

# Is there more to learn from the high-x data ?

Structure Function

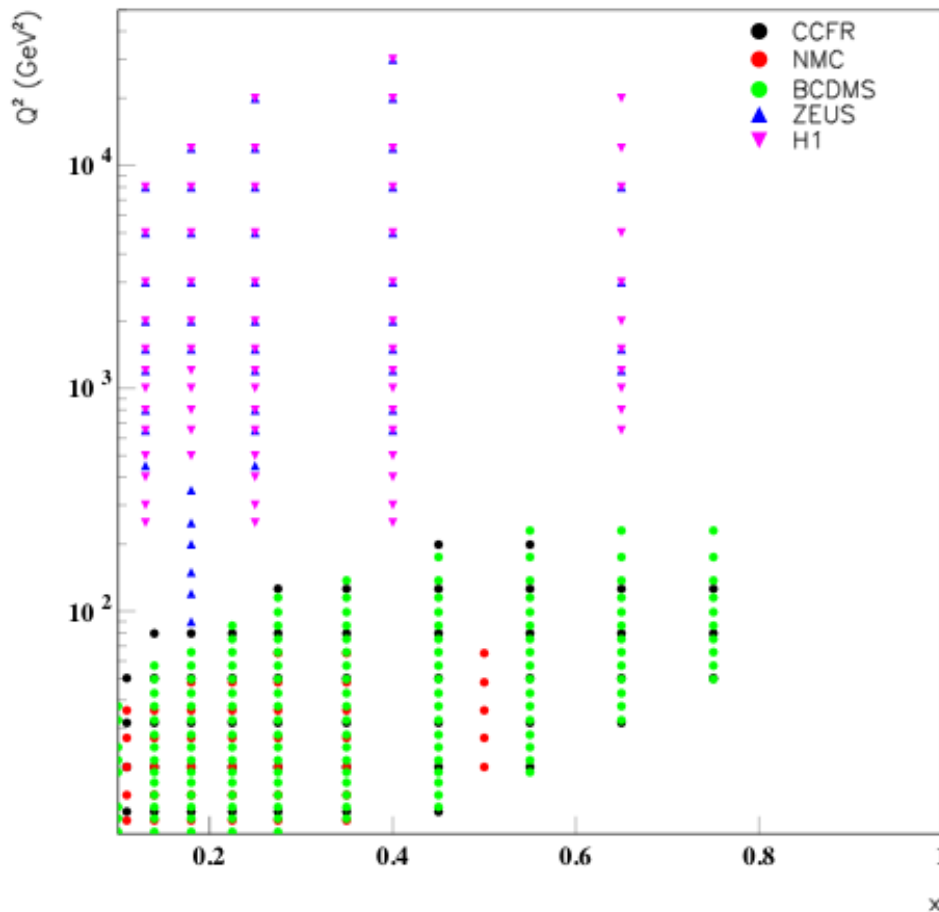
A. Caldwell

Max Planck Institute for Physics



# Motivation

Our information on the very high  $x$  behavior of the parton densities is primarily theoretical.



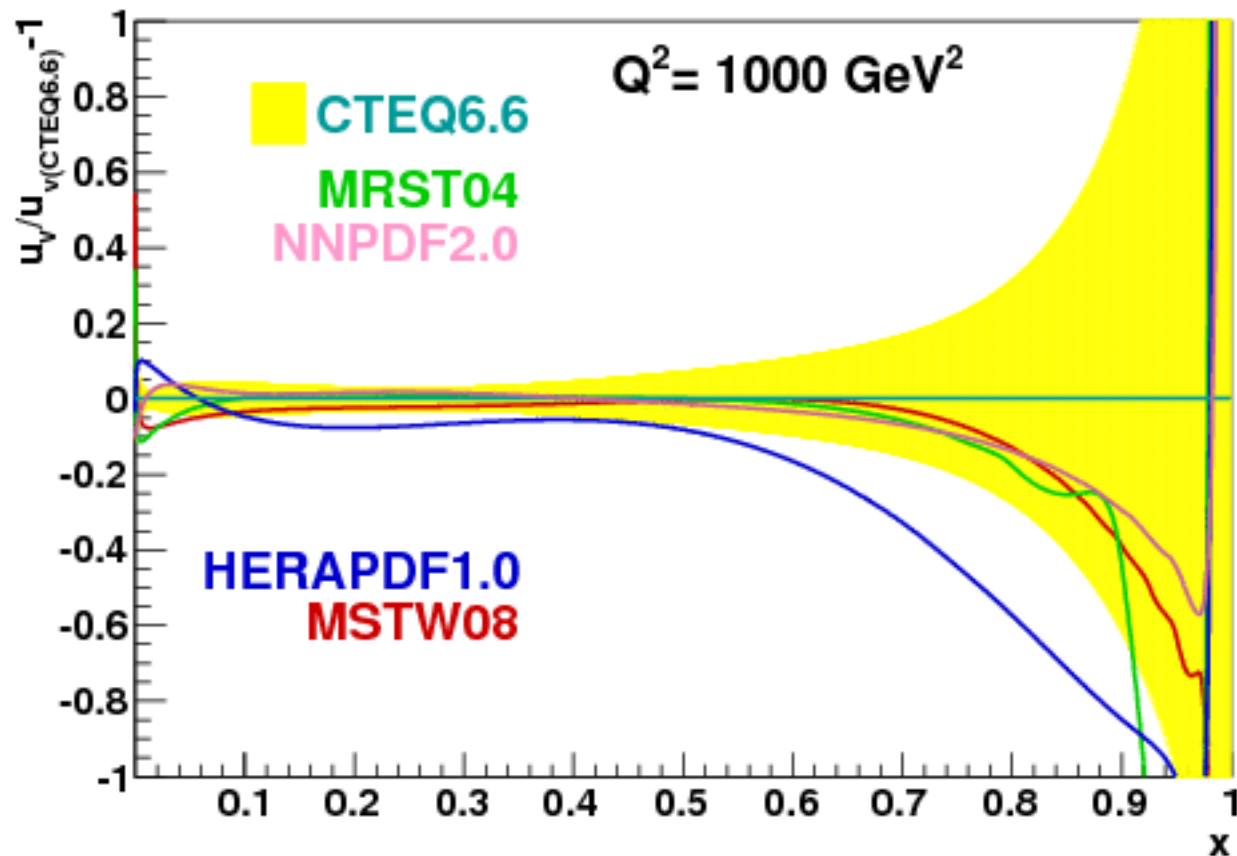
BCDMS has measured  $F_2$  up to  $x=0.75$

