

Systematic Uncertainties from Synthetic Datasets

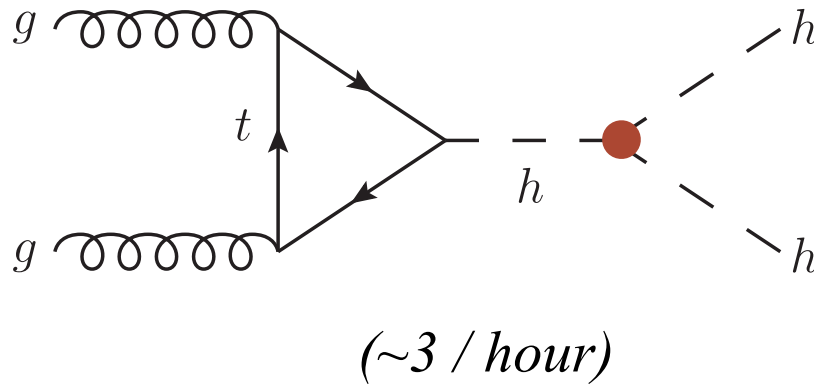
A case study with $HH \rightarrow 4b$

John Alison

Carnegie Mellon University

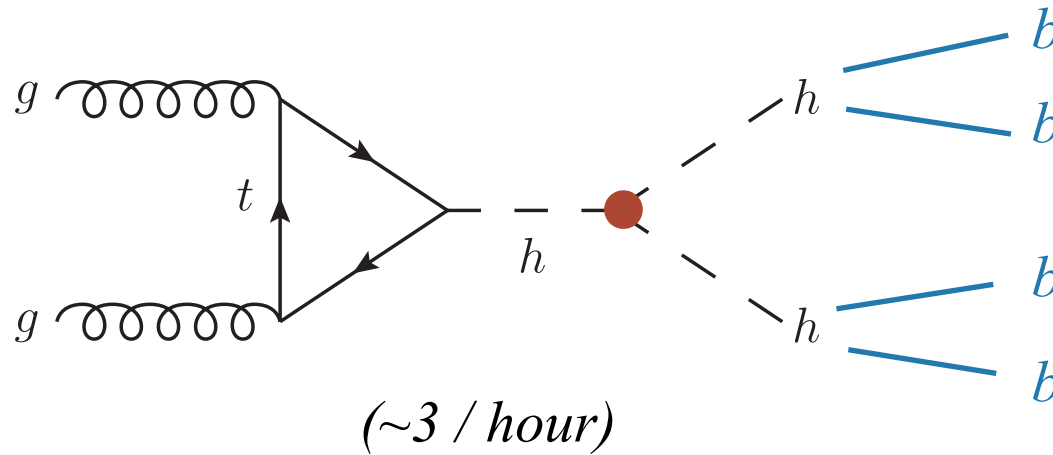
Motivation

Measuring the Higgs self-coupling λ major goal of the HL-LHC
Di-Higgs production most direct and most sensitivity way measure λ



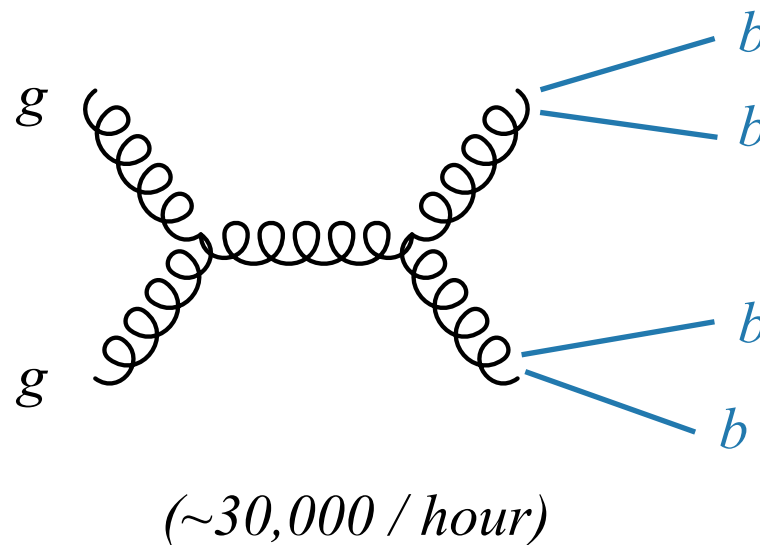
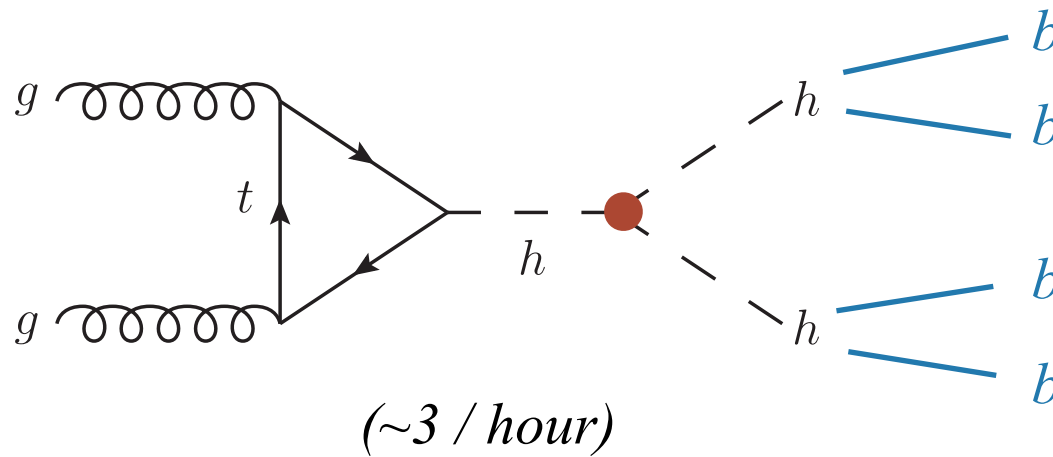
Motivation

Measuring the Higgs self-coupling λ major goal of the HL-LHC
Di-Higgs production most direct and most sensitivity way measure λ



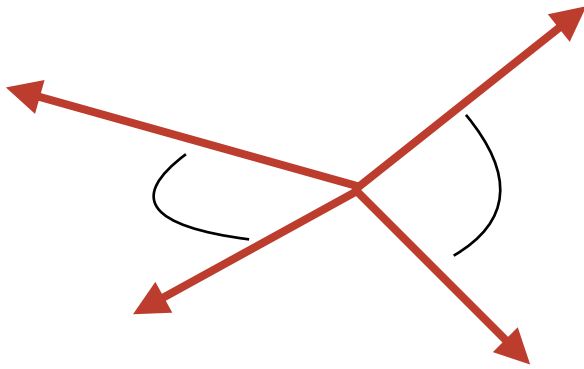
Motivation

Measuring the Higgs self-coupling λ major goal of the HL-LHC
Di-Higgs production most direct and most sensitivity way measure λ

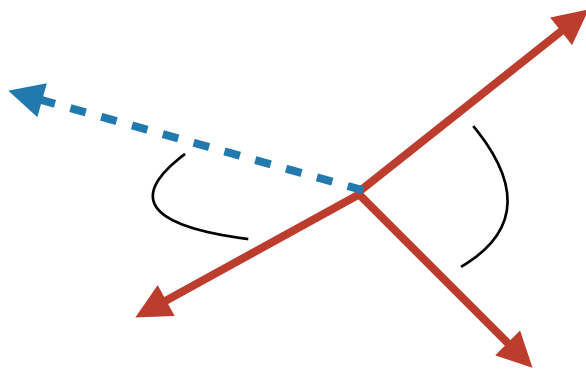


Data Driven Background: *ABCD*

“4b”

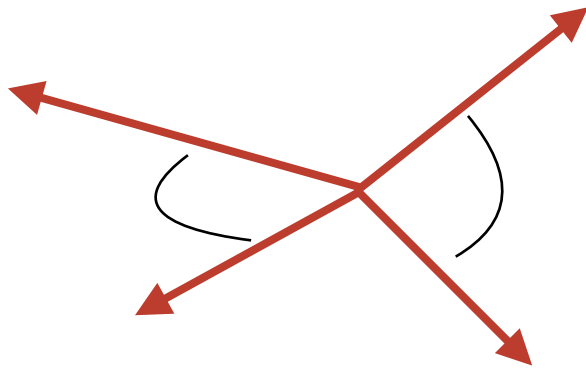


“3b”

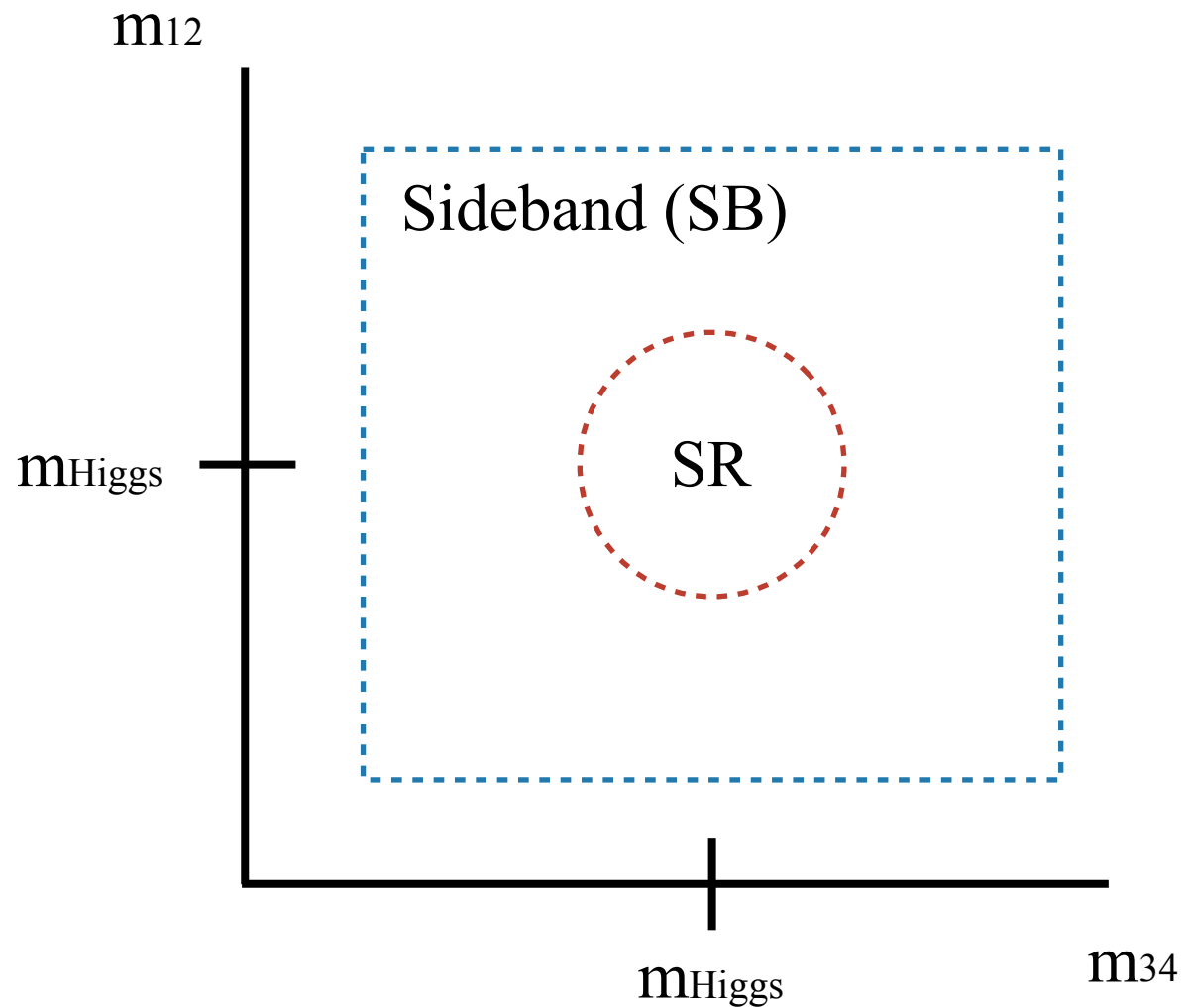
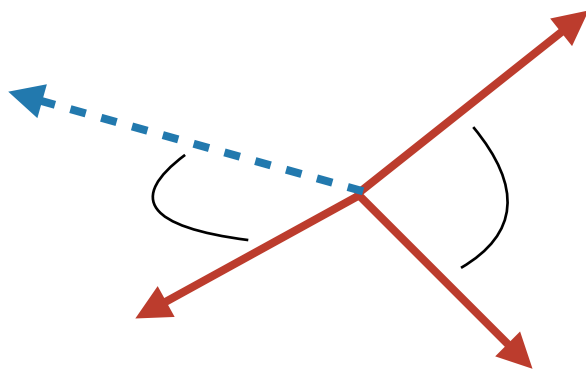


Data Driven Background: *ABCD*

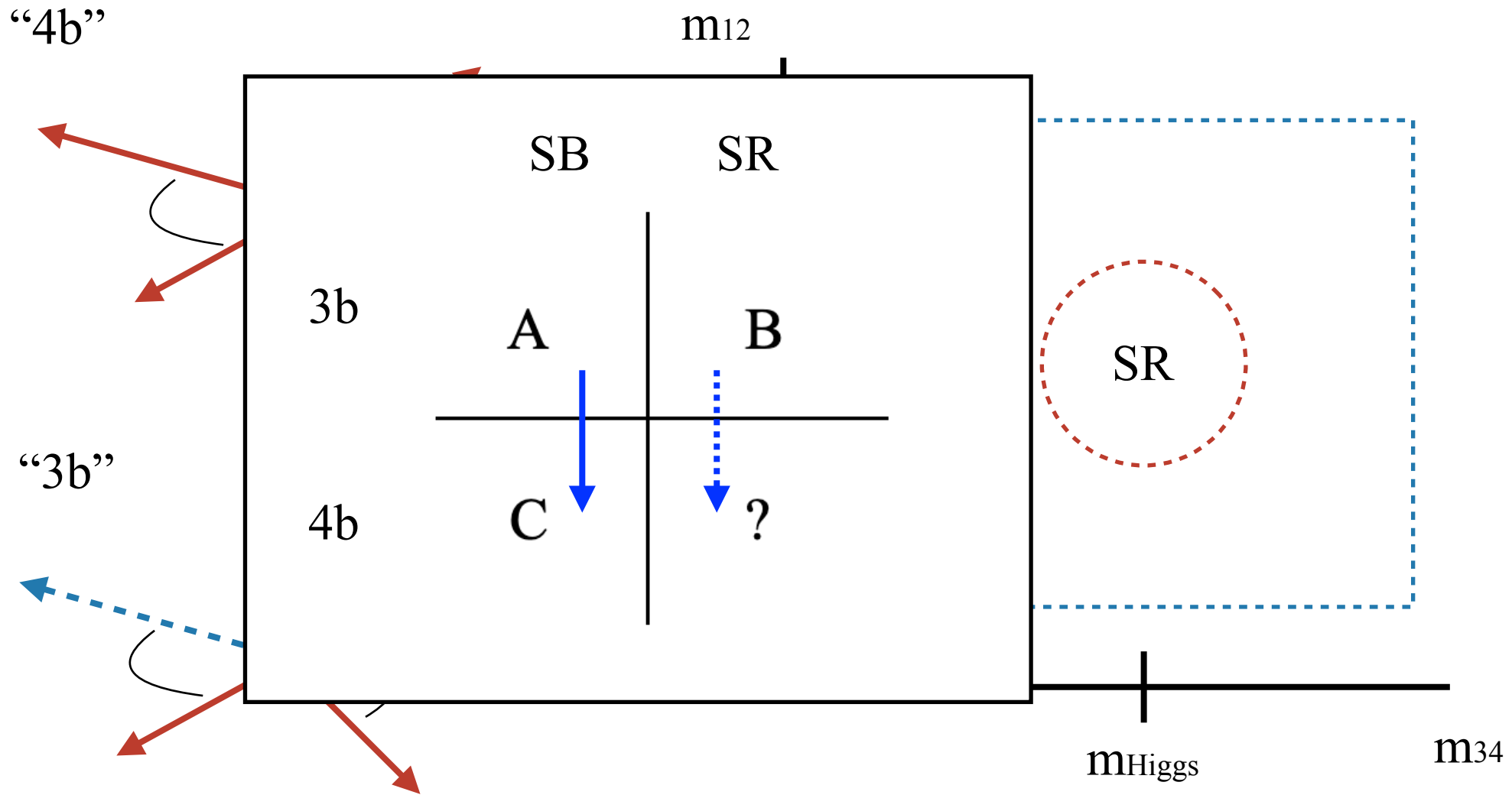
“4b”



