 + \underbrace{(n+1)(n+2)/6}_{\text{cubic only}}$$

- Oversampling by a factor of three has proven to be much better

Num params, $P$	$N_2^{(P)}$ (2nd order)	$N_3^{(P)}$ (3rd order)
1	3	4
2	6	10
4	15	35
6	28	84
8	45	165
9	55	220

$$\vec{c}_b = \tilde{\mathcal{I}}[\tilde{\mathbf{P}}]\vec{v}$$

- Use Singular Value Decomposition (SVD), a general diagonalisation for all normal matrices  $M: M = U\Sigma V^*$
- Method available in SciPy.linalg
- Minimal number of runs = number of coefficients in  $\vec{c}_b$ :

$$N_{\min}^{(n)} = 1 + n + n(n+1)/2 + \underbrace{(n+1)(n+2)/6}_{\text{cubic only}}$$

- Oversampling by a factor of three has proven to be much better

Num params, $P$	$N_2^{(P)}$ (2nd order)	$N_3^{(P)}$ (3rd order)
1	3	4
2	6	10
4	15	35
6	28	84
8	45	165
9	55	220

$$\vec{c}_b = \tilde{\mathcal{I}}[\tilde{\mathbf{P}}]\vec{v}$$

- Use Singular Value Decomposition (SVD), a general diagonalisation for all normal matrices  $M: M = U\Sigma V^*$
- Method available in SciPy.linalg
- Minimal number of runs = number of coefficients in  $\vec{c}_b$ :  
$$N_{\min}^{(n)} = 1 + n + n(n+1)/2 + \underbrace{(n+1)(n+2)/6}_{\text{cubic only}}$$
- Oversampling by a factor of three has proven to be much better

Num params, $P$	$N_2^{(P)}$ (2nd order)	$N_3^{(P)}$ (3rd order)
1	3	4
2	6	10
4	15	35
6	28	84
8	45	165
9	55	220