H1, ZEUS have measured  $F_2$  up to  $x=0.65$

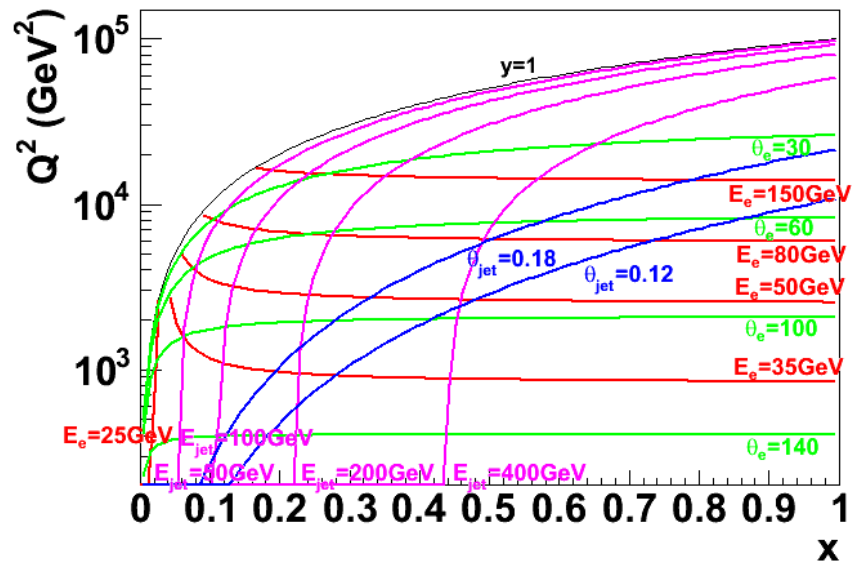
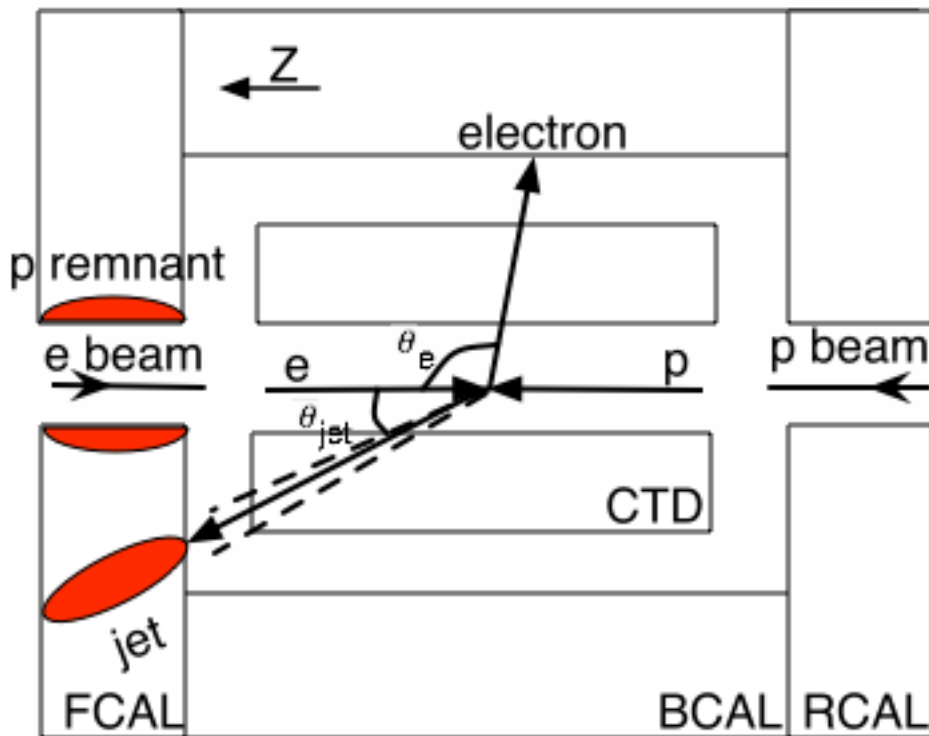
(not including dedicated ZEUS analysis at high- $x$ )

# Motivation

The PDF's are poorly determined at high- $x$ . Sizeable differences despite the fact that fits use similar parametrization  $xq \propto (1-x)^\eta$ . Is it possible to improve this situation ?