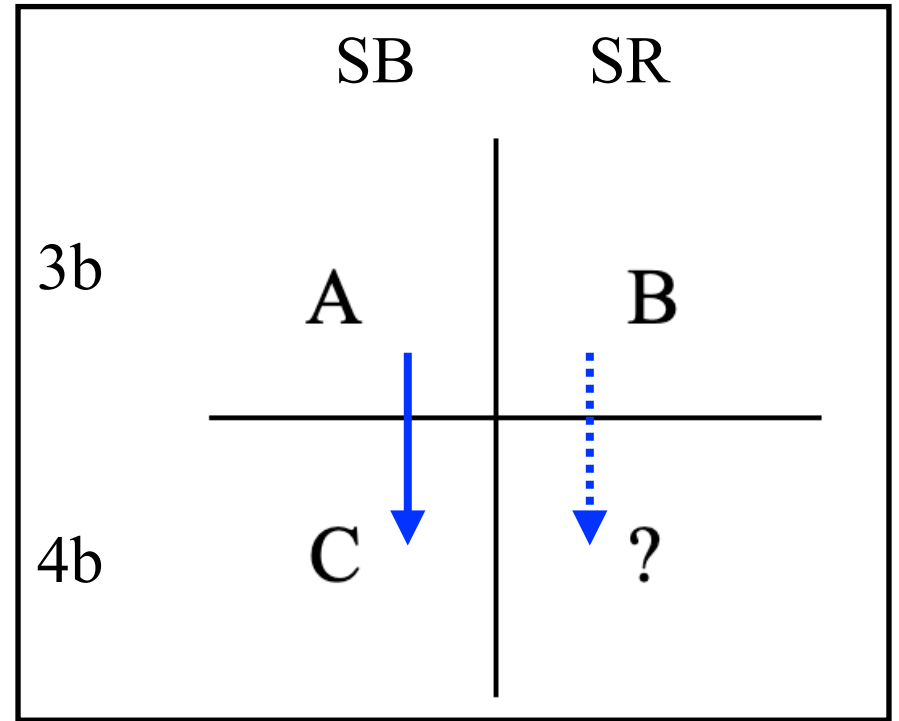
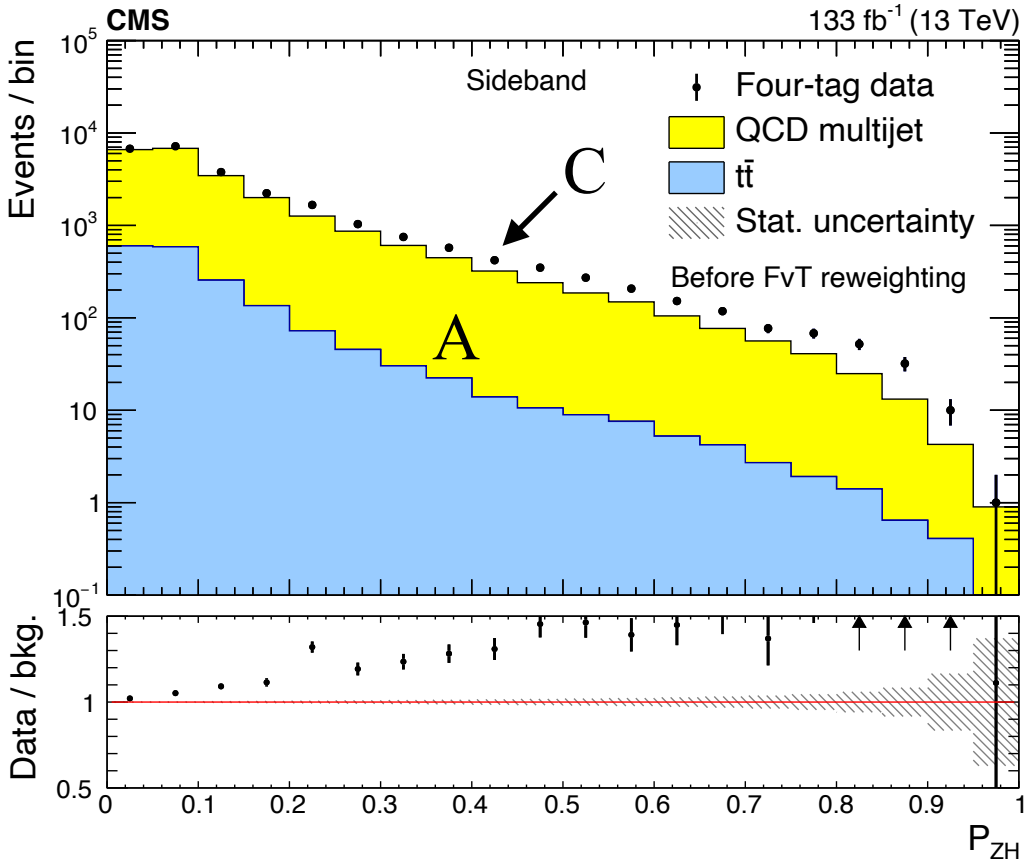
“3b”



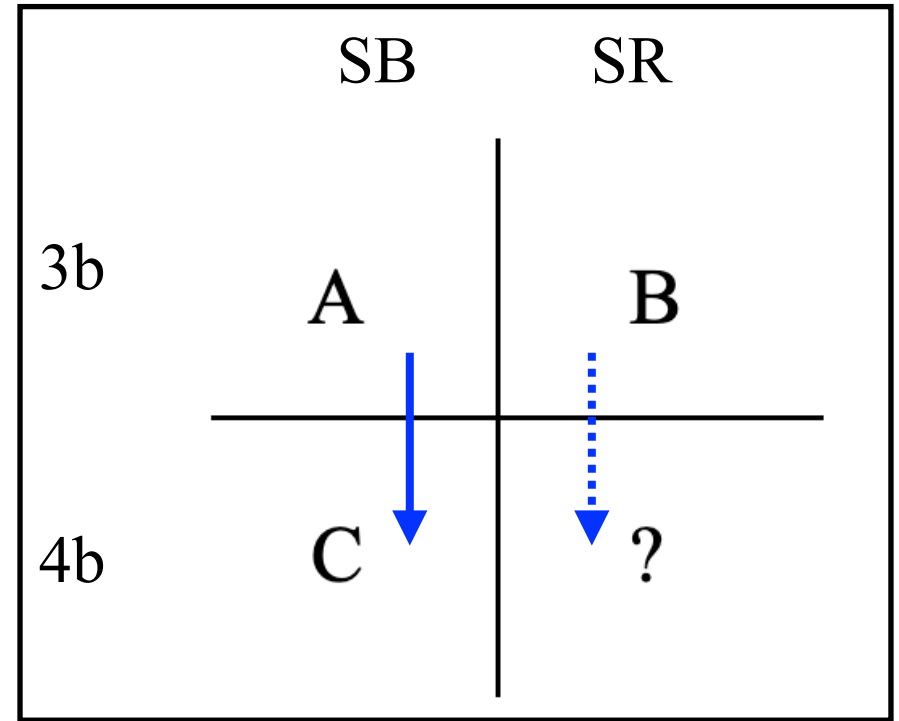
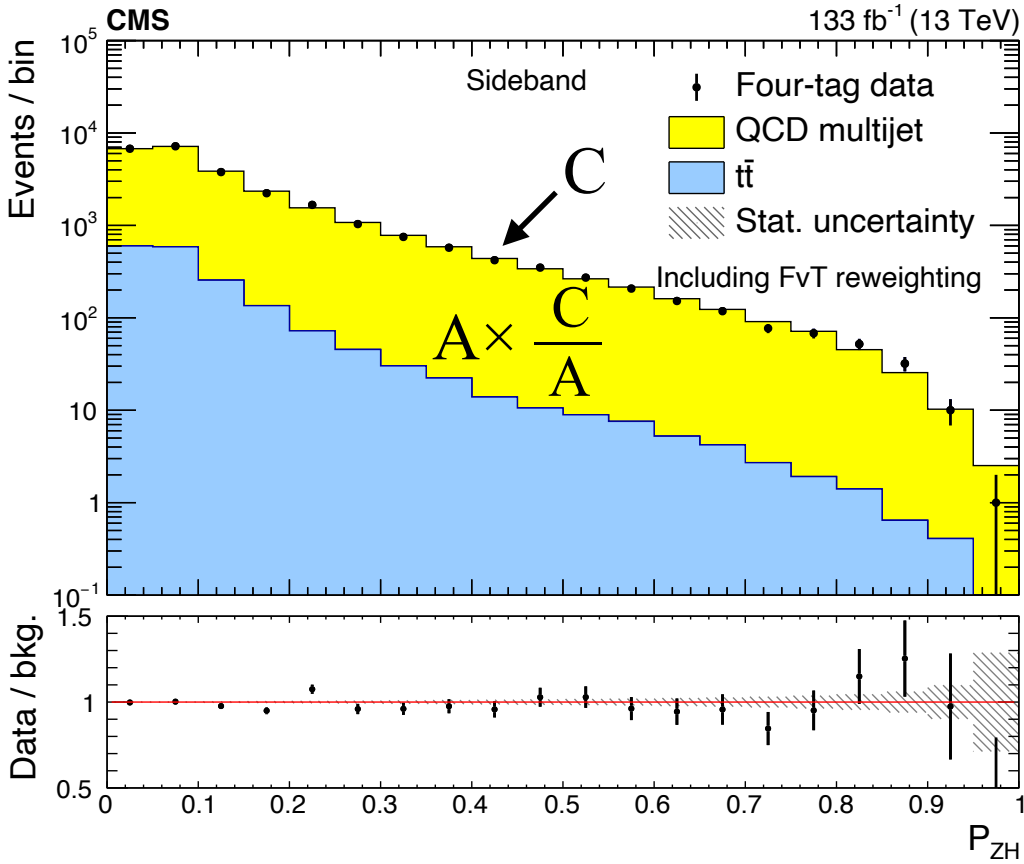
Data Driven Background: *ABCD*



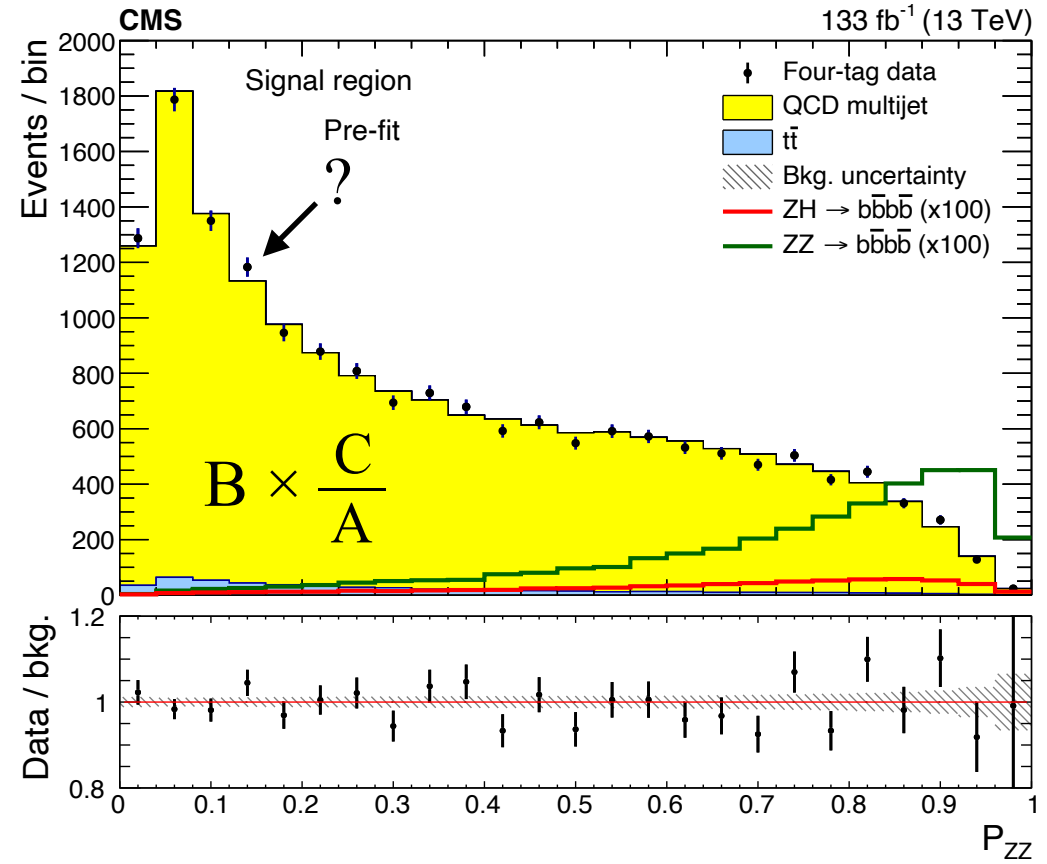
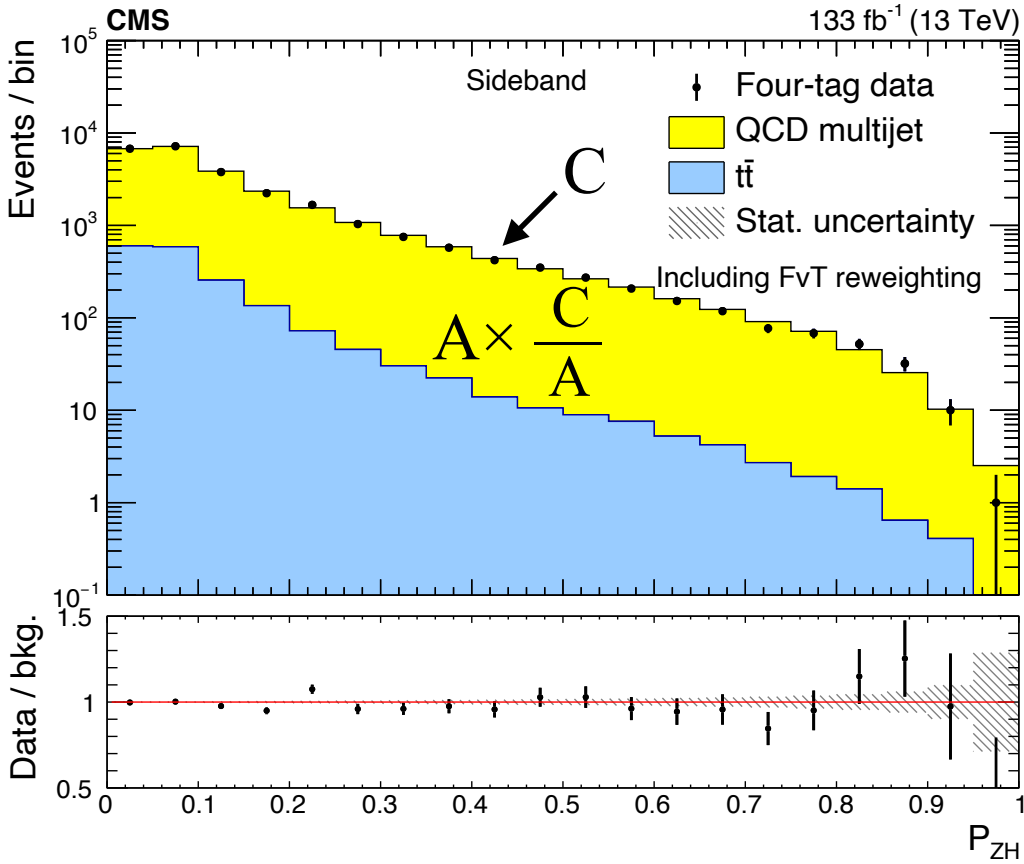
ABCD Method Works



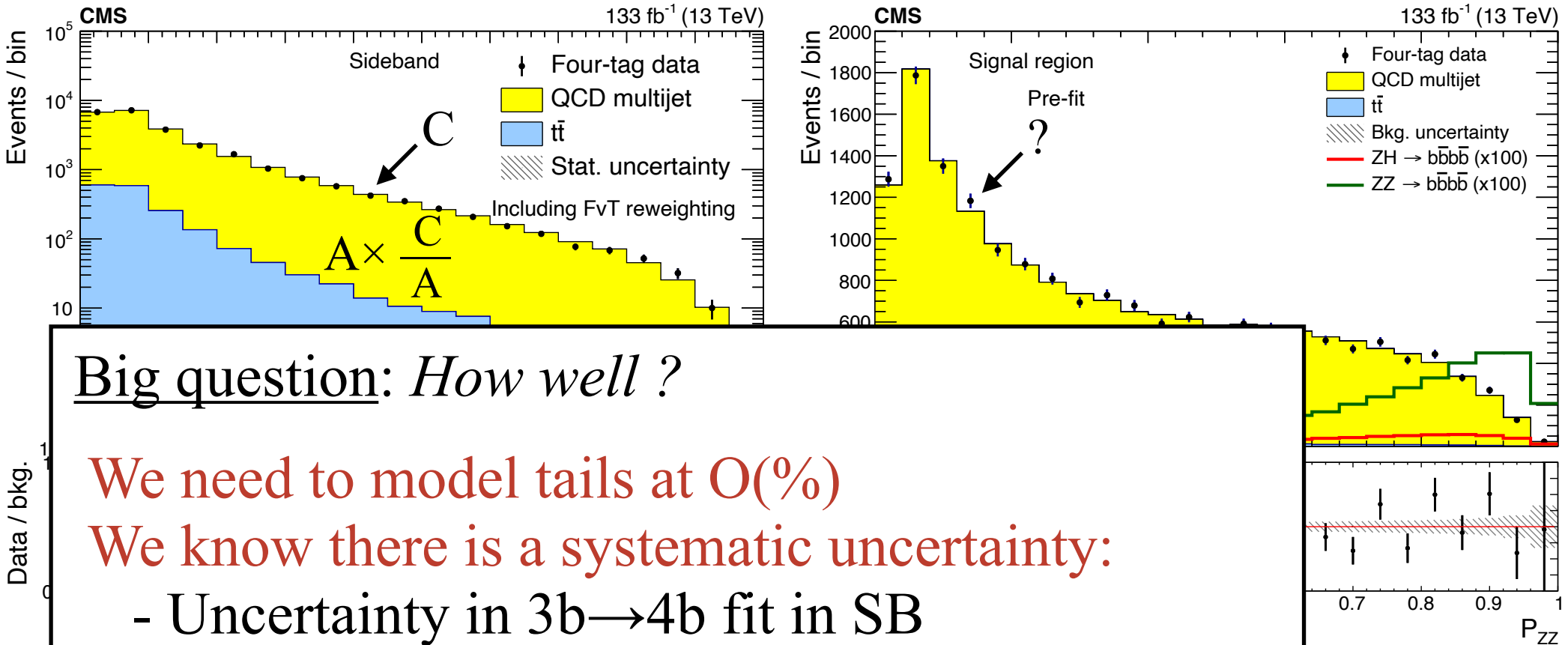
ABCD Method Works



ABCD Method Works



ABCD Method Works



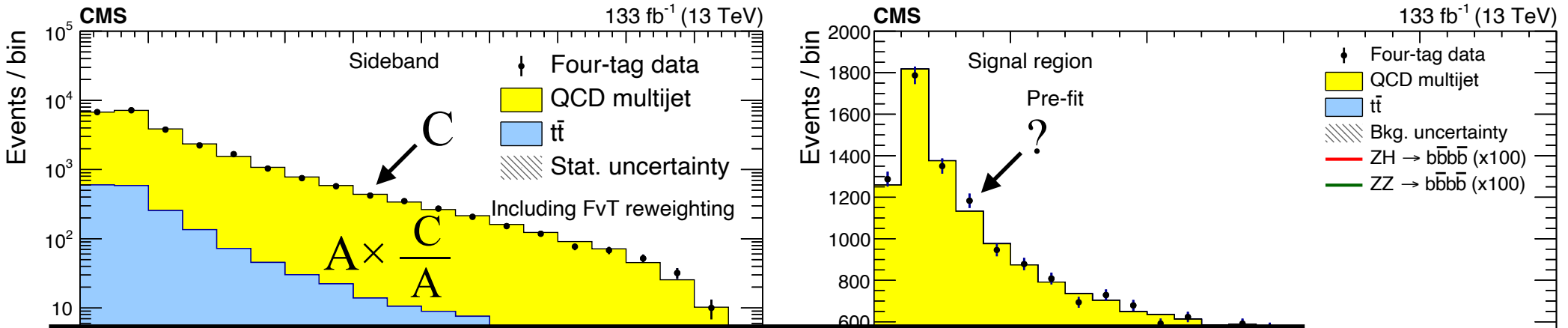
Big question: *How well?*

We need to model tails at O(%)

We know there is a systematic uncertainty:

- Uncertainty in 3b \rightarrow 4b fit in SB
- Extrapolation (domain shift) from CR \rightarrow SR
- ...

ABCD Method Works



Big question: *How well?*

We need to model tails at O(%)

We know there is a systematic uncertainty:

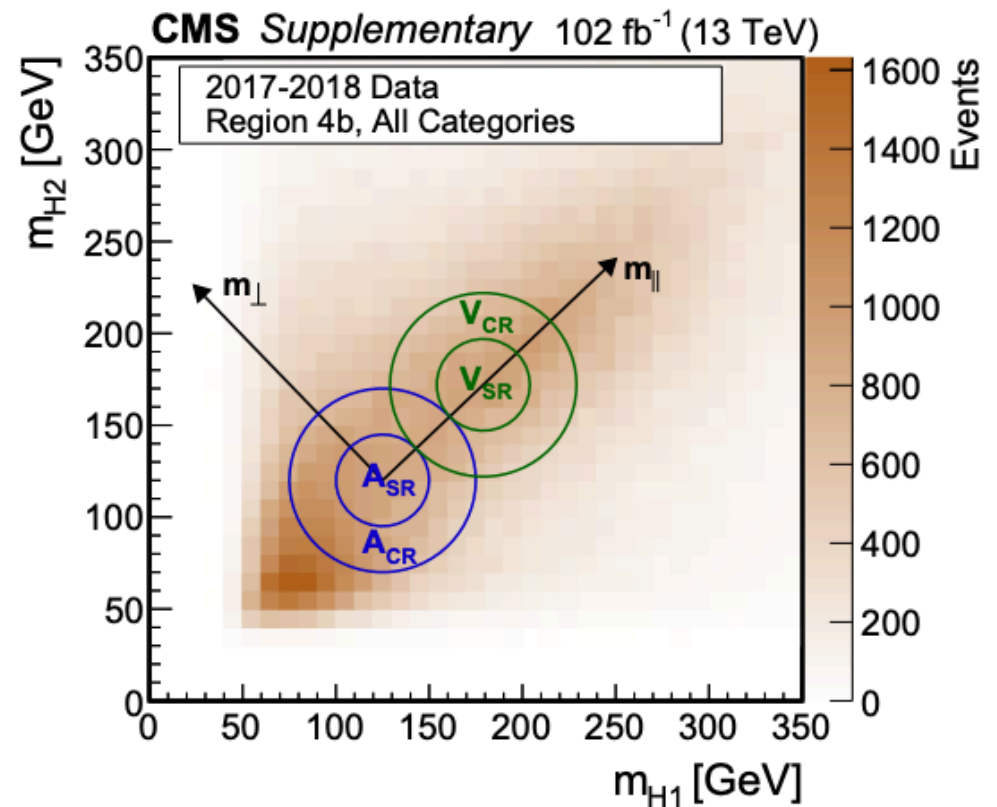
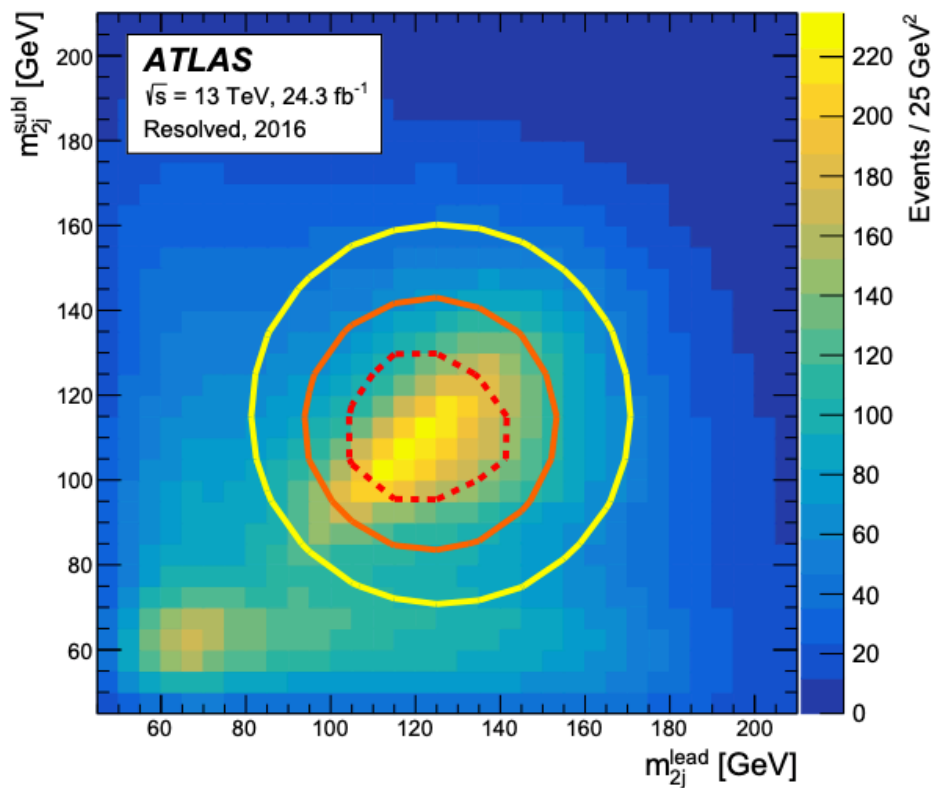
- Uncertainty in $3b \rightarrow 4b$ fit in SB
- Extrapolation (domain shift) from CR \rightarrow SR

Rest of this talk is an idea for validating this uncertainty

Standard Solution: Validation Region

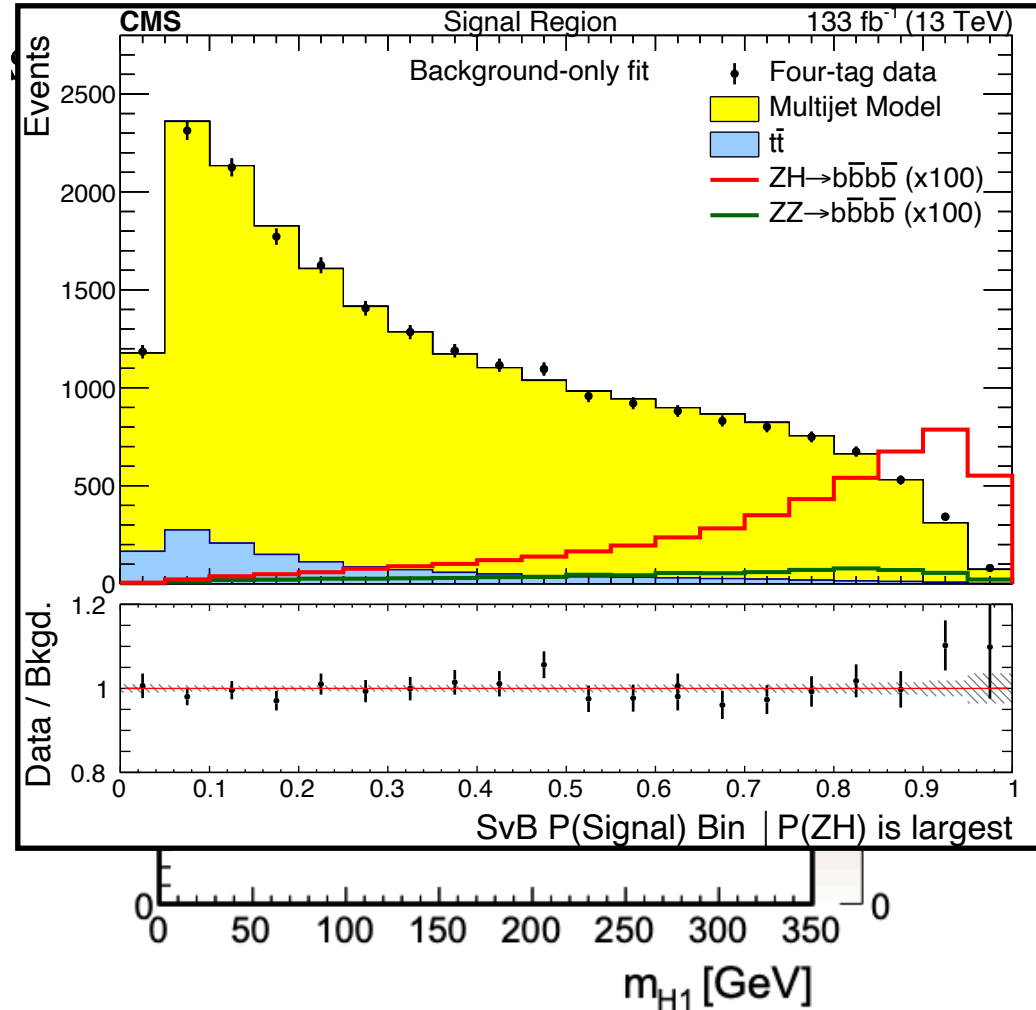
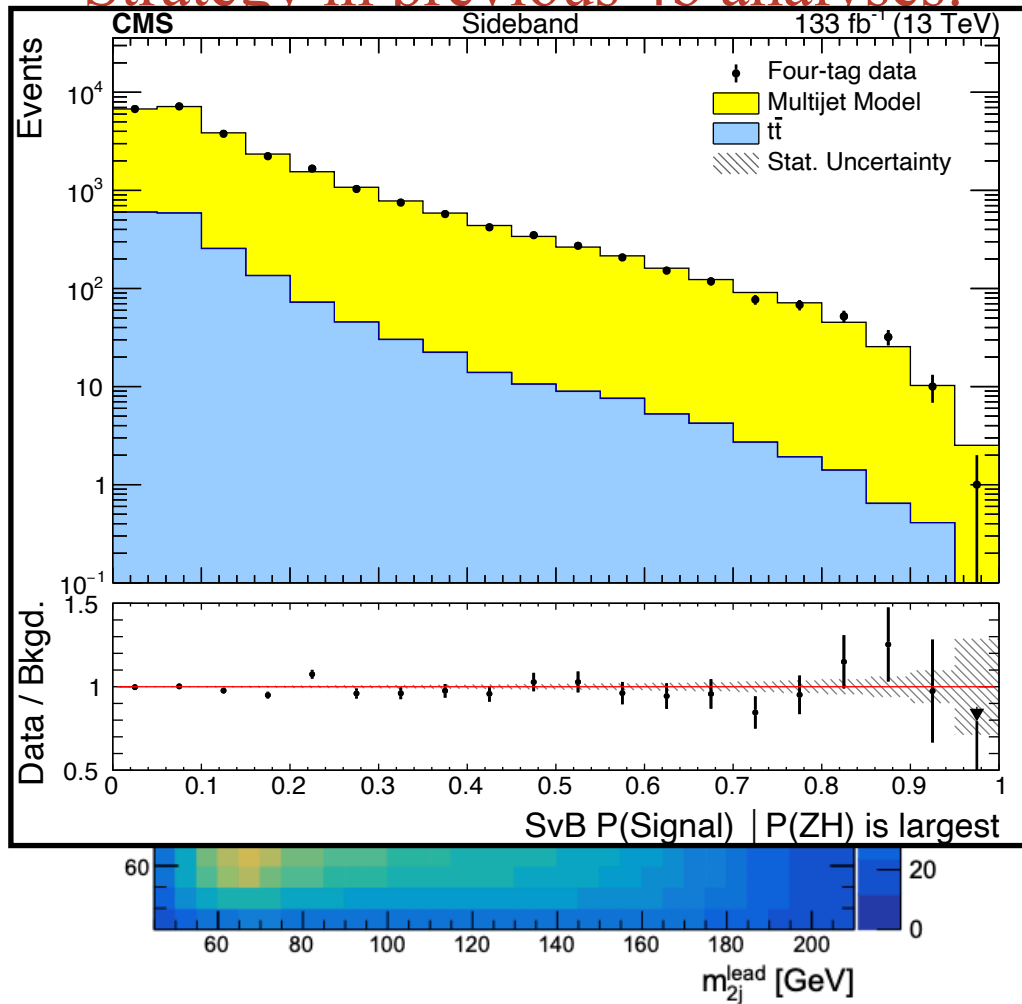
Strategy in previous 4b analyses:

Validated prediction in alternative signal-free region



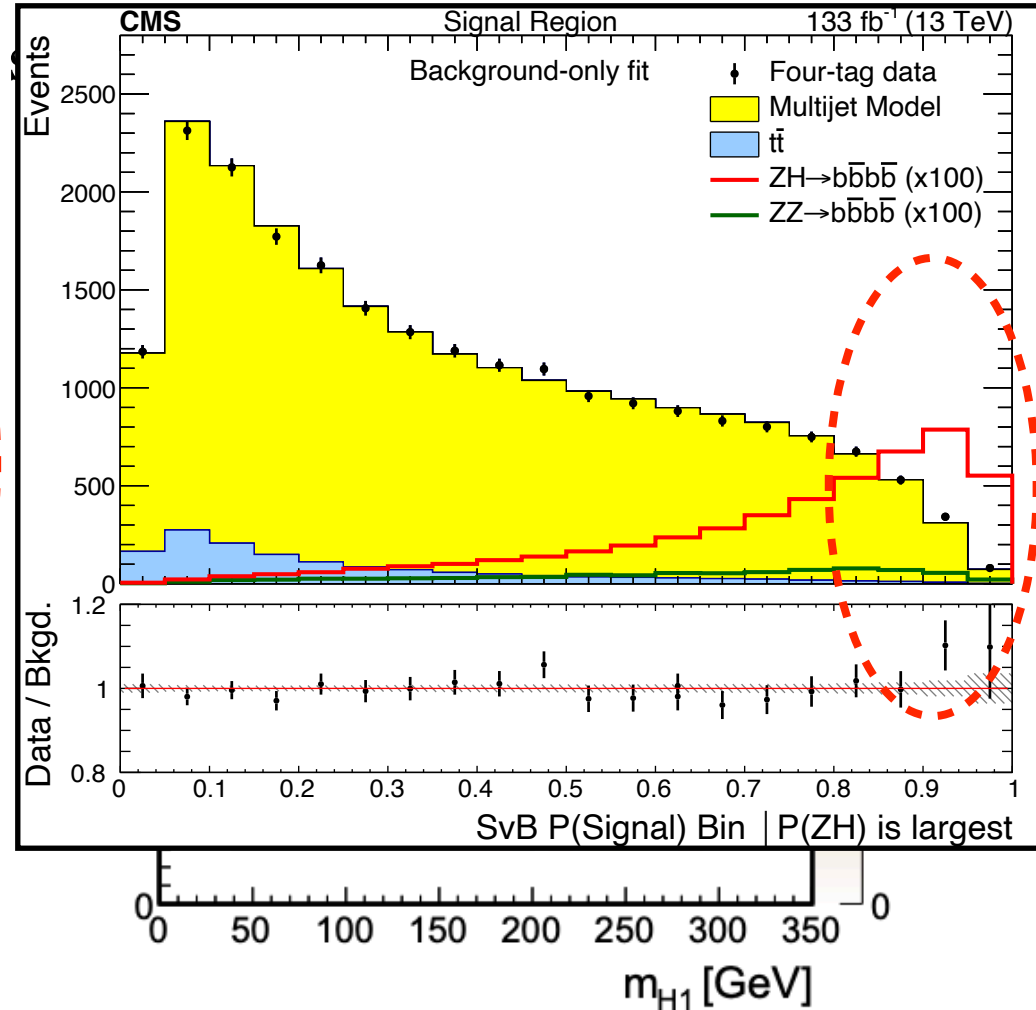
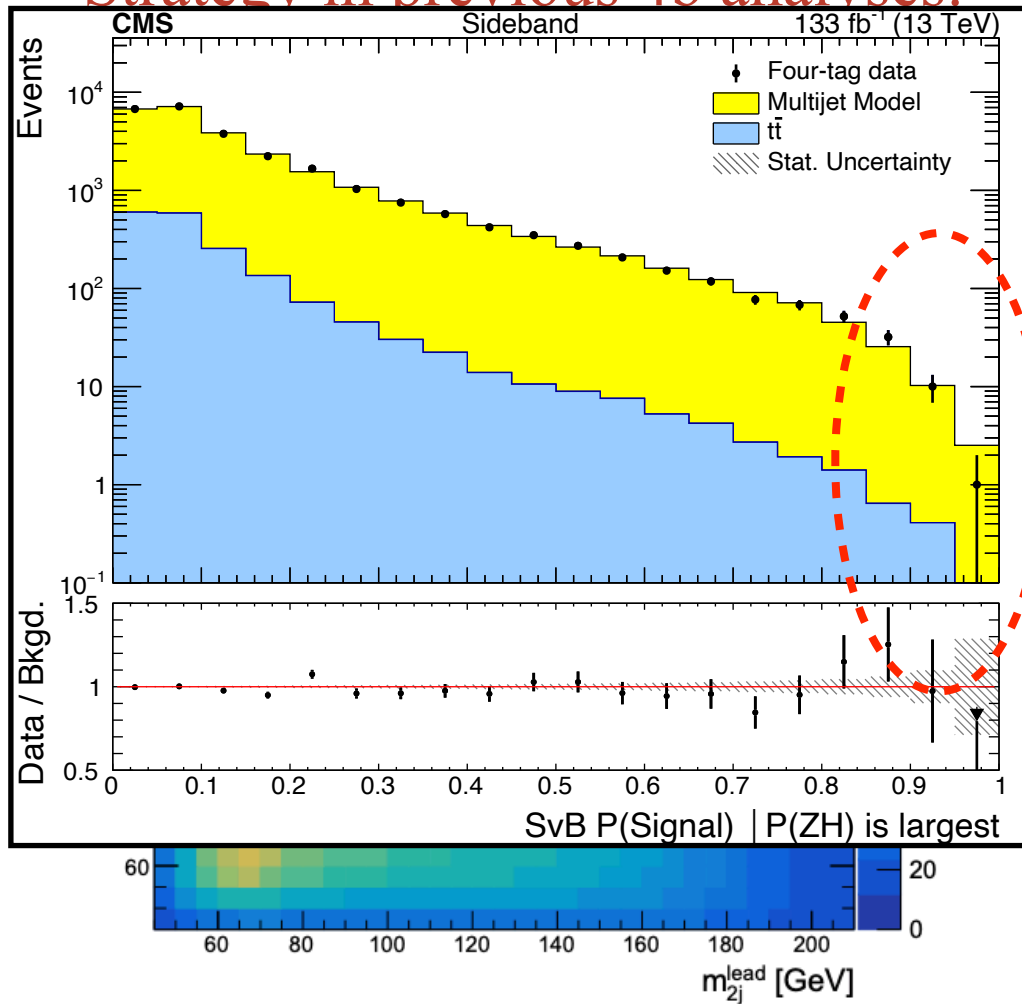
Standard Solution: Validation Region