# ZEUS high-x analysis



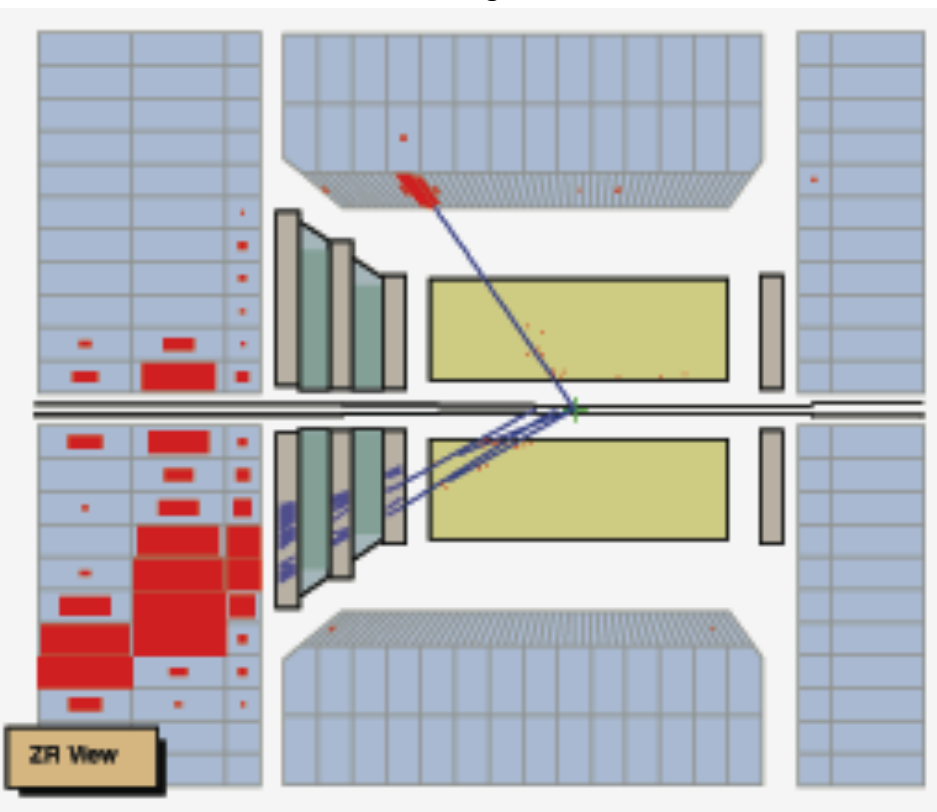
- For not too high x, measure x from jet:  $\frac{d^2\sigma}{dx dQ^2}$

- For  $x > x_{\text{Edge}}$ , measure  $\int_{x_{\text{Edge}}}^1 \frac{d^2\sigma}{dx dQ^2} dx$