Strategy in previous 4b analyses:



Standard Solution: Validation Region

Strategy in previous 4b analyses:



Aside: Solution with Optimal Transport

BACKGROUND MODELING FOR DOUBLE HIGGS BOSON PRODUCTION: DENSITY RATIOS AND OPTIMAL TRANSPORT

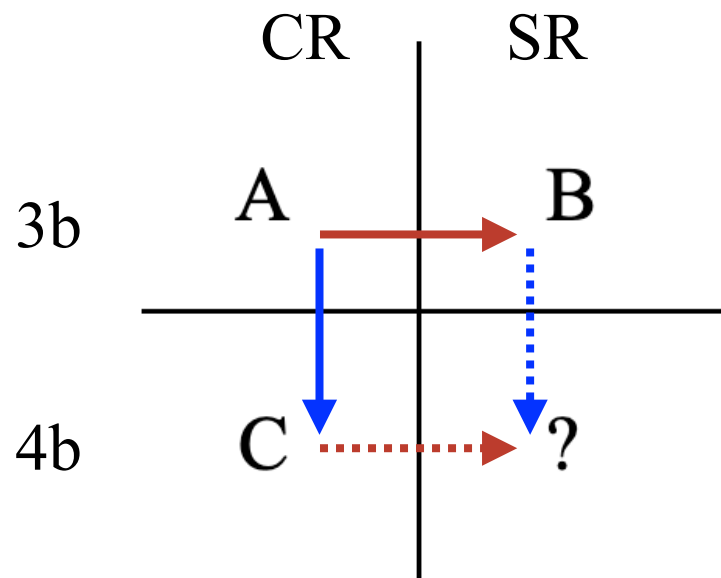
BY TUDOR MANOLE^a, PATRICK BRYANT^d, JOHN ALISON^e, MIKAEL
KUUSELA^b, AND LARRY WASSERMAN^c

*Department of Statistics and Data Science and NSF AI Planning Institute for Data-Driven Discovery in Physics,
Carnegie Mellon University*

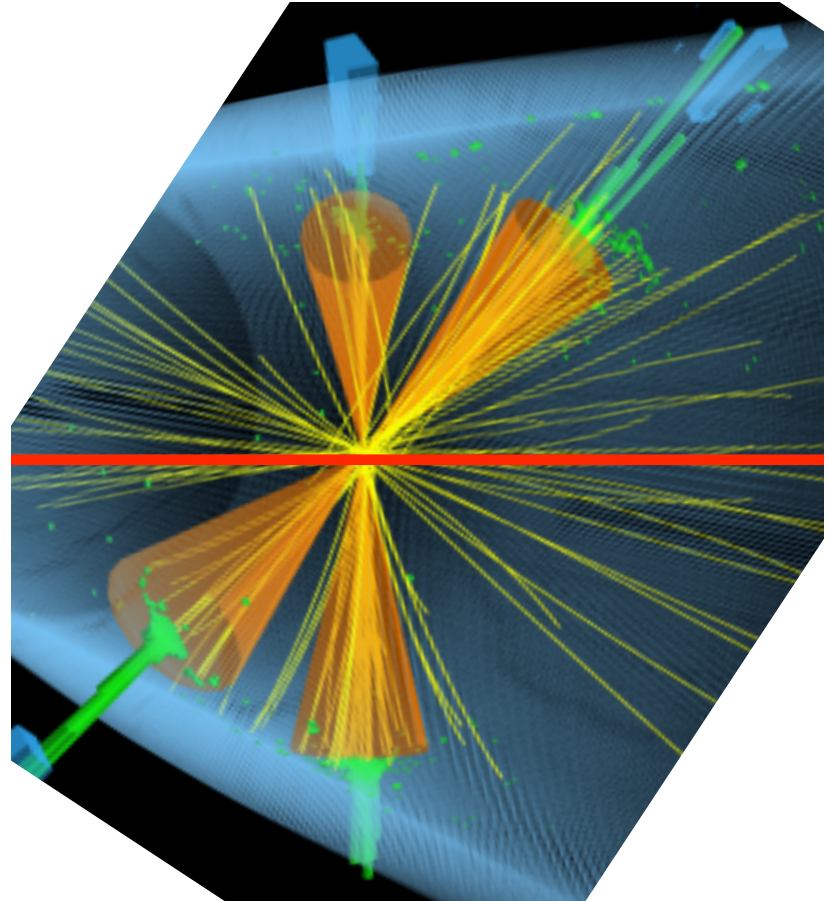
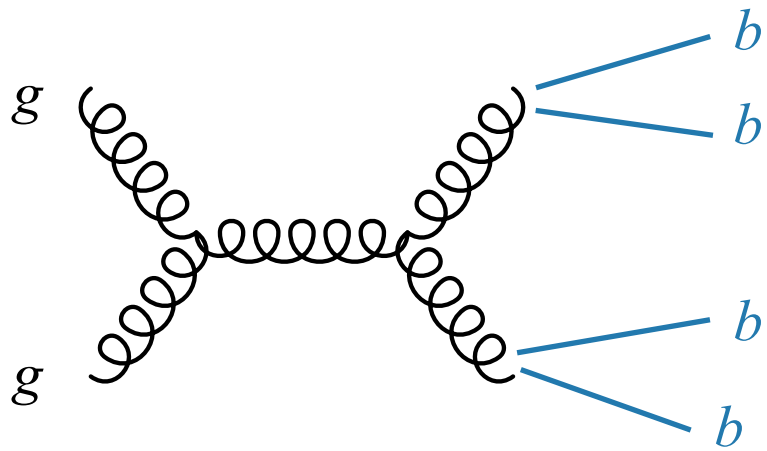
^atmanole@andrew.cmu.edu; ^bmkuusela@andrew.cmu.edu; ^clarry@stat.cmu.edu

*Department of Physics and NSF AI Planning Institute for Data-Driven Discovery in Physics,
Carnegie Mellon University*

^dpbryant2@andrew.cmu.edu; ^ejohnalison@cmu.edu

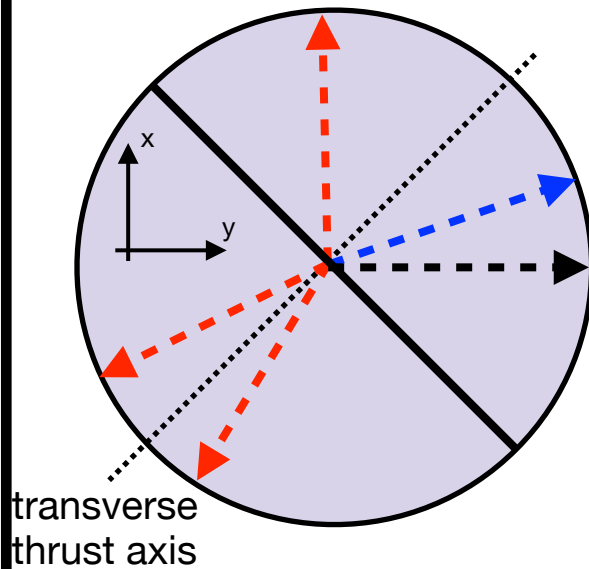


Synthetic Datasets: Event Mixing

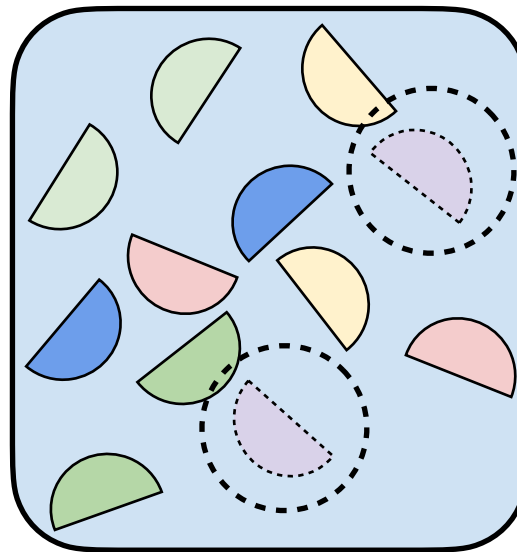


Synthetic Datasets: Event Mixing

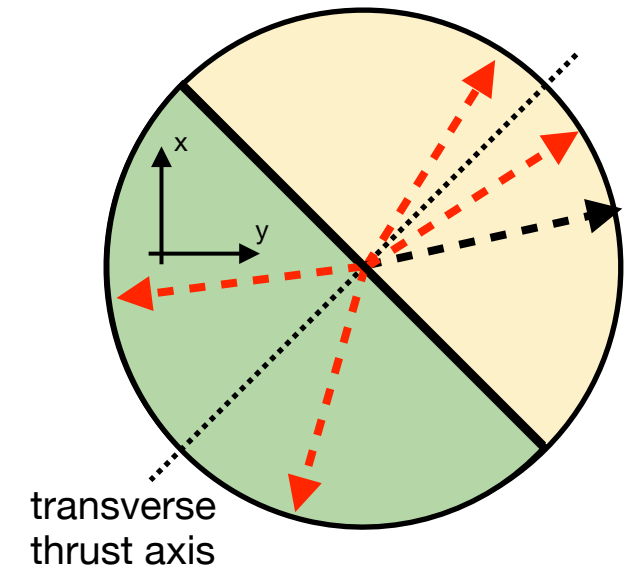
Original three-tag event
split into two hemispheres



Hemisphere library
made from four-tag events
filled in 1st pass, queried on 2nd



Mixed Event
using replaced hemispheres

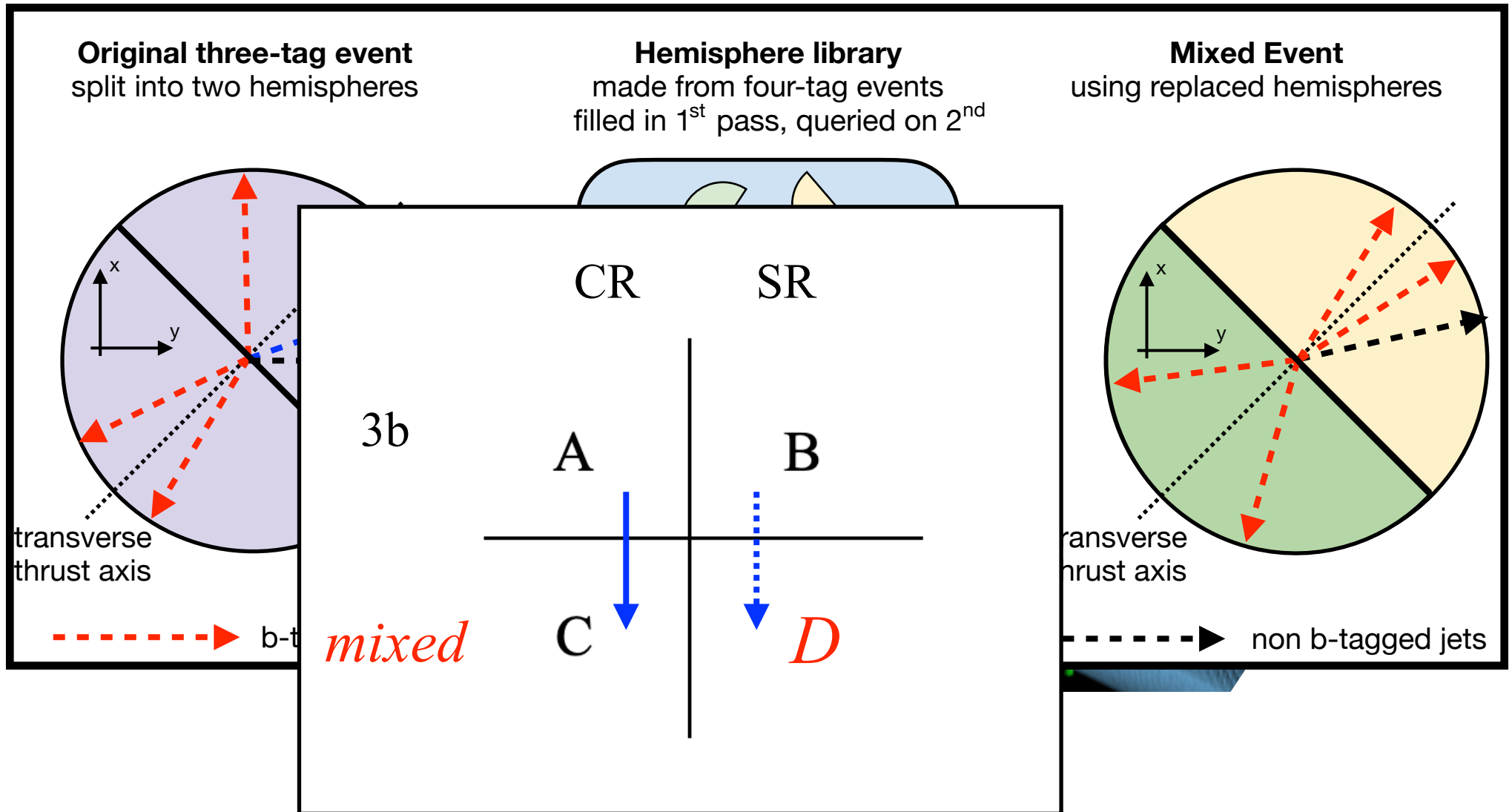


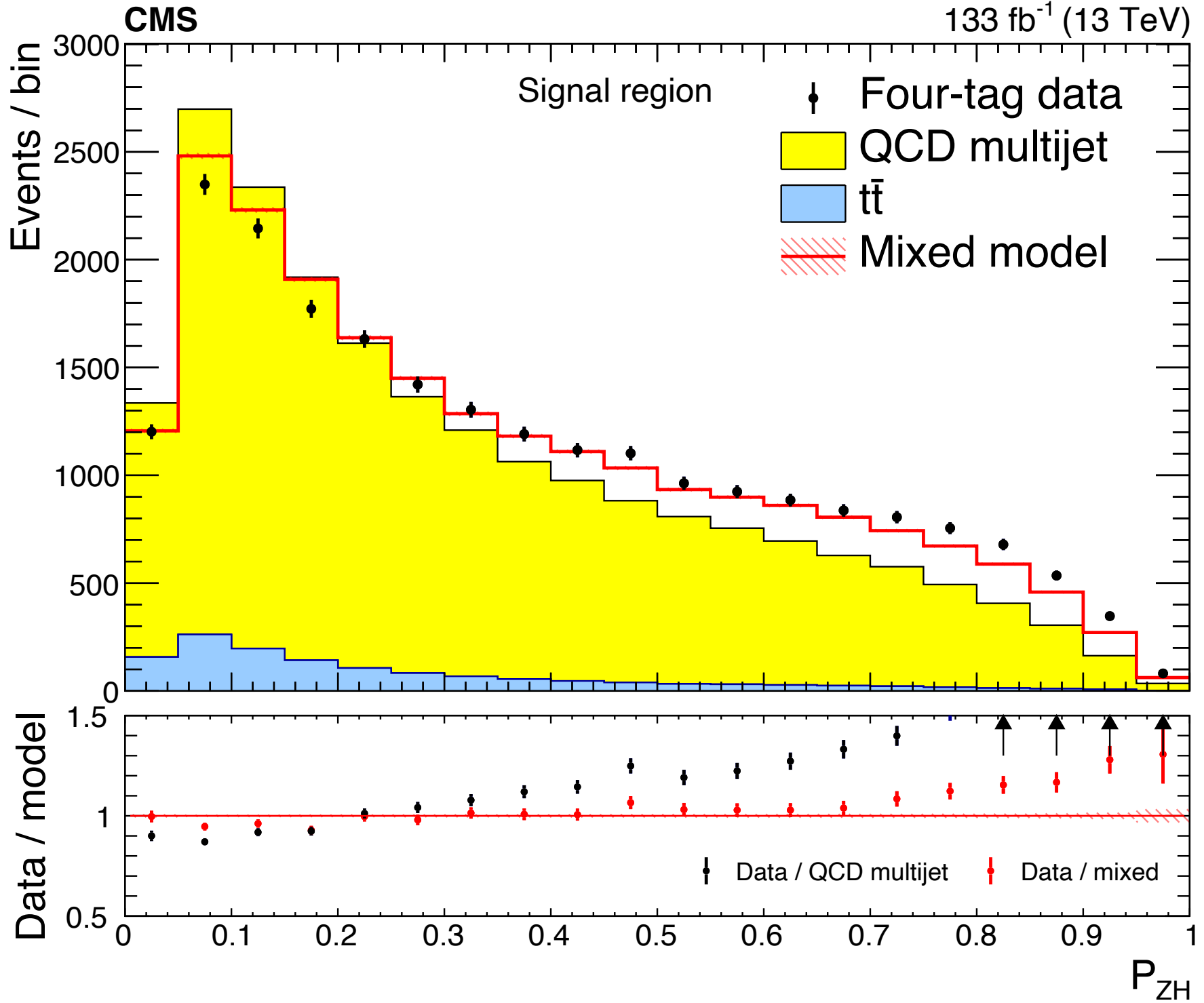
---▶ b-tagged jets

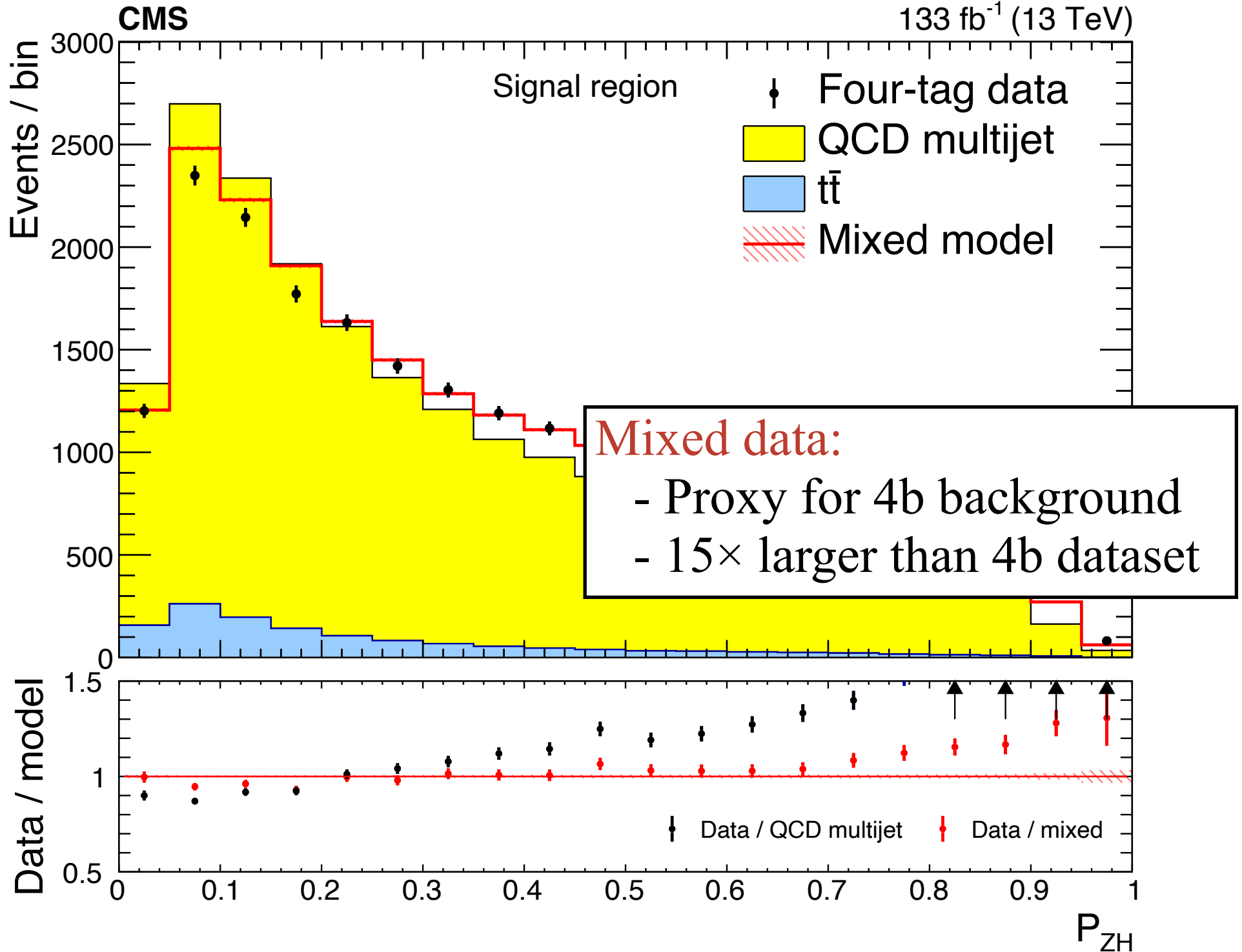
---▶ pseudo-tagged jets

---▶ non b-tagged jets

Synthetic Datasets: Event Mixing







Systematics with Mixed Data

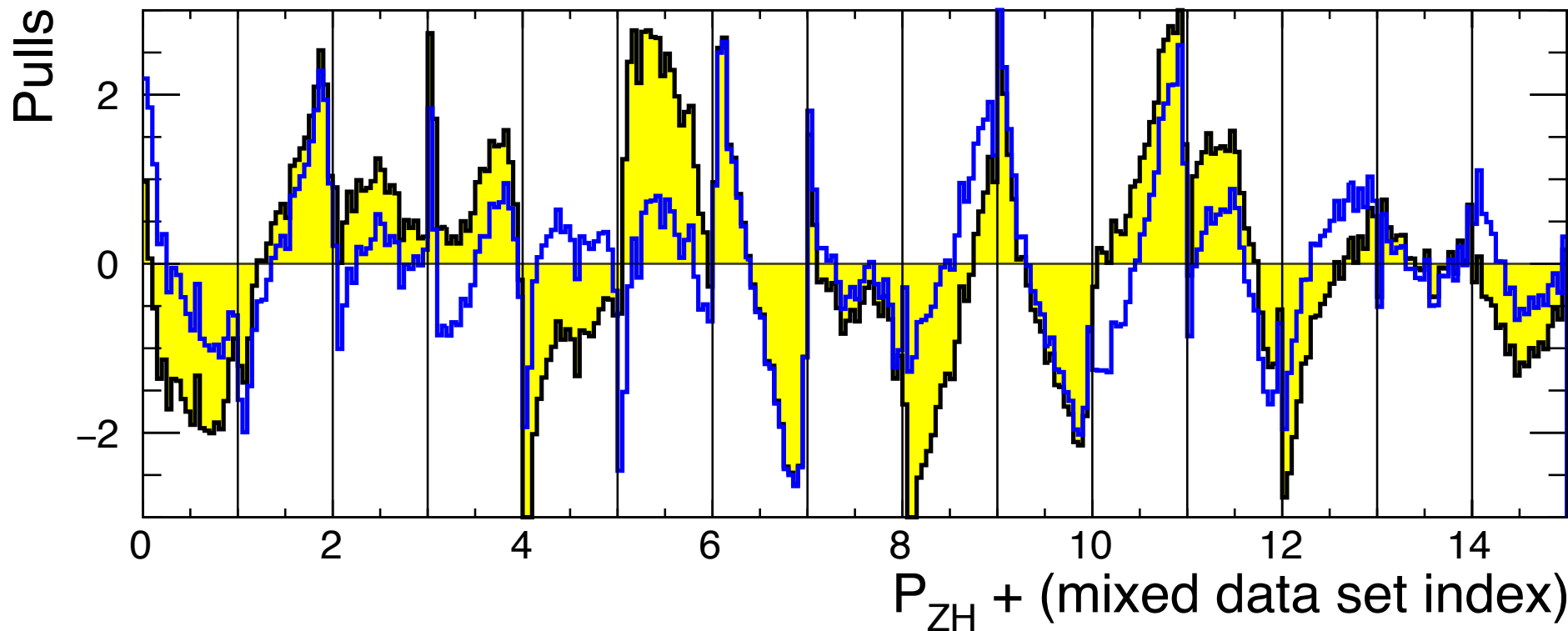
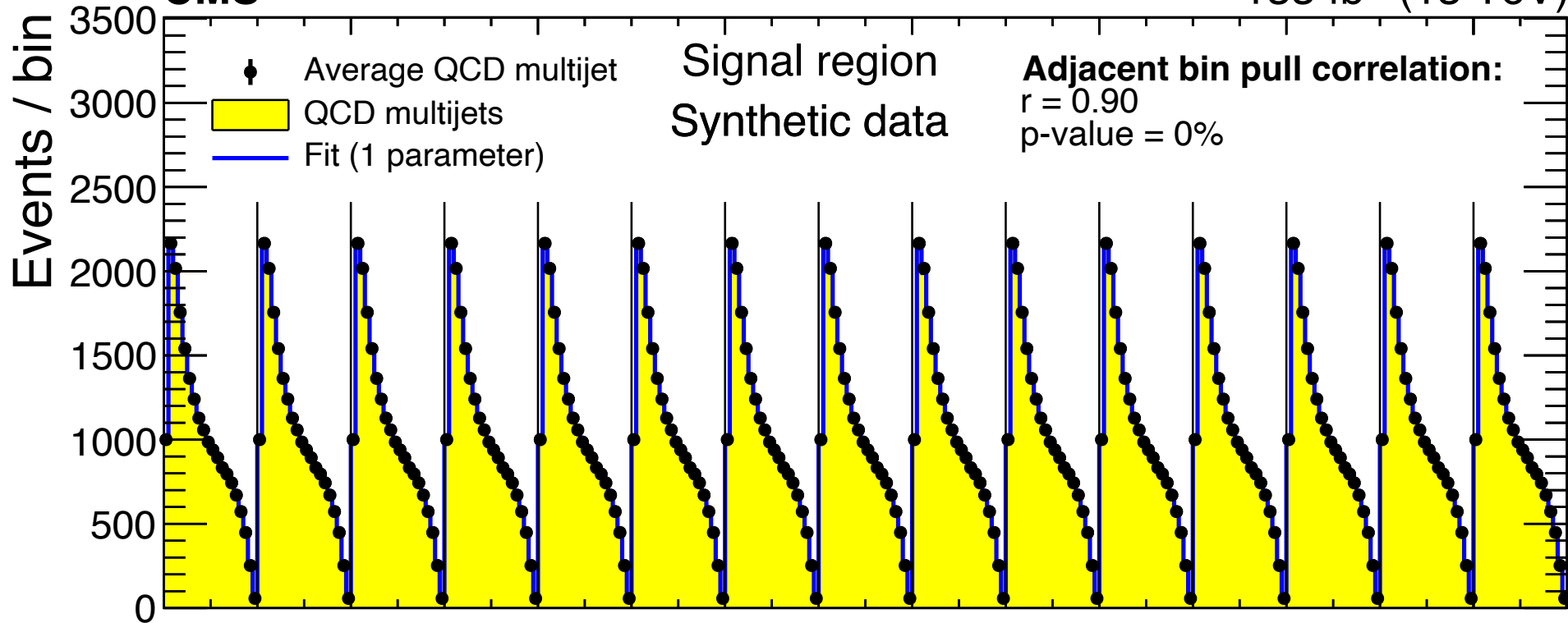
Consider three sources of potential systematic uncertainty

Variance: Arises from multi-variate classifier fit finite dataset CR

Bias: Assumption that density ratio measured in CR is same as SR

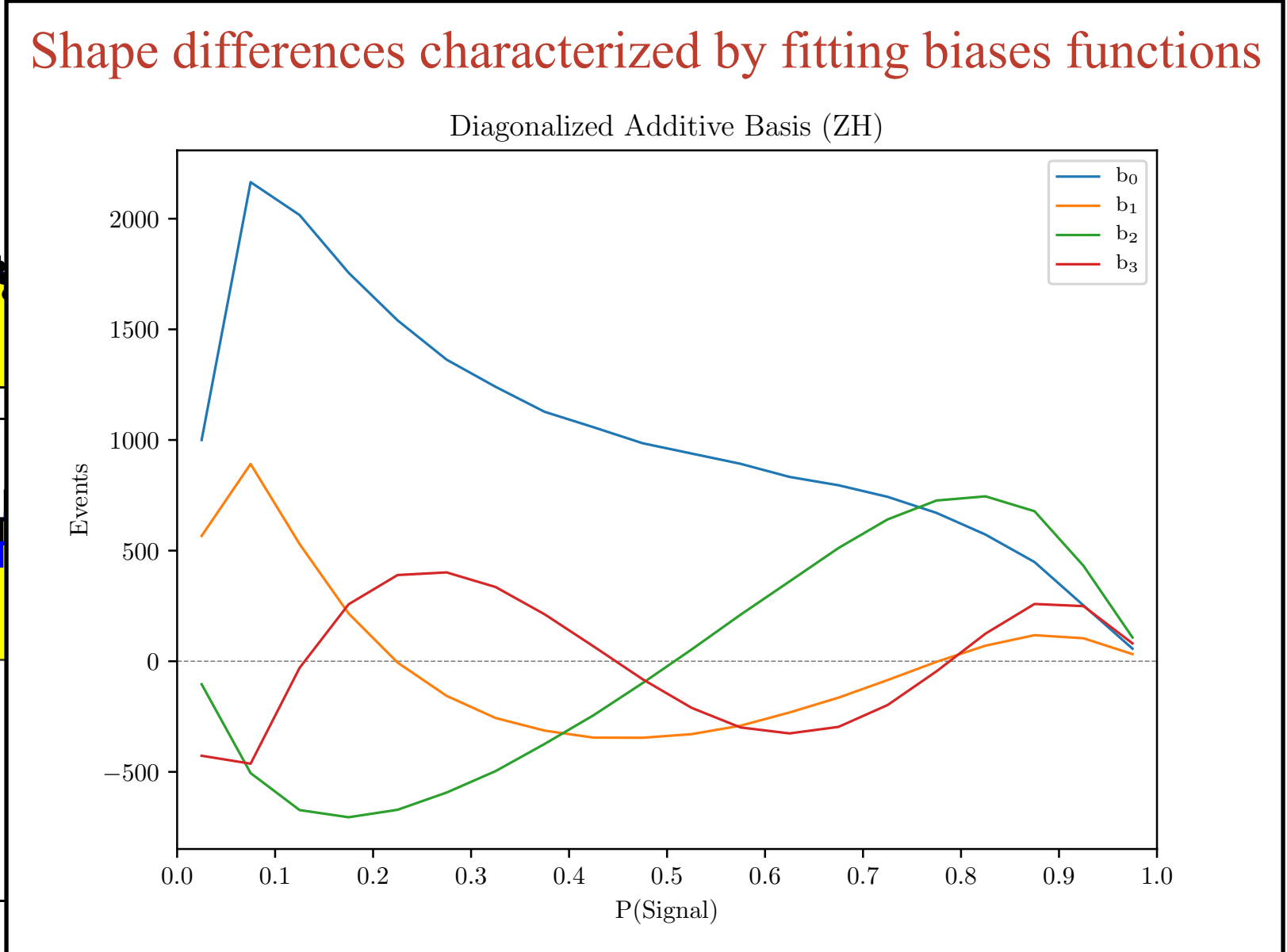
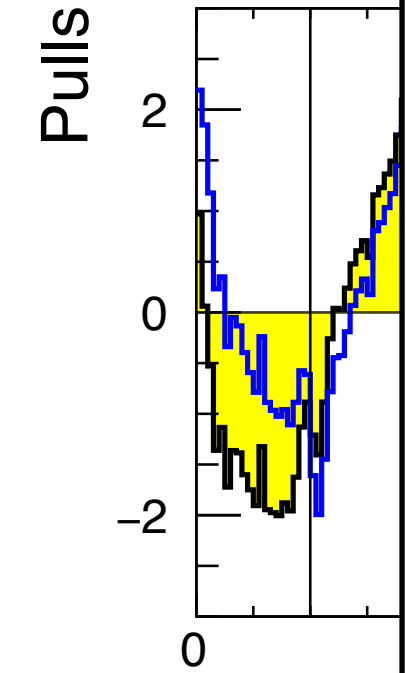
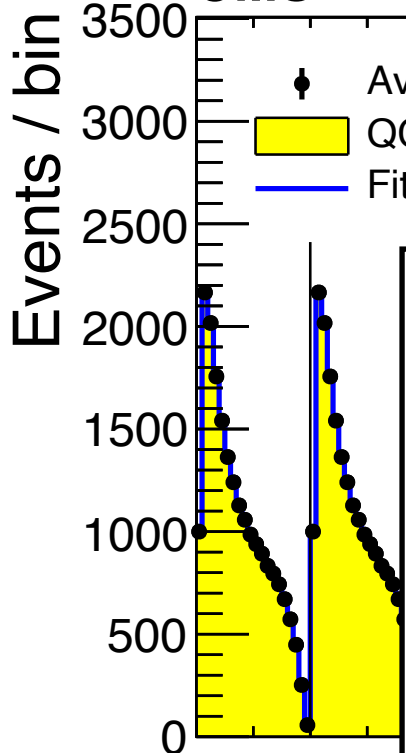
Spurious Signal Can misspecification of the background model
look like a signal under the null hypothesis ?
(*See backup*)

Assumptions rigorously defined for stats audience: [arXiv:2208.02807](https://arxiv.org/abs/2208.02807)



Signal region
Synthetic data

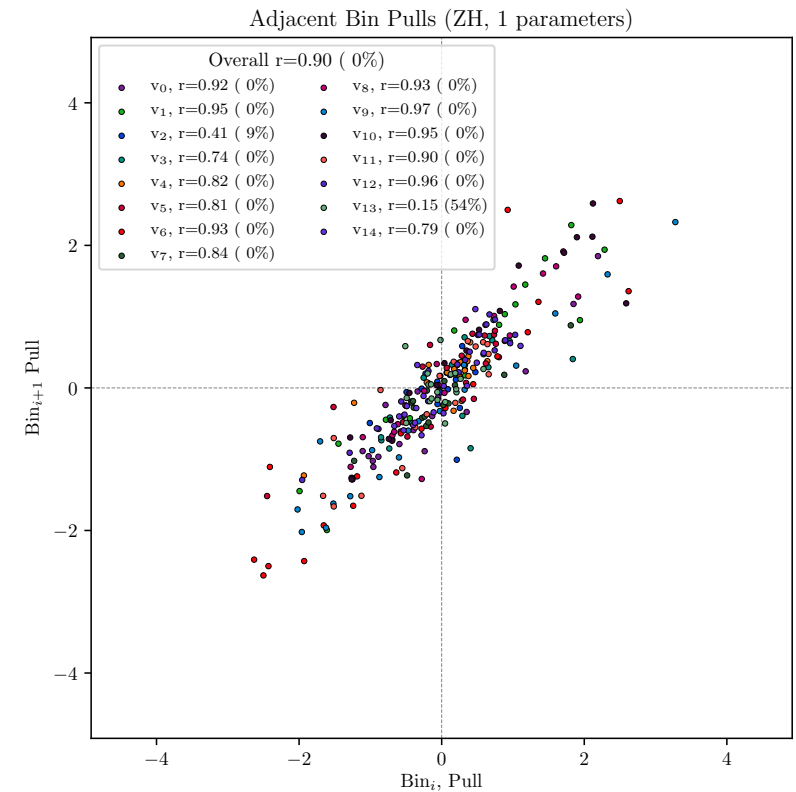
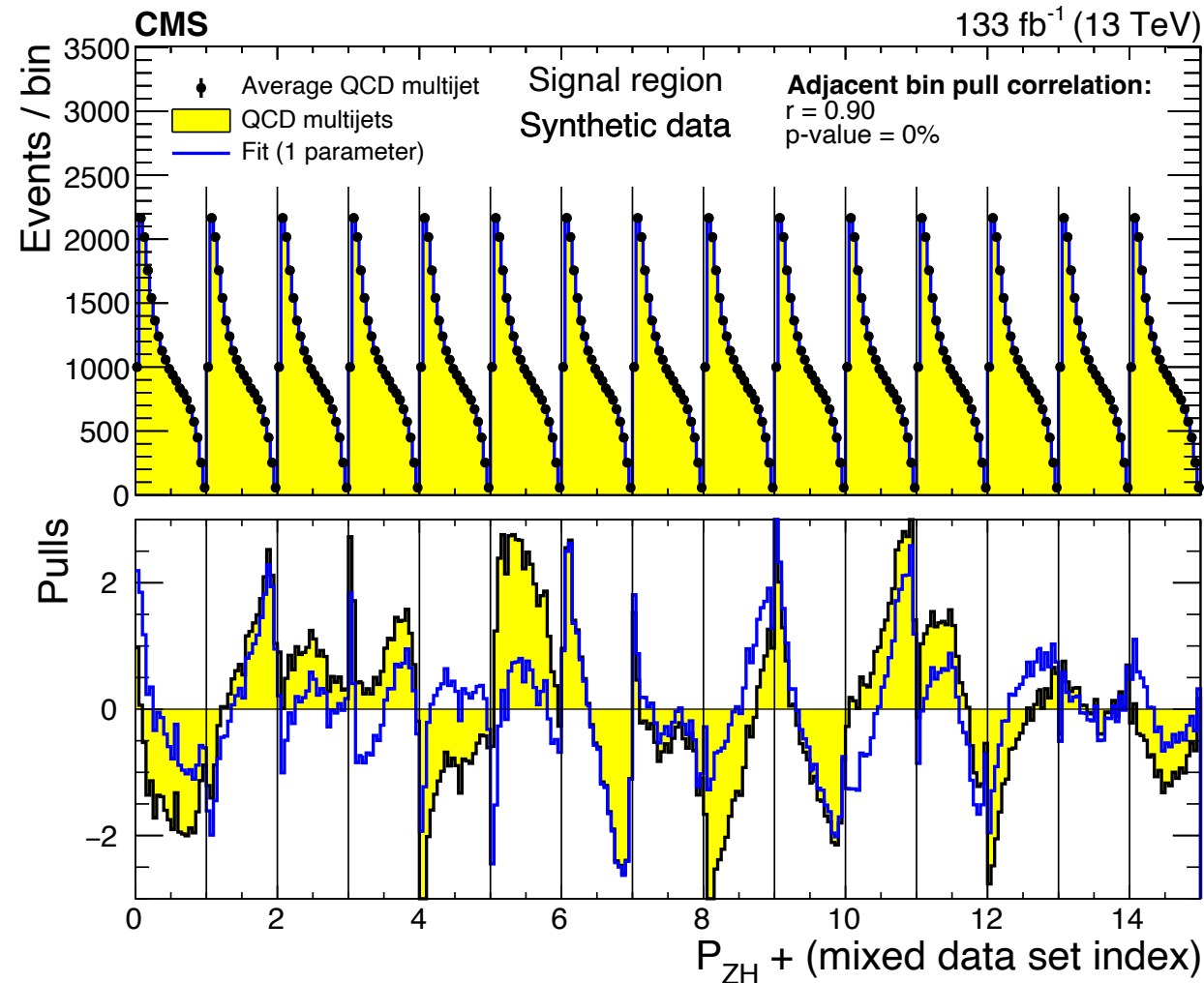
Adjacent bin pull correlation:
r = 0.90
p-value = 0%