# HERA kinematics

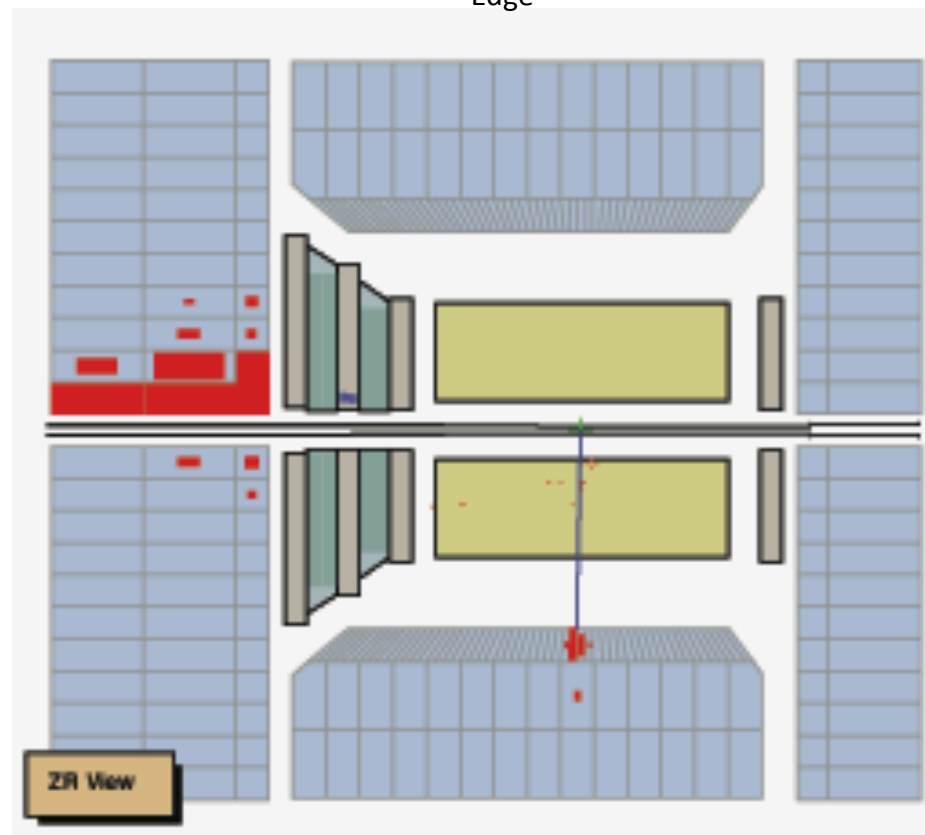
Jet found

$$x < x_{\text{Edge}}$$



No jet found

$$x > x_{\text{Edge}}$$



Jet definition:  $E_T > 10 \text{ GeV}$ ,  $\theta_{\text{jet}} > 0.12$

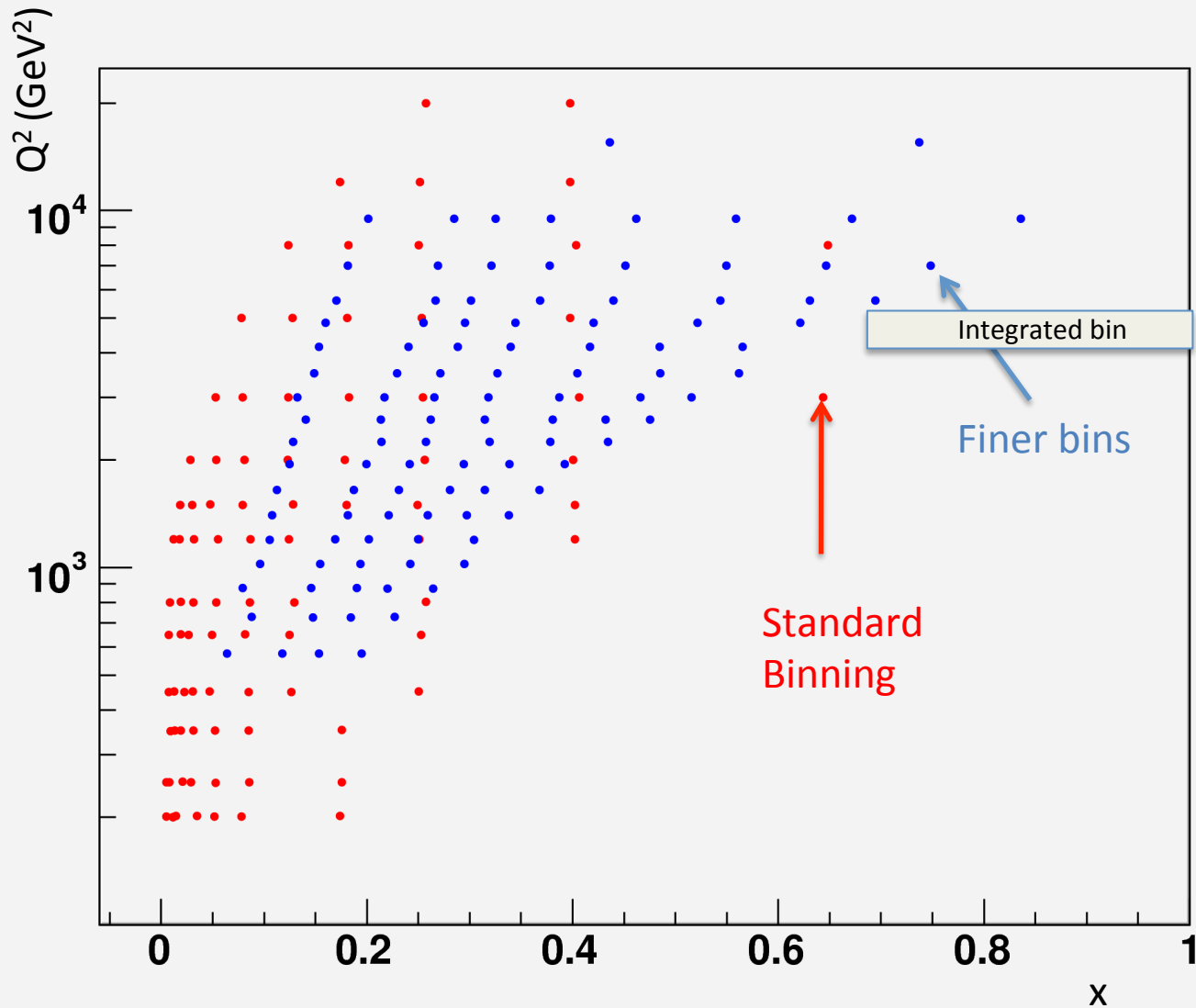
November 12, 2014

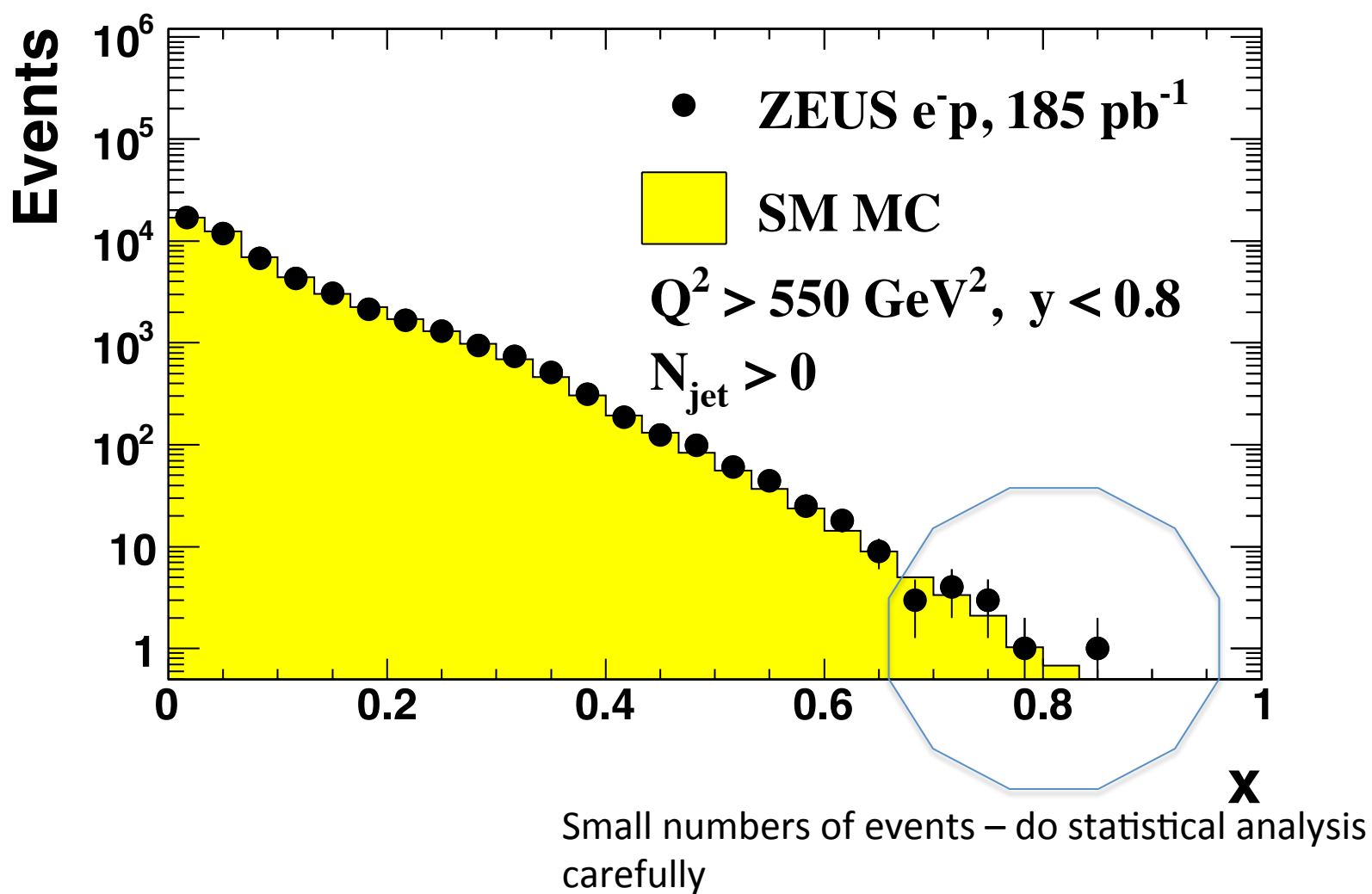
only 0,1 jet events used

Allen Caldwell

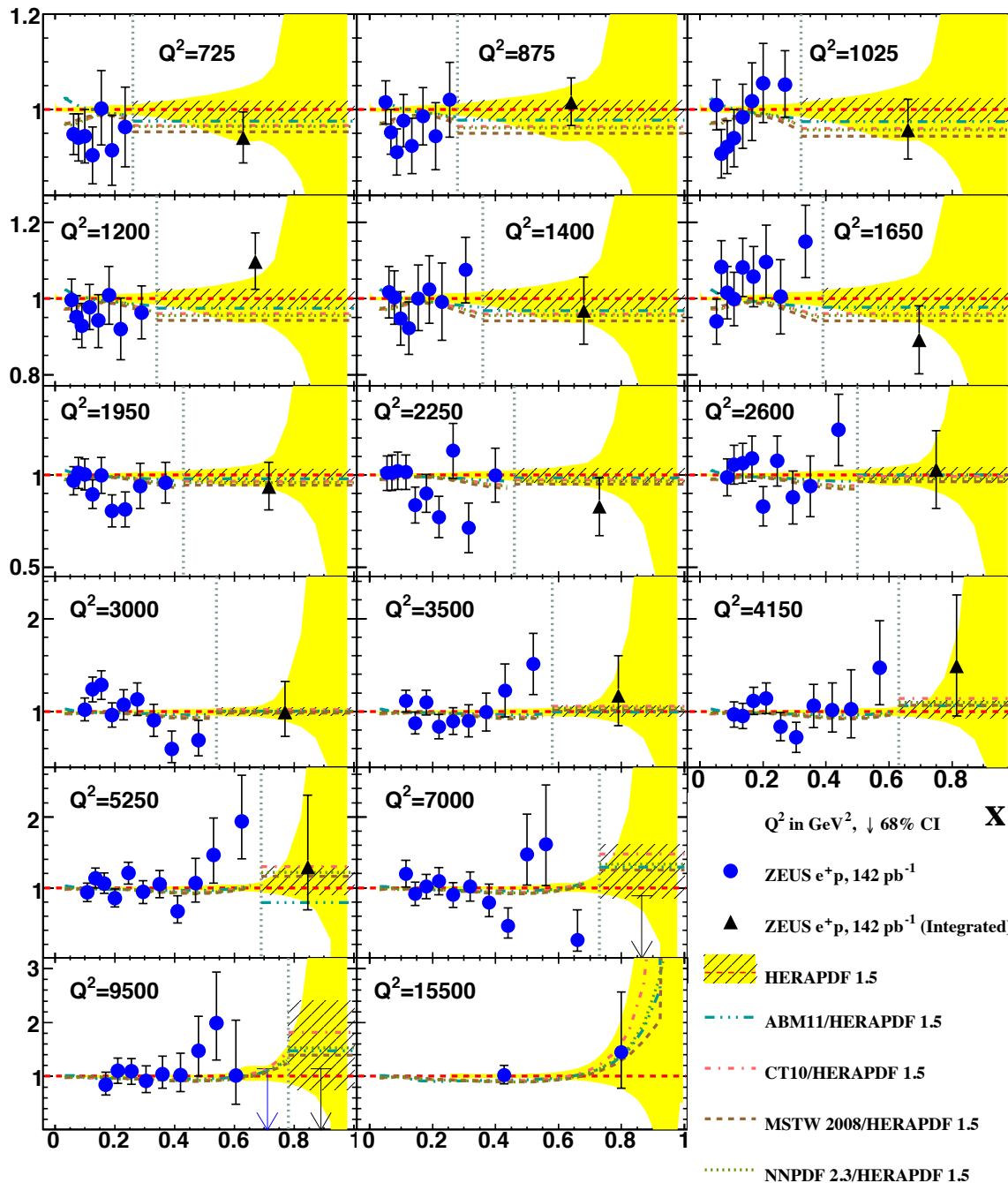
HERA Workshop

# Fine-grained cross section measurements





## ZEUS



Error bars indicate the range of probable values for the underlying cross section given the measured data. How to use this information in a fit ?

Use the observed number of events & calculate the probability to see this number given the model expectation.

5250	0.62	5	$1.76e-04$	+55.2 -35.2	+11.1 -10.5	+9.2 -9.1	-3.2 +4.6	+0.1 +0.1	+3.7 -3.7	+0.6 -0.6
7000	0.12	93	$1.61e-02$	+10.4 -10.4	+4.0 -5.1	+3.3 -3.6	-1.1 +0.8	+0.8 -0.5	-0.7 +0.7	+0.0 -0.0
7000	0.14	89	$1.25e-02$	+10.6 -10.6	+3.7 -5.2	+3.4 -3.5	-1.3 +1.2	-0.5 +0.2	-0.9 +0.9	+0.0 -0.0
7000	0.18	68	$7.02e-03$	+12.1 -12.1	+3.9 -3.6	+3.4 -3.4	-0.6 +0.6	-0.6 +0.4	-0.4 +0.4	+0.0 -0.0
7000	0.22	56	$5.60e-03$	+13.4 -13.4	+4.2 -4.2	+3.9 -3.9	-1.4 +1.1	-0.4 +1.0	-0.4 +0.4	+0.0 -0.0
7000	0.26	49	$3.79e-03$	+14.3 -14.3	+4.6 -4.8	+3.9 -4.0	-0.2 +2.1	+0.2 -0.2	-0.5 +0.5	+0.0 -0.0
7000	0.32	41	$2.70e-03$	+15.6 -15.6	+5.1 -4.7	+5.3 -4.5	-1.4 +0.8	-0.4 +0.4	-0.2 +0.2	+0.0 -0.0
7000	0.38	23	$1.52e-03$	+20.9 -20.9	+6.4 -6.2	+5.5 -5.5	-1.8 +1.7	-0.7 +0.4	+2.0 -2.0	+0.0 -0.0
7000	0.44	17	$1.15e-03$	+27.2 -21.3	+8.4 -7.9	+7.1 -7.1	-2.7 +2.7	-0.0 +0.2	-2.4 +2.4	+0.0 -0.0
7000	0.50	8	$5.38e-04$	+41.8 -29.4	+9.7 -10.3	+9.5 -9.5	-1.6 +1.4	-0.3 +0.4	+2.4 -2.4	+0.1 -0.1
7000	0.56	4	$2.37e-04$	+63.2 -38.2	+12.3 -11.8	+11.3 -11.3	-3.4 +3.4	+0.1 -0.0	+1.2 -1.2	+0.2 -0.2
7000	0.66	10	$2.30e-04$	+36.7 -26.8	+12.6 -13.6	+12.1 -12.3	-4.3 +2.5	-0.3 -0.0	+2.0 -2.0	+0.9 -0.9
9500	0.17	76	$6.77e-03$	+11.5 -11.5	+5.6 -7.7	+4.9 -4.9	-2.0 +2.3	+0.2 -0.2	-0.6 +0.6	+0.0 -0.0
9500	0.21	53	$3.87e-03$	+13.7 -13.7	+5.8 -5.1	+4.3 -4.5	-1.1 +1.8	-0.7 +0.4	-1.1 +1.1	+0.0 -0.0
9500	0.25	40	$2.27e-03$	+15.8 -15.8	+4.8 -4.9	+4.5 -4.5	-2.0 +1.5	+0.2 +0.4	+0.1 -0.1	+0.0 -0.0
9500	0.31	27	$1.50e-03$	+19.2 -19.2	+5.7 -8.1	+5.2 -5.3	-2.6 +1.5	-0.5 +0.2	+1.5 -1.5	+0.0 -0.0
9500	0.36	19	$8.89e-04$	+25.5 -20.3	+6.6 -6.1	+5.9 -5.9	-1.0 +1.9	-0.4 +0.2	-0.3 +0.3	+0.0 -0.0
9500	0.42	12	$5.64e-04$	+33.1 -24.8	+11.3 -7.5	+13.4 -7.3	-1.0 +2.4	-0.8 +0.5	-0.8 +0.8	+0.0 -0.0
9500	0.48	8	$3.63e-04$	+41.7 -29.4	+10.5 -10.4	+9.2 -9.2	-2.6 +2.4	-0.4 +0.6	-3.4 +3.4	+0.0 -0.0
9500	0.54	5	$2.31e-04$	+55.2 -35.2	+14.3 -13.7	+12.5 -12.2	-1.7 +3.6	-0.2 +0.6	+5.7 -5.7	+0.1 -0.1
9500	0.61	4	$1.39e-04$	+63.3 -38.3	+15.5 -15.4	+14.6 -14.8	-4.2 +4.2	+0.0 -0.1	+0.4 -0.4	+0.4 -0.4
9500	0.71	1	$1.50e-05$	+158.0 -58.0	+21.1 -19.8	+18.9 -18.9	-3.3 +4.5	-0.4 +0.3	+4.8 -4.8	+1.3 -1.3

This uncertainty refers to how well we know the underlying cross section assuming that our only knowledge is the observed number of events. Not the uncertainty that belongs in a fit.

# Integrated bins

$Q^2$ (GeV <sup>2</sup> )	$x_{\text{edge}}$	$N$	$I(x)$ (pb/GeV <sup>2</sup> )	$\delta_{\text{stat}}$ (%)	$\delta_{\text{sys}}$ (%)	$\delta_u$ (%)	$\delta_1$ (%)	$\delta_2$ (%)	$\delta_3$ (%)	$\delta_4$ (%)
725	0.63	504	$7.71e-02$	+4.5 -4.5	+2.8 -3.2	+1.5 -1.3	+1.3 -2.0	+0.0 -0.1	+1.9 -1.9	+0.3 -0.3
875	0.64	671	$5.12e-02$	+3.9 -3.9	+2.3 -1.9	+1.2 -1.2	-0.1 +1.0	+0.0 -0.0	+1.3 -1.3	+0.5 -0.5
1025	0.66	414	$2.75e-02$	+4.9 -4.9	+3.4 -3.6	+1.5 -1.5	-1.6 +0.8	+0.0 -0.0	+2.7 -2.7	+0.6 -0.6
1200	0.67	368	$1.80e-02$	+5.2 -5.2	+3.7 -2.9	+1.7 -1.6	-1.5 +2.5	+0.0 -0.0	+1.4 -1.4	+0.8 -0.8
1400	0.68	202	$1.04e-02$	+7.0 -7.0	+3.5 -3.8	+2.1 -2.0	-1.8 +1.3	+0.0 -0.0	+2.1 -2.1	+1.0 -1.0
1650	0.69	173	$5.91e-03$	+7.6 -7.6	+4.4 -4.1	+2.3 -2.2	-1.7 +2.0	+0.1 -0.1	+2.6 -2.6	+1.2 -1.2
1950	0.71	74	$2.51e-03$	+11.6 -11.6	+5.2 -5.1	+3.1 -3.0	-1.9 +1.9	+0.0 -0.1	+3.1 -3.1	+1.6 -1.6
2250	0.73	51	$1.84e-03$	+14.0 -14.0	+7.1 -7.6	+4.1 -4.1	-3.0 +2.2	+0.0 -0.0	+4.9 -4.9	+2.0 -2.0
2600	0.75	36	$9.65e-04$	+16.7 -16.7	+6.9 -6.6	+4.8 -4.8	-1.9 +2.6	+0.0 -0.0	+3.0 -3.0	+2.5 -2.5
3000	0.77	19	$4.90e-04$	+25.5 -20.3	+11.0 -10.9	+6.6 -6.6	-3.9 +4.1	+0.1 -0.1	+6.8 -6.8	+3.2 -3.2
3500	0.79	17	$3.01e-04$	+27.2 -21.3	+11.6 -11.5	+8.0 -7.8	-3.3 +3.5	+0.1 -0.2	+6.4 -6.4	+4.0 -4.0
4150	0.81	5	$8.19e-05$	+55.2 -35.2	+14.6 -15.0	+11.9 -11.8	-3.7 +2.1	+0.0 -0.3	+6.4 -6.4	+5.1 -5.1
5250	0.85	3	$1.98e-05$	+75.7 -42.3	+18.6 -18.1	+14.3 -14.3	-2.9 +4.7	+0.2 +0.0	+8.1 -8.1	+6.9 -6.9
7000	0.87	1	$5.56e-06$	+158.0 -58.0	+29.0 -26.4	+24.3 -24.1	-3.7 +10.7	-0.4 +0.0	-1.3 +1.3	+9.4 -9.4
9500	0.89	1	$5.60e-06$	+158.0 -58.0	+58.3 -63.4	+53.5 -53.5	-19.0 -0.7	+0.0 +0.0	+19.2 -19.2	+13.0 -13.0

## 2. WHY USE THE BIN COUNTS

The standard way to estimate the differential cross section at an  $(x, Q^2)$  point within a bin is given by

$$\frac{d\sigma^{\text{Data}}}{dxdQ^2} = \frac{d\sigma^{\text{theory}}}{dxdQ^2} \frac{\sigma^{\text{Data}}}{\sigma^{\text{theory}}}$$

with

$$\sigma^{\text{Data}} = \frac{N}{\mathcal{L}a}$$

where we will choose for concreteness  $N = 2$  as the observed number of events in our example below.  $\mathcal{L}$  is the integrated luminosity in the data set, and  $a$  is the acceptance in the bin:

$$a = \frac{N_{\text{reconstructed}}^{\text{MC}}}{N_{\text{generated}}^{\text{MC}}} \quad .$$

(1) The 'old-fashioned-standard-prescription'  $\sigma \pm \delta_\sigma$ :

$$\delta_\sigma = \frac{\sqrt{N}}{\mathcal{L}a} = \frac{\sigma}{\sqrt{N}} .$$

For  $N = 2$ , we have a fractional uncertainty of  $\delta_\sigma/\sigma = 0.71$ . For  $N = 0$ , this prescription breaks down and it has often been the case that the uncertainty for  $N = 1$  is taken.

(2) with confidence levels  $\sigma_{-\delta^{\text{down}}}^{+\delta^{\text{up}}}$  (central interval definition):

$$\begin{aligned} \delta^{\text{up}} &\rightarrow \sum_{i=0}^{i=N-1} P(N|\sigma + \delta^{\text{up}}) \leq 0.16 \\ \delta^{\text{down}} &\rightarrow \sum_{i=N+1}^{i=\infty} P(N|\sigma - \delta^{\text{down}}) \leq 0.16 \end{aligned}$$

(3) with credibility intervals  $\sigma_{-\delta^{\text{down}}}^{+\delta^{\text{up}}}$ :

$$\int_{\sigma - \delta^{\text{down}}}^{\sigma + \delta^{\text{up}}} P(\sigma'|N) d\sigma' = 0.68$$

where  $P(\sigma|N)$  is the posterior probability density for  $\sigma$  given the measured number of events. For  $N = 2$  and using a flat prior on the expected number of events gives  $\delta^{\text{down}} = 0.63$  and  $\delta^{\text{up}} = 2.62$

$2^{+2.62}_{-0.63}$

Now imagine that we have two pdf sets, one which predicts  $\nu_1 = 0.01$  and a second which predict  $\nu_2 = 4.5$  events in the bin. The standard fitting approach is to calculate a  $\chi^2$  by comparing the measured differential cross section to the predicted one, using the fractional uncertainty ON THE DATA quoted above. The resulting contributions to  $\chi^2$  using the three prescriptions given above are:

(1) 'old-fashioned-standard-prescription' :

$$\chi_1^2 = (N - \nu_1)^2 / N = 2.0$$

$$\chi_2^2 = (N - \nu_2)^2 / N = 3.1$$

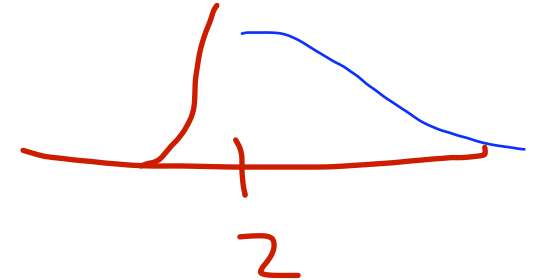
We see that the smaller prediction is preferred.

(2) confidence levels (similar to credibility below - not yet done):

(3) Credibility interval:

$$\chi_1^2 = (N - \nu_1)^2 / (0.63)^2 = 10.0$$

$$\chi_2^2 = (N - \nu_2)^2 / (2.62)^2 = 0.9$$



In this case, there is a strong preference for the higher prediction.

Now we see what happens if instead of performing a  $\chi^2$  fit using the differential cross section, we instead compare the predicted and observed event counts in the bin. The probabilities are

$$\begin{aligned}P(N = 2|\nu_1 = 0.01) &= \frac{e^{-0.01}0.01^2}{2!} = 5 \cdot 10^{-5} \\P(N = 2|\nu_2 = 4.5) &= \frac{e^{-4.5}4.5^2}{2!} = 0.11 \ .\end{aligned}$$

The preference for the larger prediction is even more pronounced. We see that in this case, there is very strong discriminating power between the different predictions, and very different probabilities would be obtained for  $\nu = 0.1, 0.01, 0.001$ . This is because we are making use of the probability to observe the given number of events assuming the expectation - the statistical fluctuations are handled correctly. Note also that  $N = 0$  poses no problem in the fitting.

### 3. INFORMATION REQUIRED FOR EVENT COUNT FITTING

In order to make full use of the data as outlined above, more information has to be provided than is normally the case. We need to supply the information which allows the pdf fitters to calculate a predicted number of events in each of the bins in which we report an event count. This prediction has to be calculated for every instance of the pdf parameters by integrating over the full kinematic phase space:

$$\nu(\Delta x, \Delta Q^2) = \int_{\Delta x, \Delta Q^2} \int \frac{d\nu^{pred}(x', Q'^2)}{dx' dQ'^2} P(x, Q^2 | x', Q'^2) dx' dQ'^2 dx dQ^2$$

where  $\frac{d\nu^{pred}(x', Q'^2)}{dx' dQ'^2}$  is the probability density to have an event with true kinematics at  $(x', Q'^2)$  and  $P(x, Q^2 | x', Q'^2)$  is the probability density to reconstruct an event at  $(x, Q^2)$  given the true kinematics at  $(x', Q'^2)$ . We approximate this integral with

$$\nu_j = \sum_i a_{ij} \nu_{\underline{i}}$$

$$a_{ij} = \frac{\sum_{k=1}^{M_i} \omega_k I(k \in j)}{\sum_{k=1}^{M_i} \omega_k}$$

with  $M_i$  the number of events generated in bin  $i$ ,  $\omega_k$  the weight given to the  $k^{\text{th}}$  event, and  $I(k \in j) = 1$  if event  $k$  is reconstructed in bin  $j$ , else  $I(k \in j) = 0$ .

The matrix  $a_{ij}$  should be provided with the bins  $i$  covering the full phase space which can give non-negligible contribution to our predictions. I.e., the bins defined in  $i$  should go beyond the bins defined in  $j$ , and can be made finer to reduce migration uncertainties.

The probability for the model to yield the data is

$$P(D|\vec{\lambda}) = \prod_j \frac{e^{-\nu_j(\vec{\lambda})} \nu_j(\vec{\lambda})^{n_j}}{n_j!}$$

is being used.

To take account of the systematic uncertainties, we need to further provide the following matrices:

$\delta a^{\text{unc}}$  The matrix of uncorrelated systematic uncertainties (given as a difference from the nominal matrix  $a$ ). E.g., MC statistical uncertainties would enter into this matrix. The MC statistical uncertainty for  $a_{ij}$  is

$$\delta a_{ij}^{\text{unc}} = \frac{\sqrt{\sum_{k=1}^{M_i} \omega_k^2 I(k \in j)}}{\sum_{k=1}^{M_i} \omega_k}$$

Note that we should have probably at least 25 MC events in all bins and no large spread of weights for this to be reasonably accurate.

$\delta a^l$  The correlated systematic uncertainties. For each of our correlated systematic uncertainties,  $l$ , we generate a new matrix  $a^l$  given from a one sigma deviation in the quantity of interest, and then calculate

$$\delta a_{ij}^l = a_{ij}^l - a_{ij} \quad .$$

The prediction including systematic uncertainties is

$$\nu'_j = \nu_j + \sum_i x_i \delta a_{ij}^{\text{unc}} \nu_i + \sum_l x_l \sum_i \delta a_{ij}^l \nu_i \quad .$$

The  $x$ 's are drawn from a Normal distribution of width 1, and a penalty is added to the  $\chi^2$  of  $\sum_k x_k^2$ .

## The dream:

- Find someone interested in this project who is also technically savvy with the data
- Discuss the procedure to see if we believe it can be made to work. Should include discussions with the fitting teams.
- Run some simple examples to get a feeling for how much we can learn from 'doing it right'. E.g., translate  $P(\text{Data} | \text{prediction})$  into chi squared for fits.
- If the gain is significant, then carry this out.
- And hopefully someone in H1 would also be interested ...
- worst case – give some prescriptions to fitting teams how better to account for the 'probability of the data given the model' in their analysis.