$P_{ZH} +$ (mixed data set index)

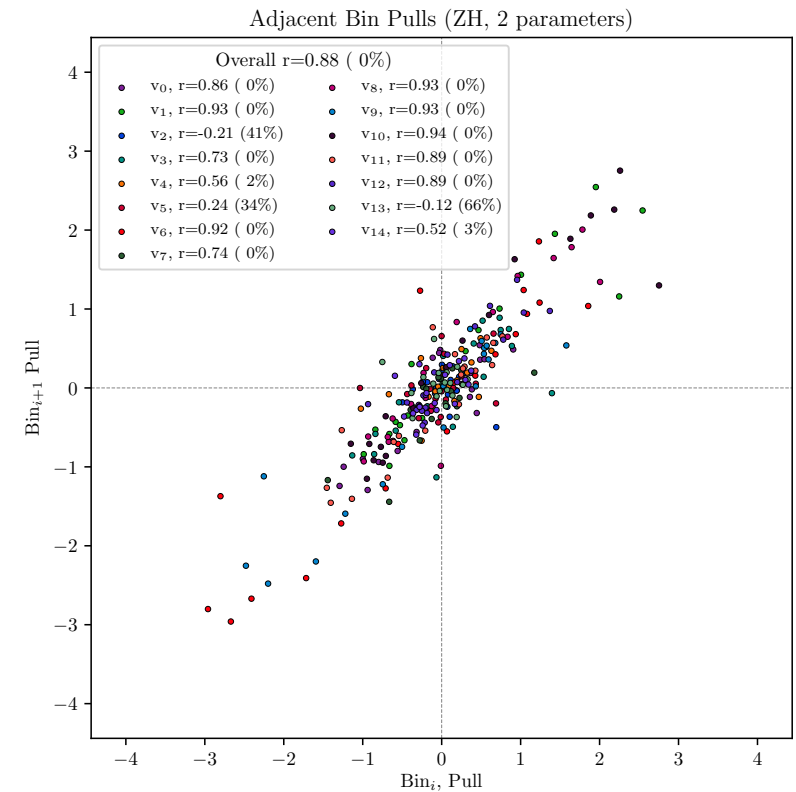
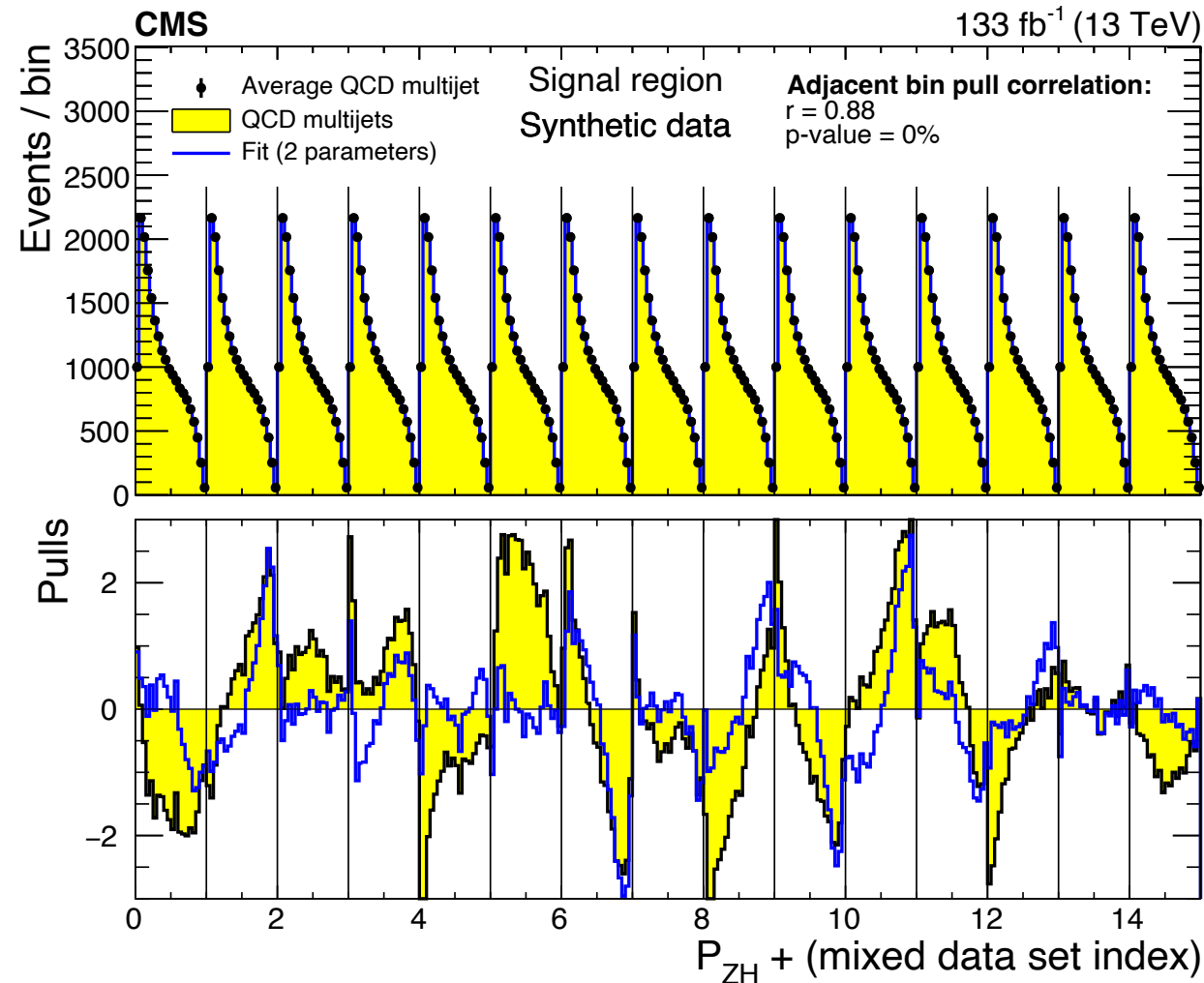
Quantifying Variance

Compare each of the background predictions to the average



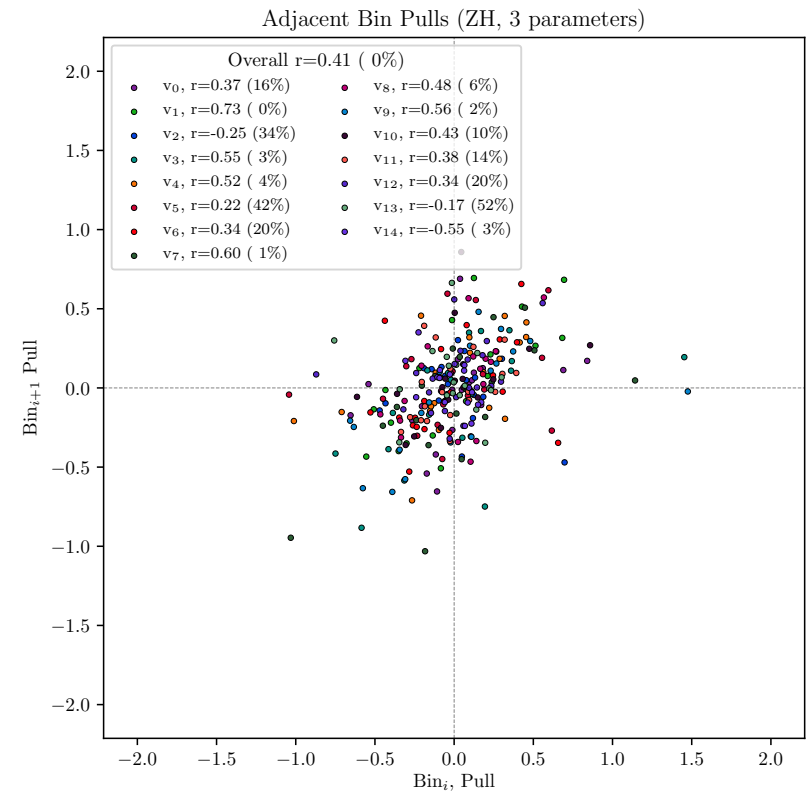
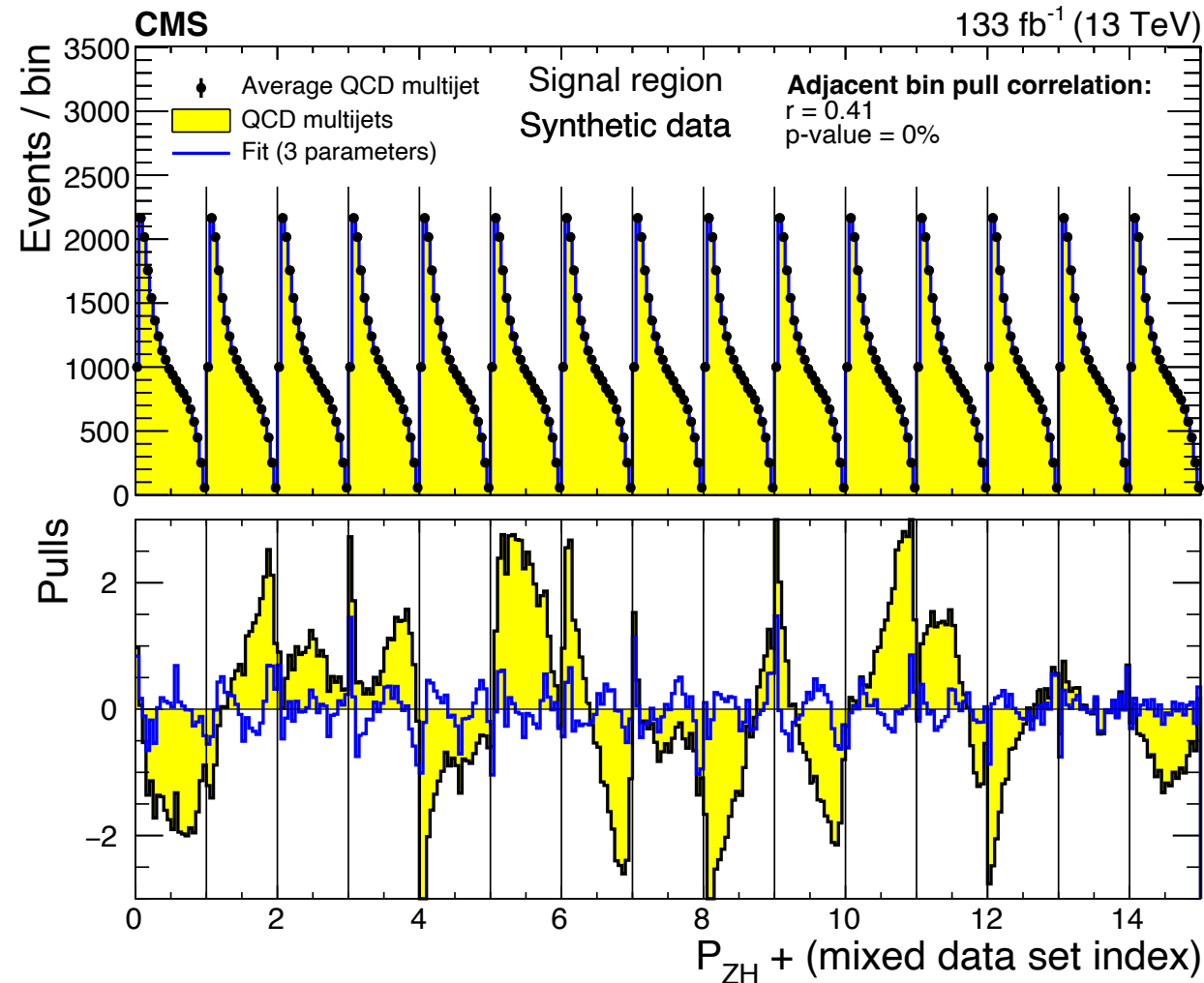
Quantifying Variance

Compare each of the background predictions to the average



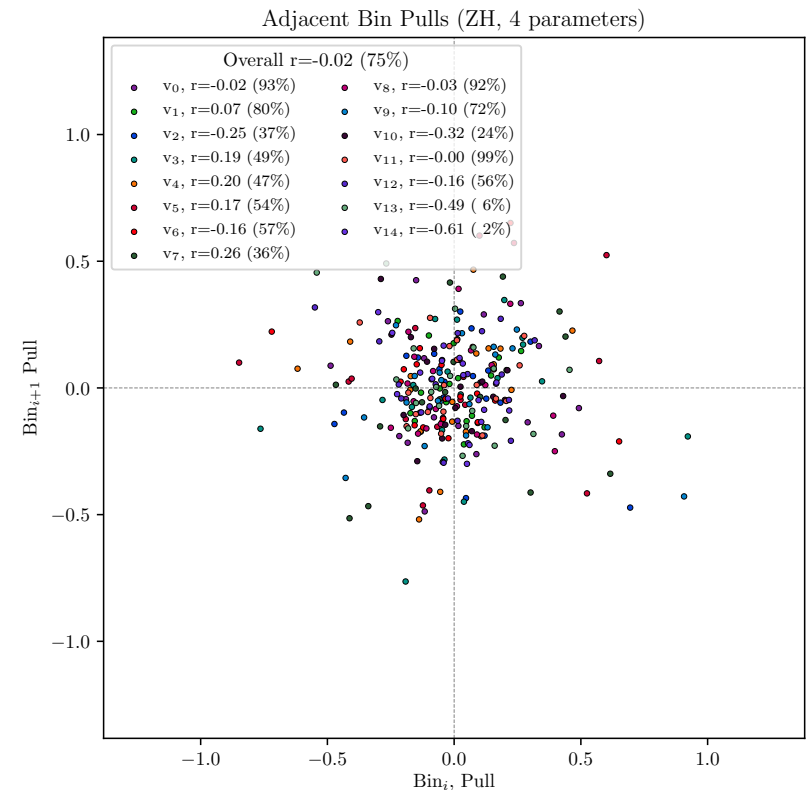
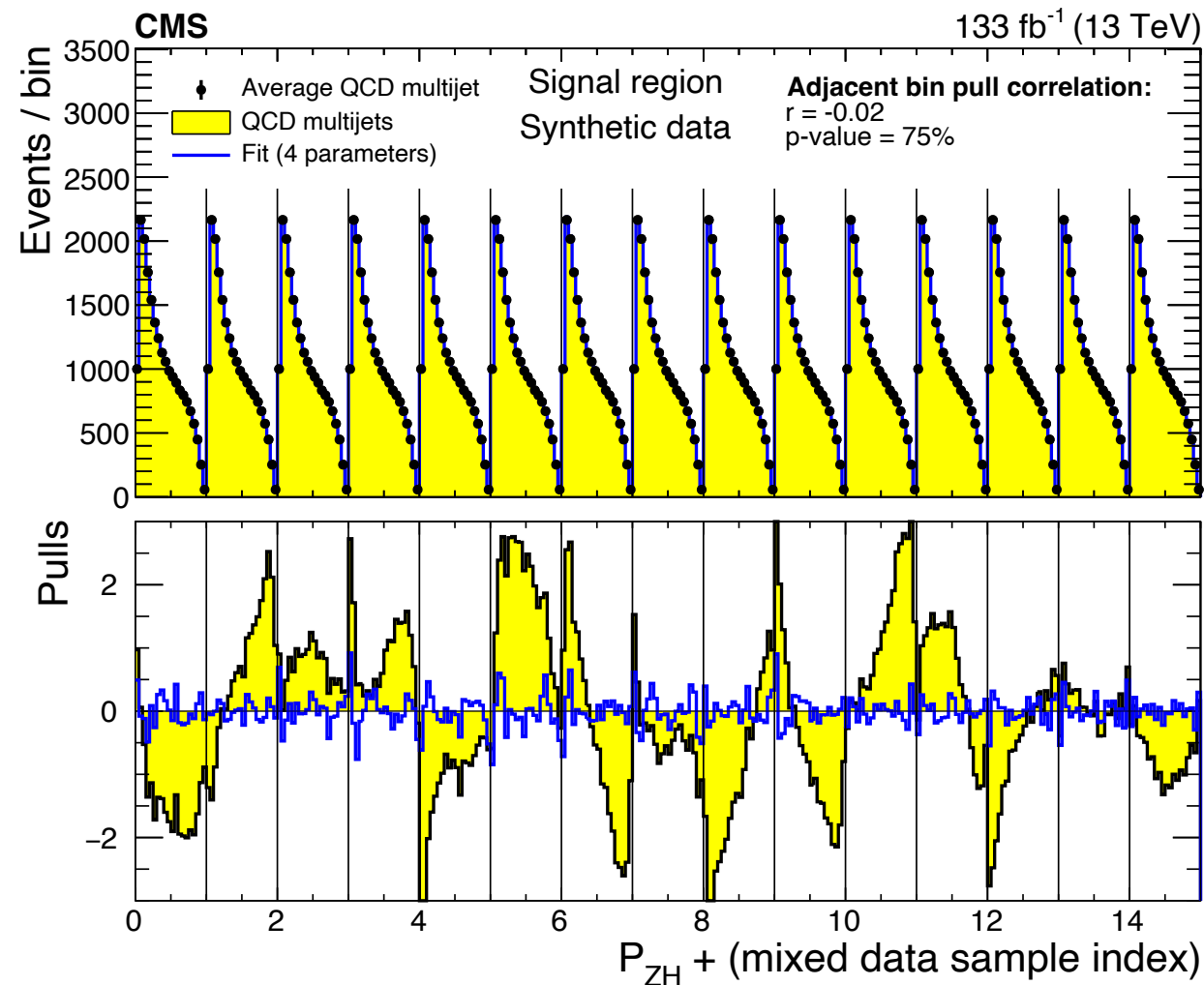
Quantifying Variance

Compare each of the background predictions to the average



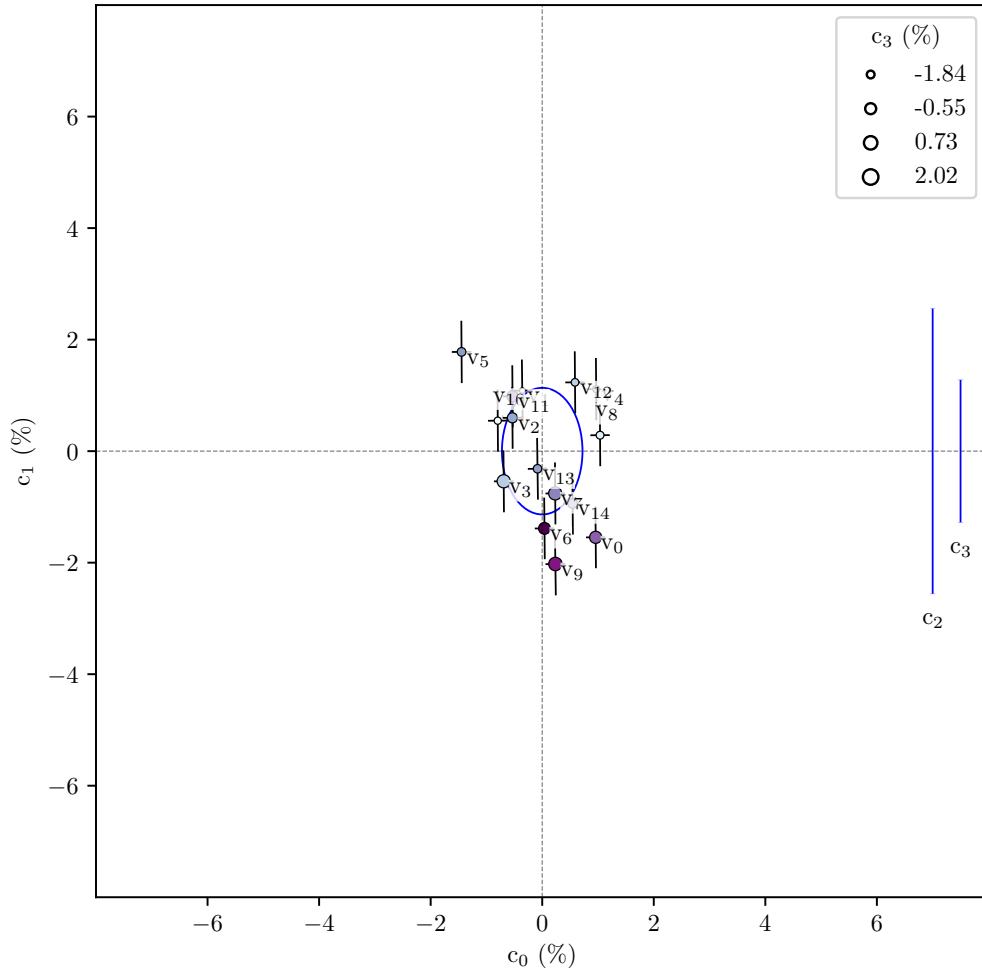
Quantifying Variance

Compare each of the background predictions to the average

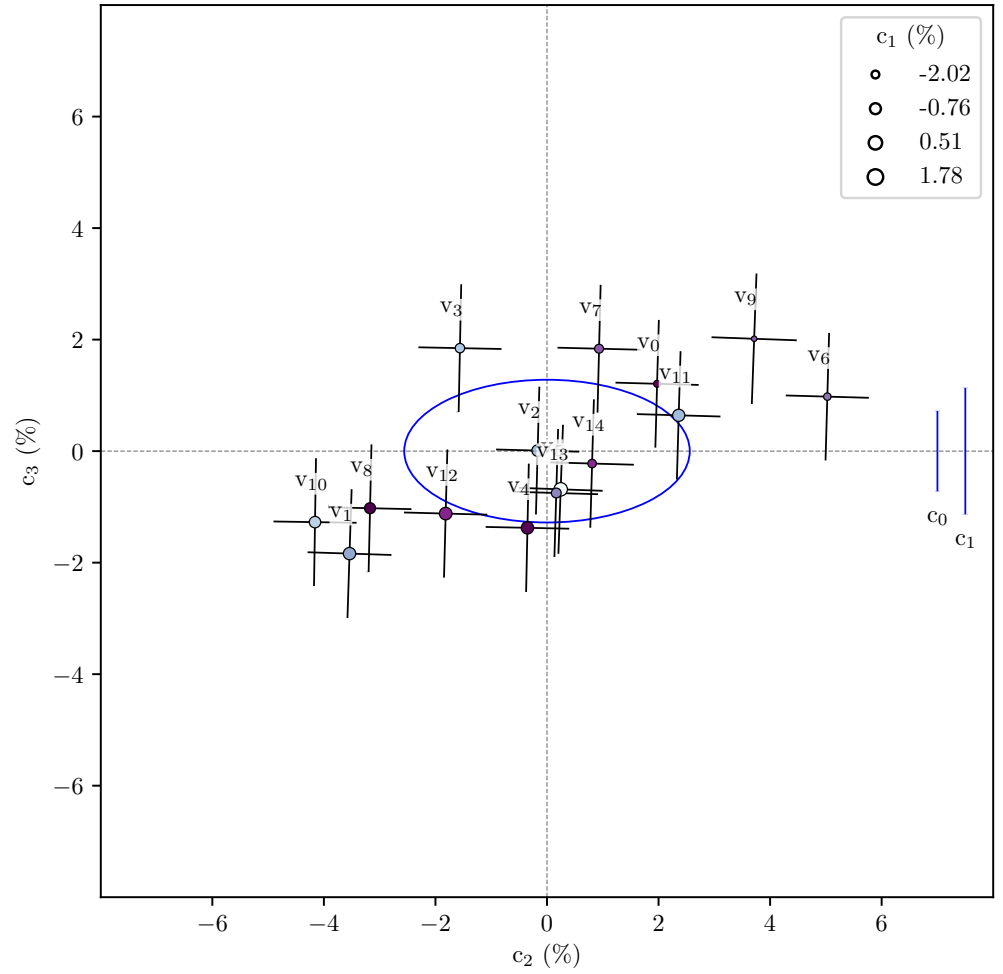


Variance Uncertainties

Multijet Model Variance Fits (ZH)

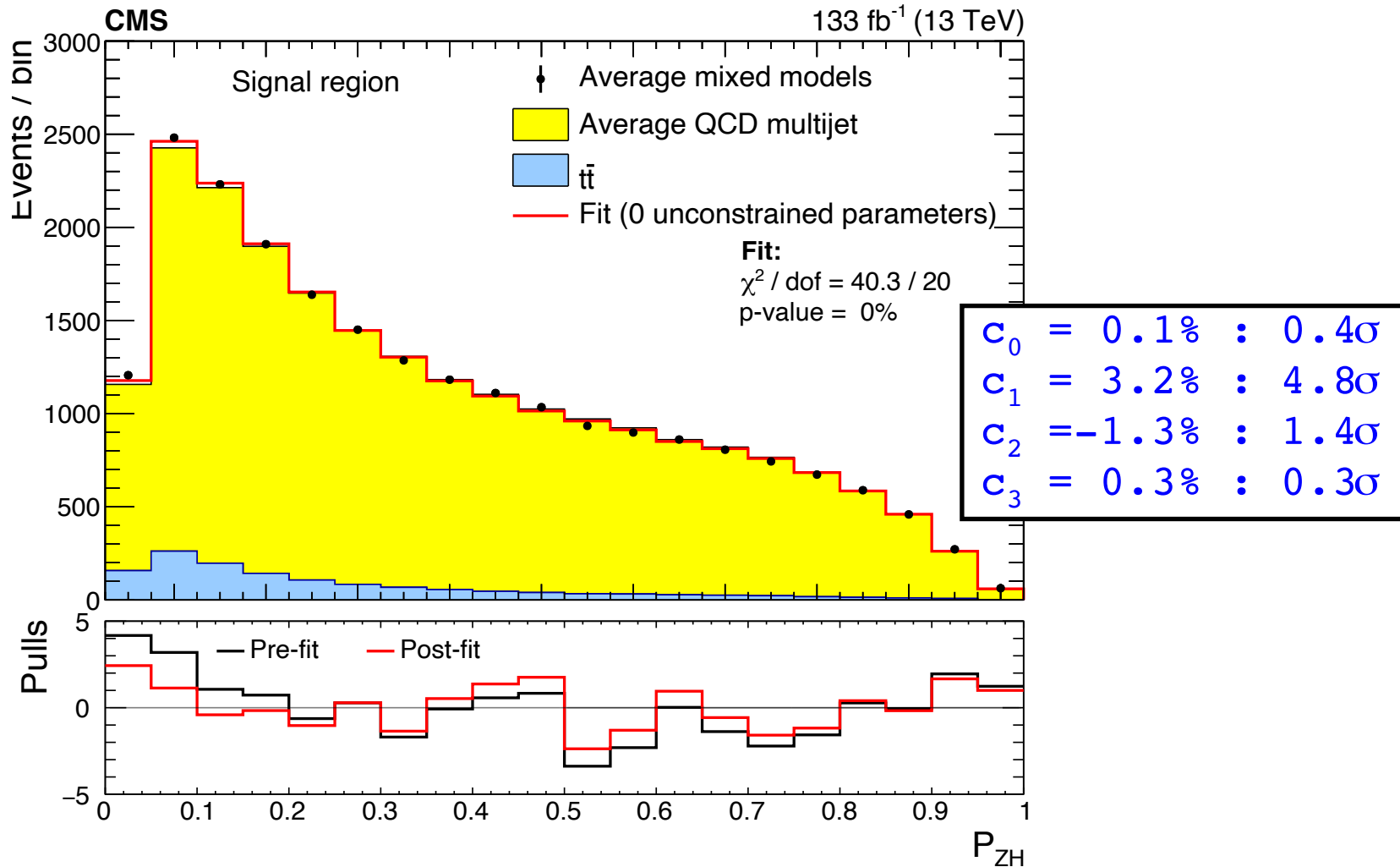


Multijet Model Variance Fits (ZH)



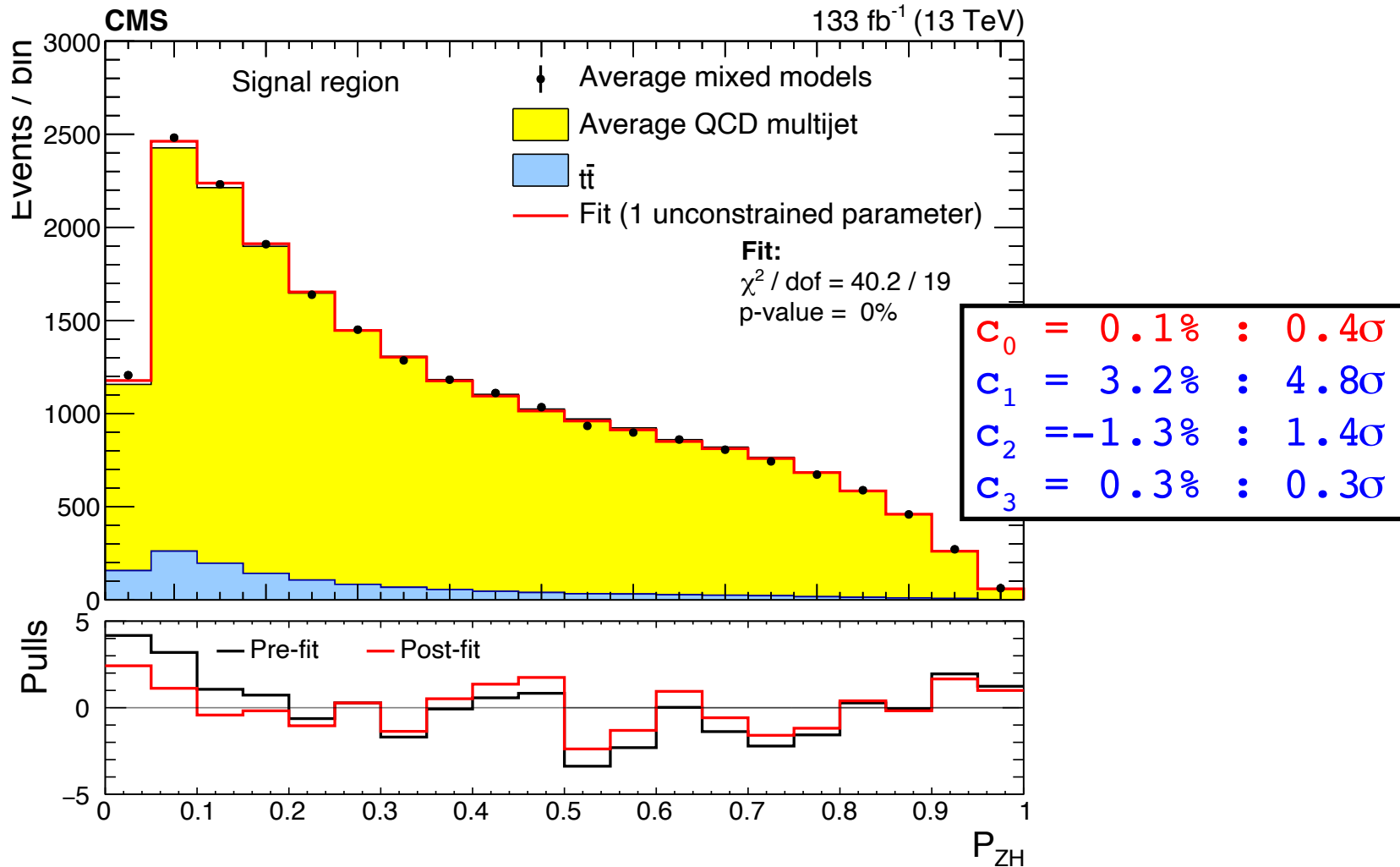
Extrapolation Uncertainty

Compare average background predictions to observed yield in (mixed-data) signal region



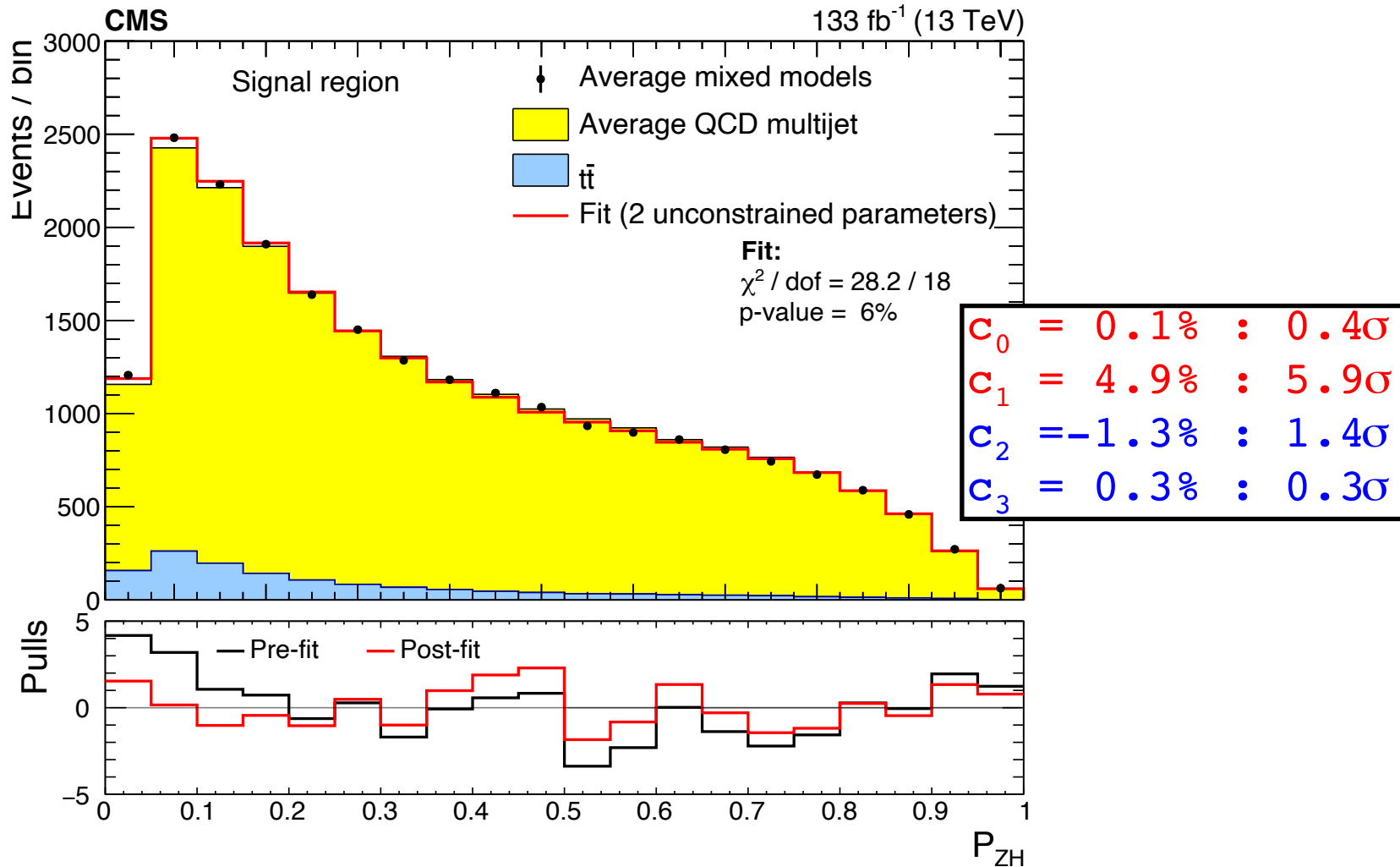
Extrapolation Uncertainty

Compare average background predictions to observed yield in (mixed-data) signal region



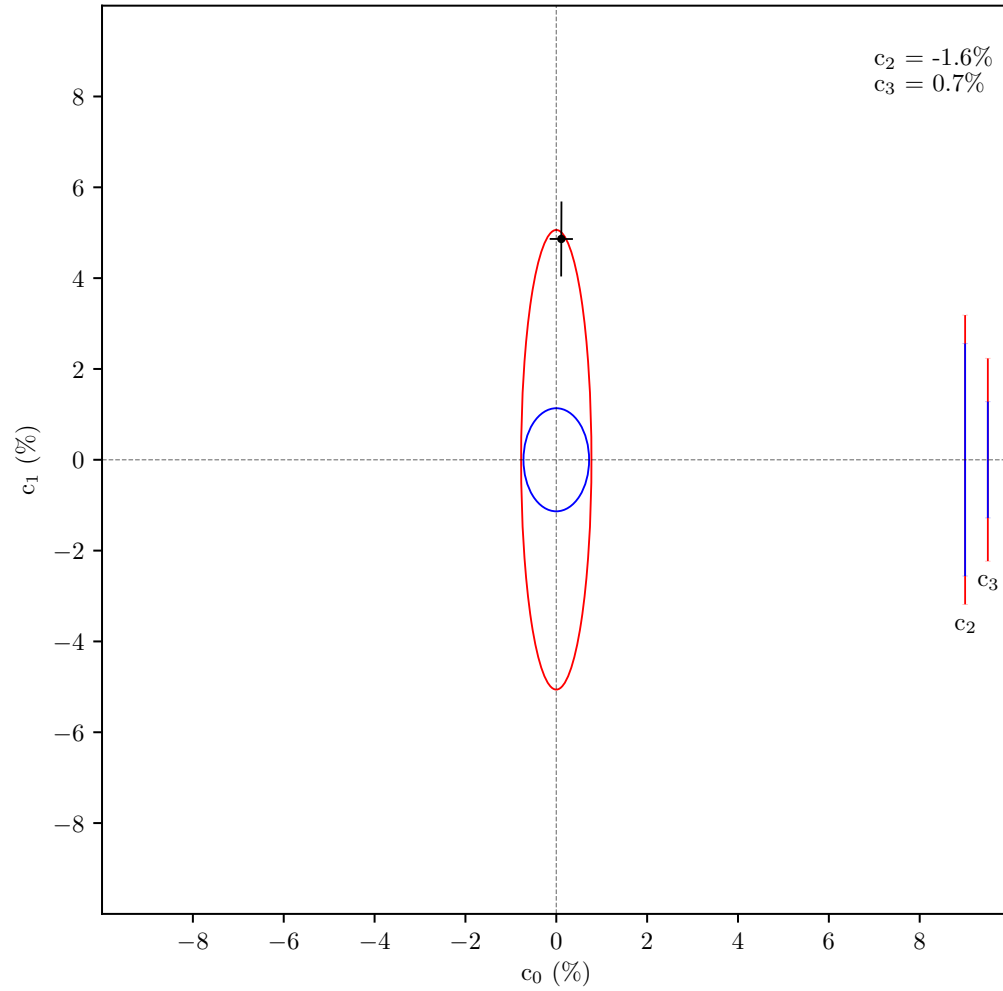
Extrapolation Uncertainty

Compare average background predictions to observed yield in (mixed-data) signal region

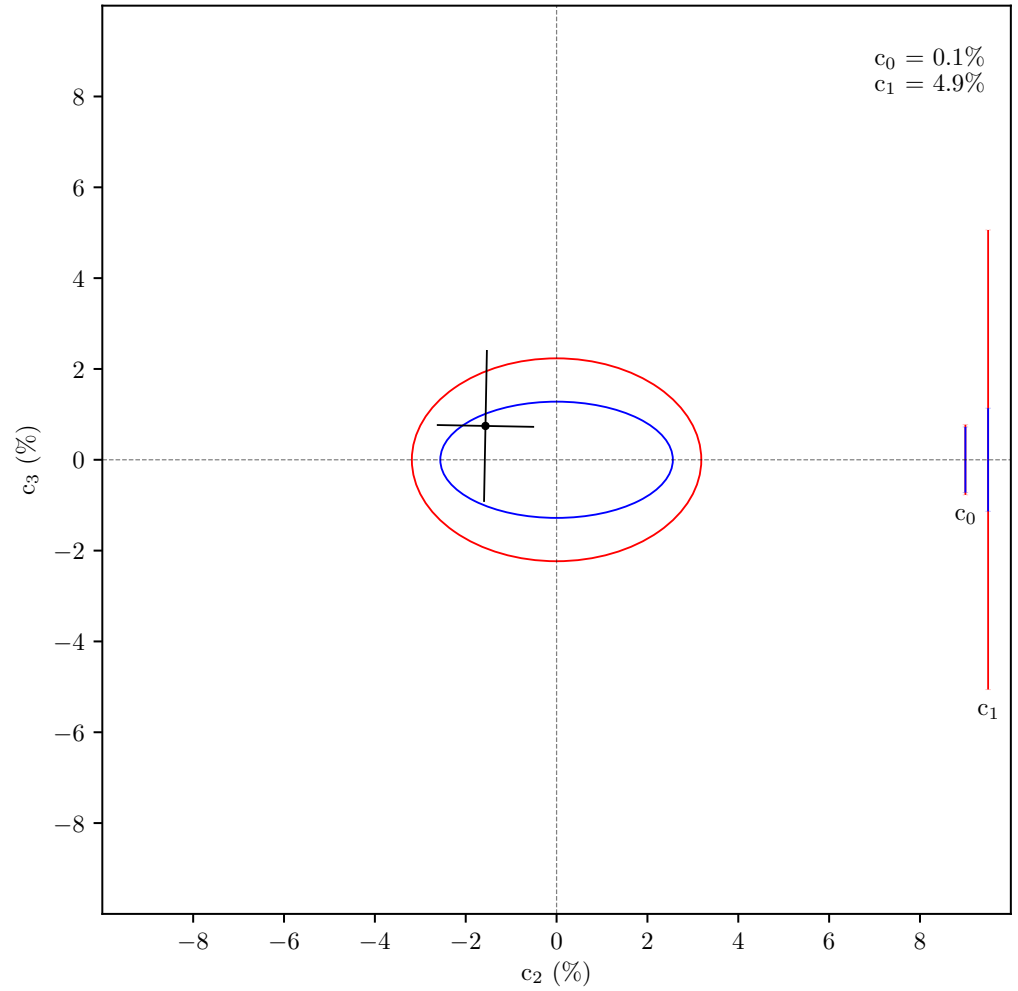


Bias Uncertainties

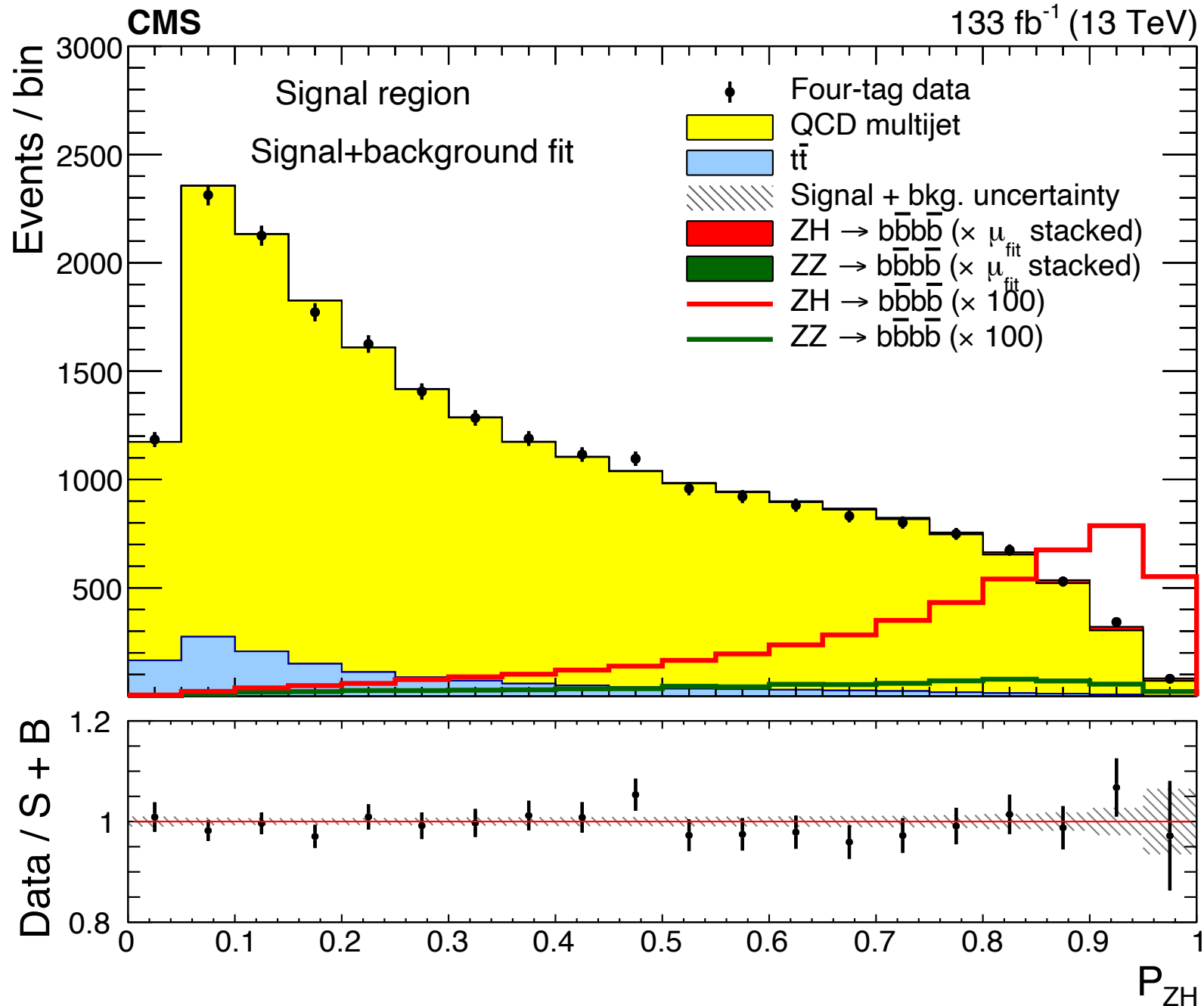
Multijet Model Bias Fit (ZH)



Multijet Model Bias Fit (ZH)



Analysis of 4b Signal Region



Conclusions

Data-driven background ubiquitous in particle physics

Require assumptions w/large hard-to-quantify systematic uncertainties

Synthetic datasets can provide more principled assessment of systematics

Believe synthetic datasets will be increasing important in future

Case study in search for $HH \rightarrow 4b$ more details: [arXiv:2403.20241](https://arxiv.org/abs/2403.20241)

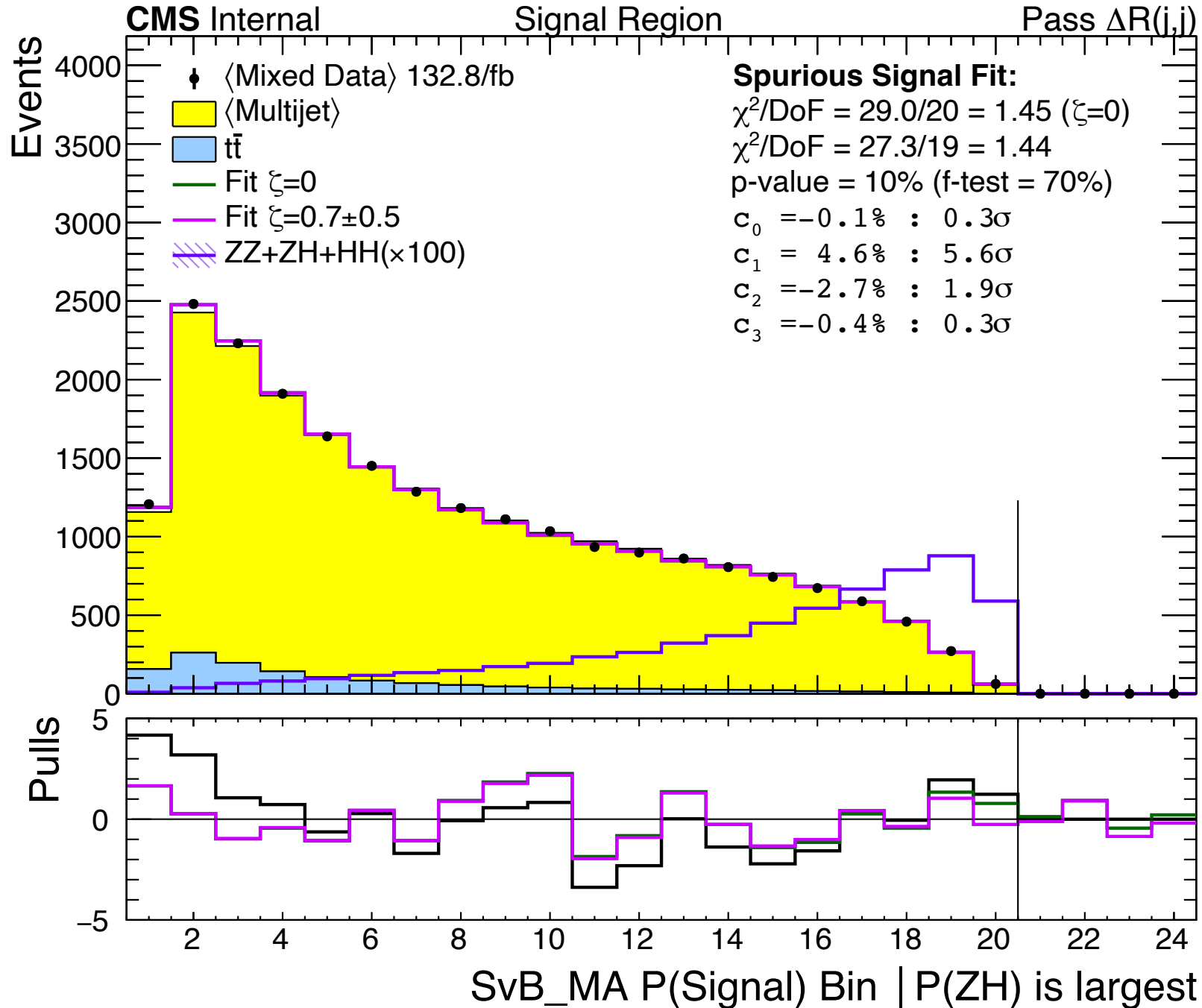
Believe concept can be generalized beyond HH and high-energy physics

Future directions:

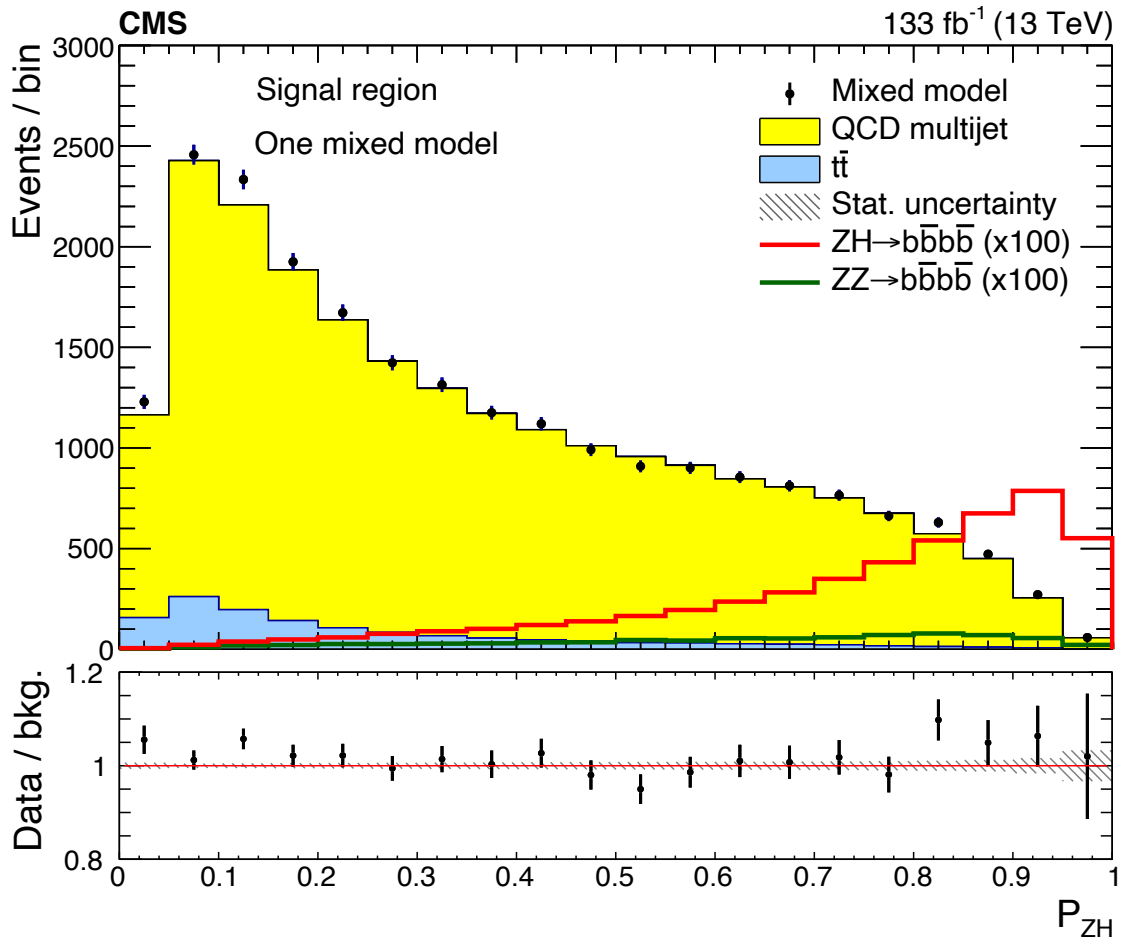
- Reduce variance by k-folding
- Correct bias, take smaller uncertainty
- Larger higher fidelity synthetic datasets

Backup

Spurious Signal



Mixed data:



Mixed data:

