

# Boosted $ggH \rightarrow cc$ @CMS

YSF: Search for boosted Higgs boson decays  
to charm quark pairs

**Andrzej Novak** on behalf of the CMS Collaboration



European Research Council  
Established by the European Commission



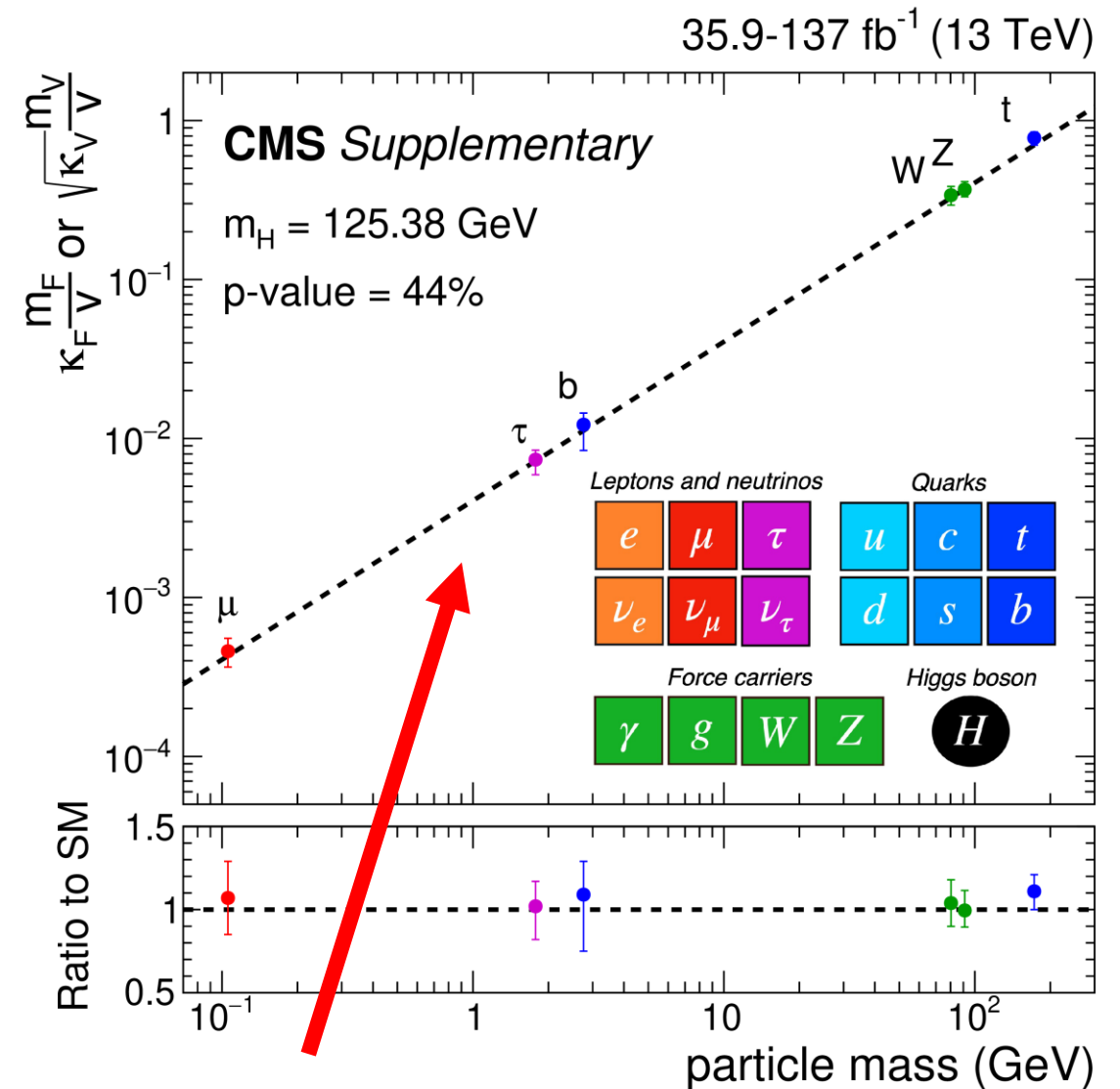
# Motivation

## Complete description of the Higgs boson

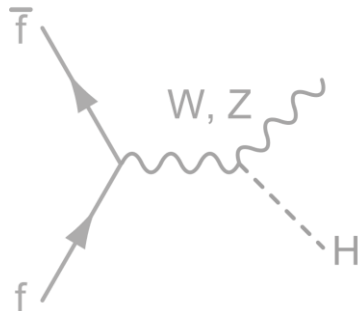
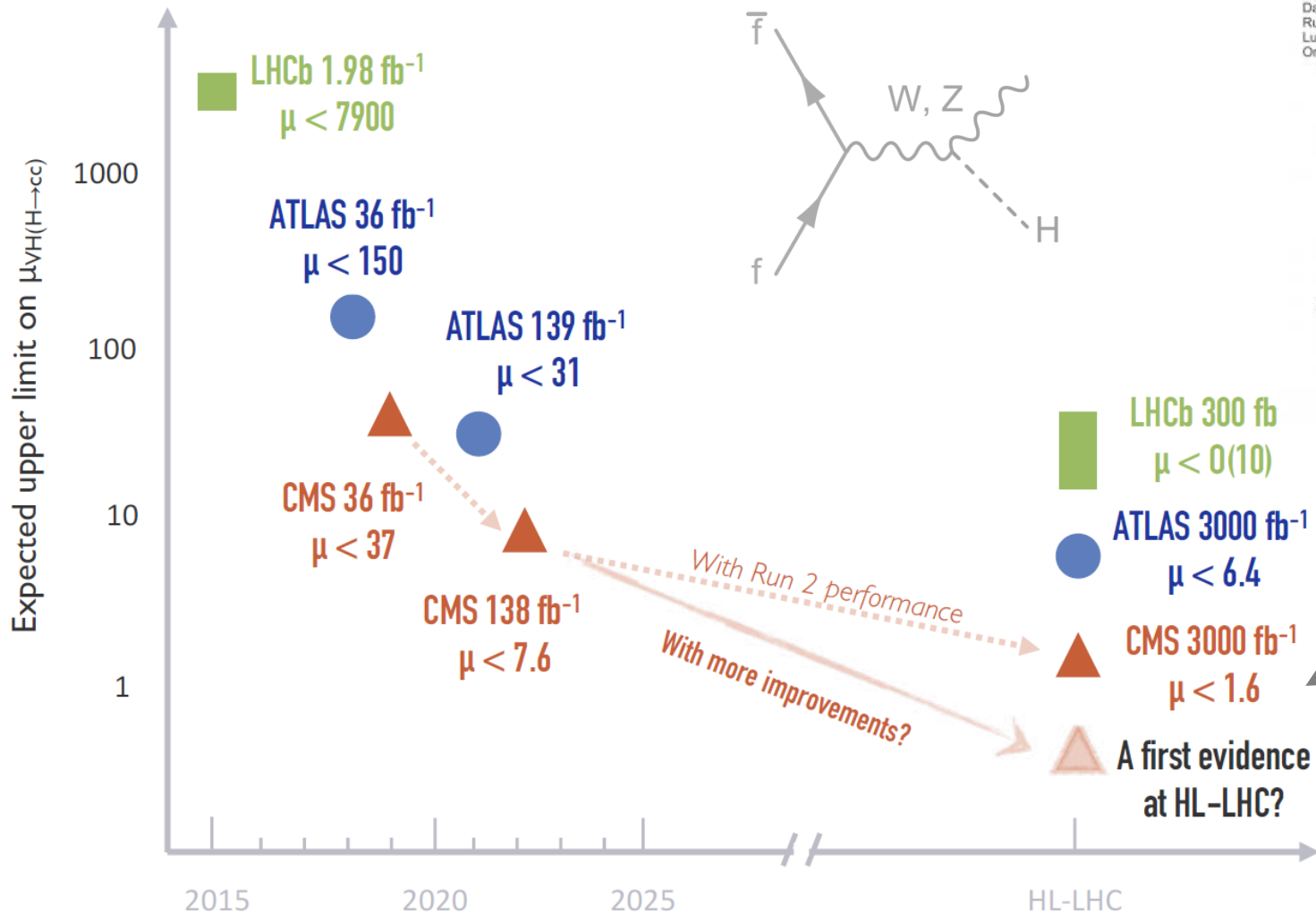
- Observed couplings ( $\sim 90\%$  of BR)
  - Vector bosons (WW, ZZ)
  - 3<sup>rd</sup> generation fermions (ttH, bb,  $\tau\tau$ )
- Evidence for 2<sup>nd</sup> generation couplings
  - $H \rightarrow \mu\mu$

## Next target – search for $H \rightarrow cc$

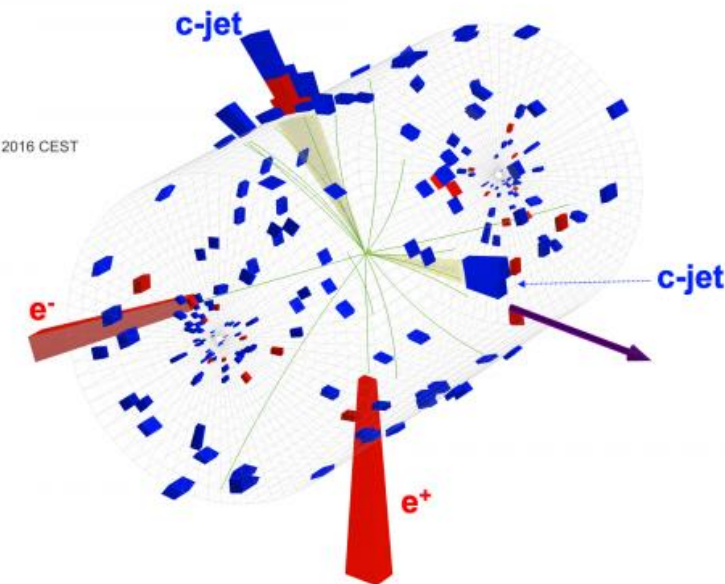
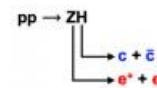
- Largest missing fraction of BR
- Establishing couplings to 2<sup>nd</sup> gen. quarks
- Measuring any enhancements to coupling would hint at new physics



# VH → cc Searches



CMS Experiment at LHC, CERN  
Data recorded: Tue May 31 11:26:24 2016 CEST  
Run/Event: 274250 / 1058807020  
Lumi section: 543  
Orbit/Crossing: 142305803 / 593



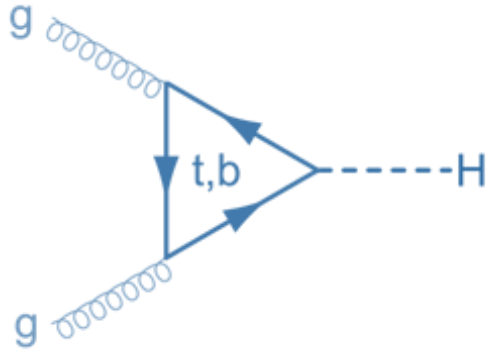
Not quite good enough

Need to improve sensitivity and/or combine with other channels

H. Qu & L. Mastrolorenzo, LHC Seminar 2022

# Direct Search Channels

## Gluon Fusion



### Premise:

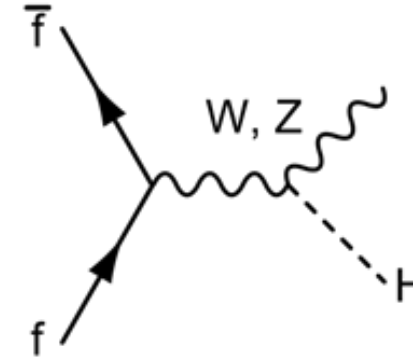
Largest Higgs production cross-section

- $\sigma(pp \rightarrow H)@13 \text{ TeV} = 48.6 \text{ pb}$

Very large QCD contamination:  $g \rightarrow qq$

- $\sigma(pp \rightarrow qq)@13 \text{ TeV} \sim 360\,000 \text{ pb}$  (HT > 300 GeV)

## Higgs-Strahlung



### Premise:

Lower cross-section

- $\sigma(pp \rightarrow VH)@13\text{TeV} = 2.25 \text{ pb}$

Clean signature: trigger on leptons from Z,W

Mostly V+jets background

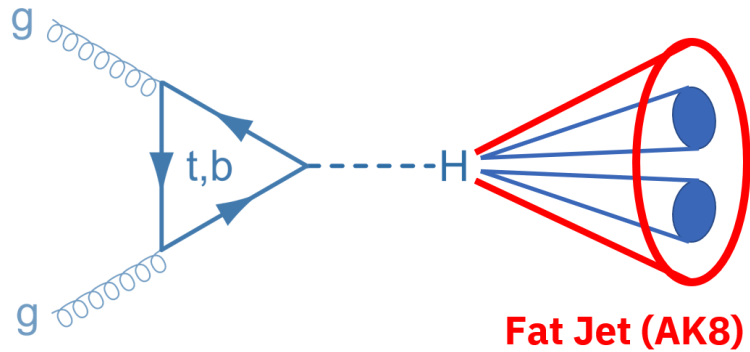
- $\sigma(pp \rightarrow V+qq)@13\text{TeV} \sim 600 \text{ pb}$  (HT > 400 GeV)

x 20

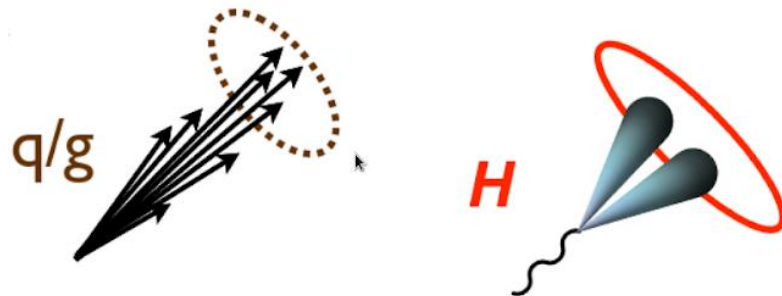
x 600

# General Strategy

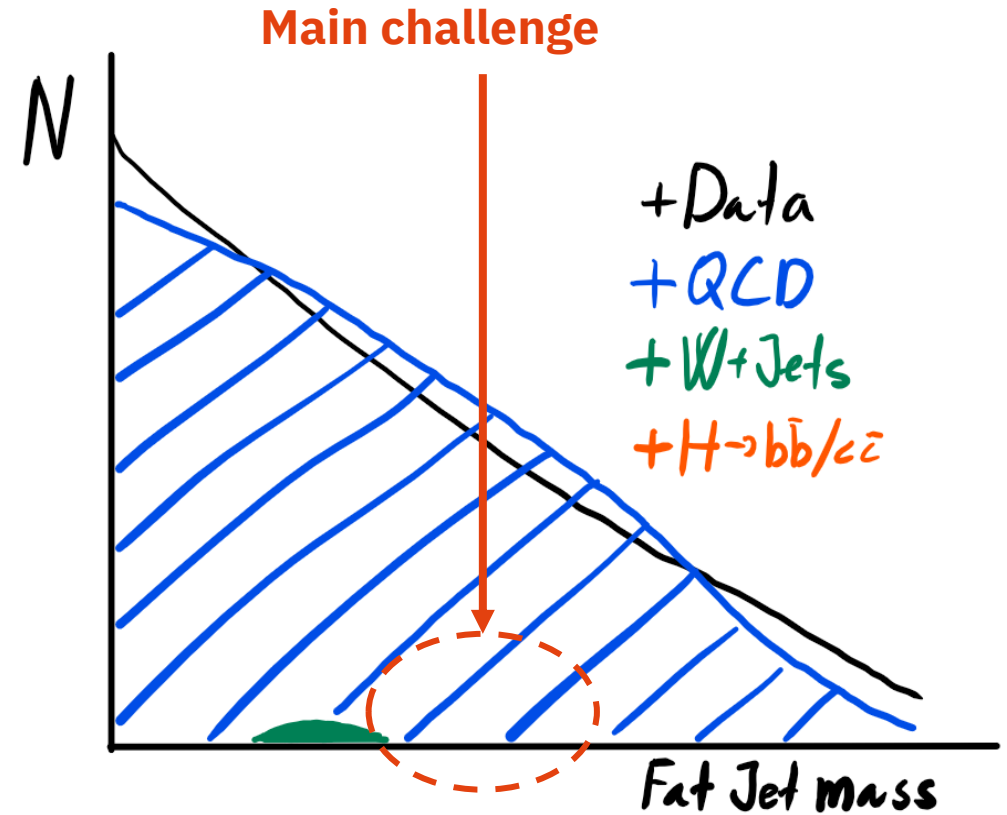
## Gluon Fusion Production Mode



## Jet Identification



## Data-driven fit in jet mass



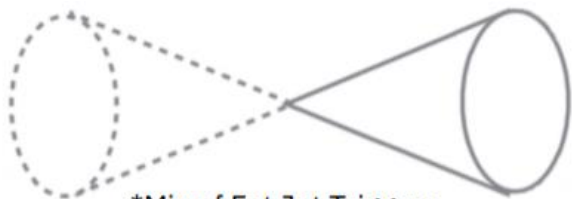
# Analysis In a Nutshell

## Event Selection

Lepton veto

Top veto

- $E_T^{miss} < 140$
- b-tag veto



\*Mix of Fat Jet Triggers

AK8 Jet (Higgs Cand.)

- $m_{SD}, p_T$  cuts
- $N_2$  substructure cut
- Charm vs. bottom cut

## Signal Region Definition

Classifier

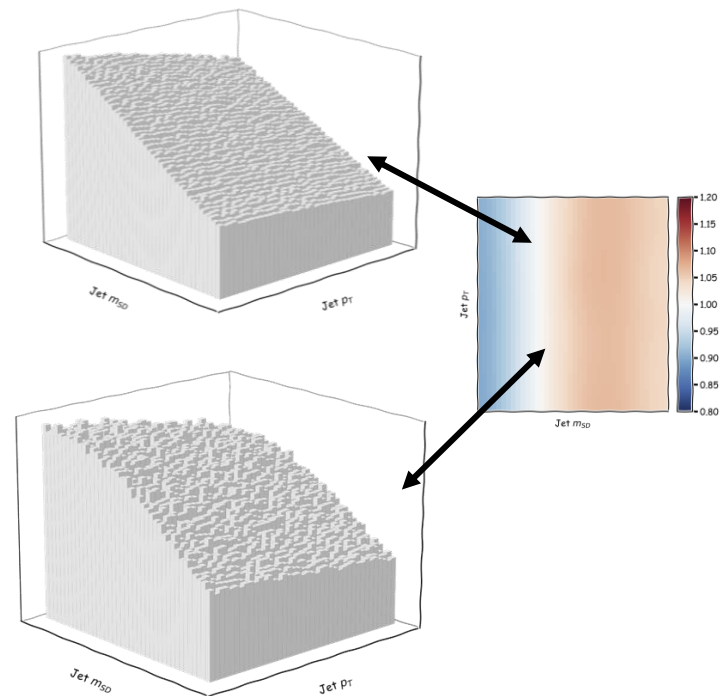
Signal Rich  
(pass)

cut

Signal Depleted  
(fail)

Preselection

## Simultaneous Fit



Fit jet soft-drop mass in bins of  $p_T$

Constrain QCD from data *in-situ*

# Analysis In a Nutshell

## Event Selection

- Lepton veto  
Top veto
- $E_T^{miss} < 140$
  - b-tag veto

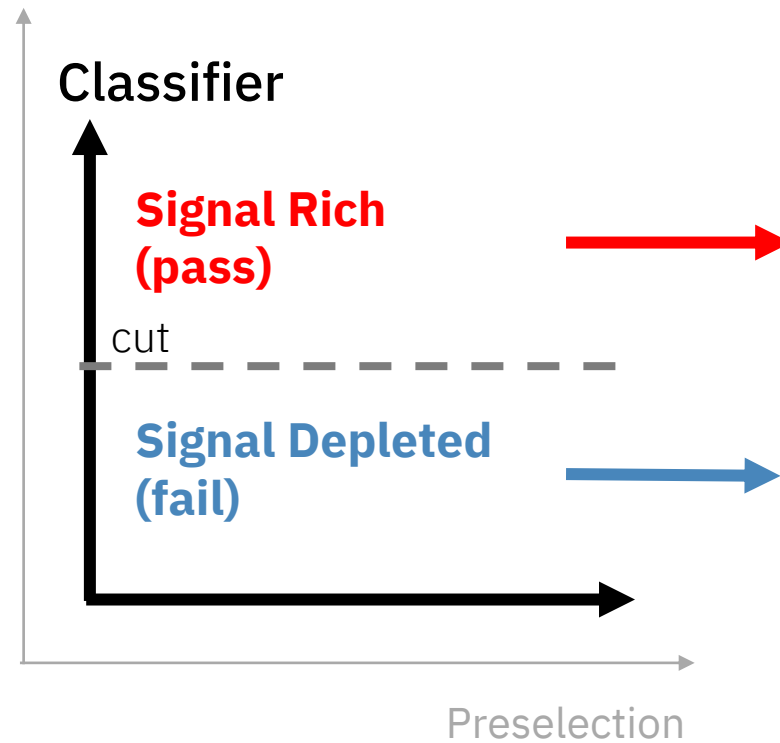


\*Mix of Fat Jet Triggers

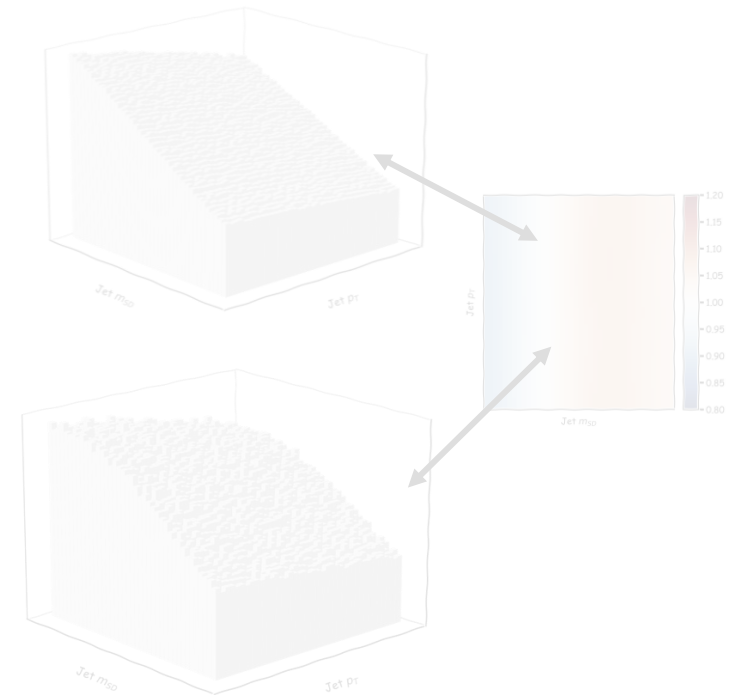
### AK8 Jet (Higgs Cand.)

- $m_{SD}, p_T$  cuts
- $N_2$  substructure cut
- Charm vs. bottom cut

## Signal Region Definition



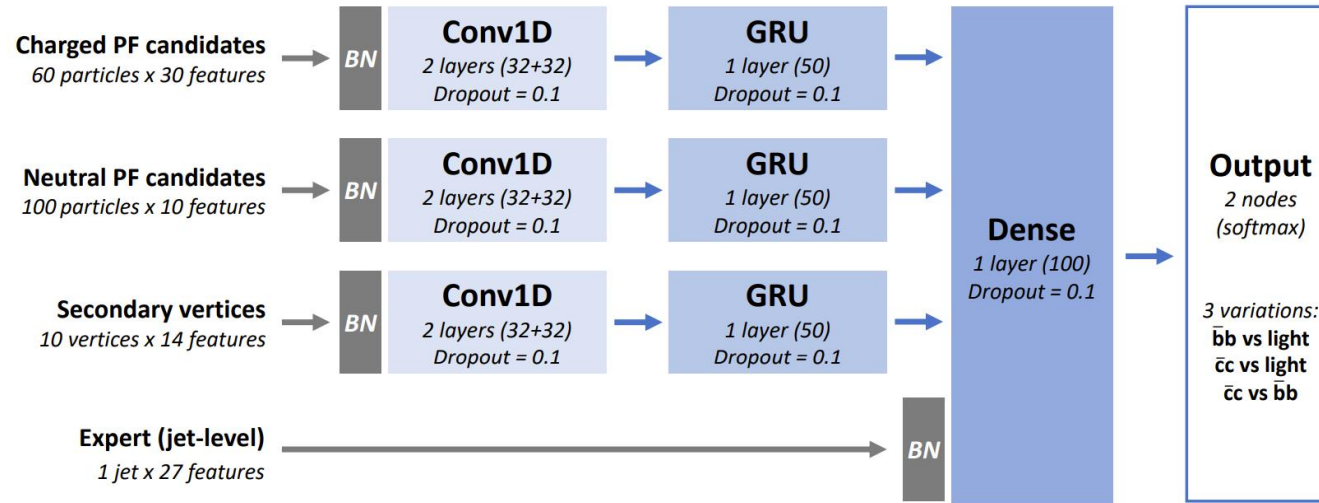
## Simultaneous Fit



Fit jet soft-drop mass in bins of  $p_T$   
Constrain QCD from data *in-situ*

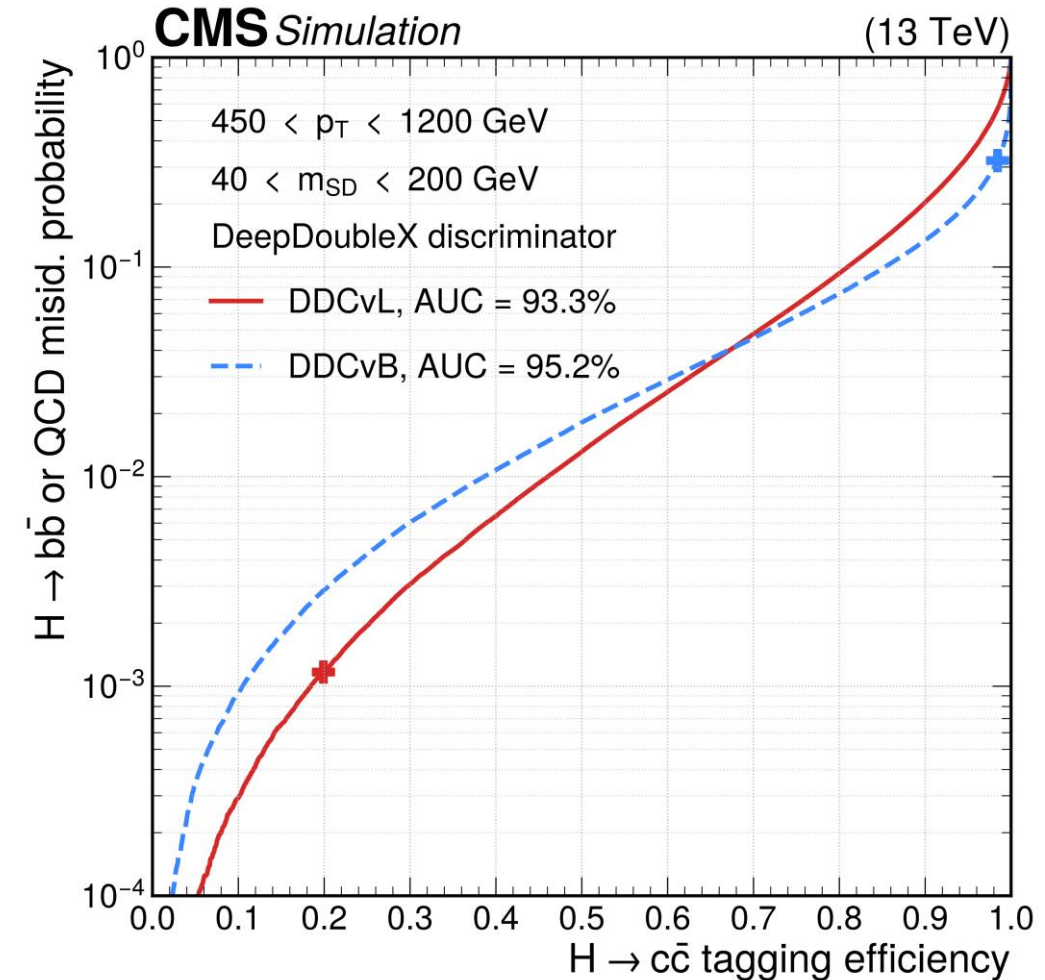
# Charm Tagging – DeepDoubleX

DeepJet (CNN) inspired architecture with extended range of low-level inputs



Input relevance studies (IG, DeepTaylor)

- ~50% reduction in input features, no loss in perf.
- Much faster inference





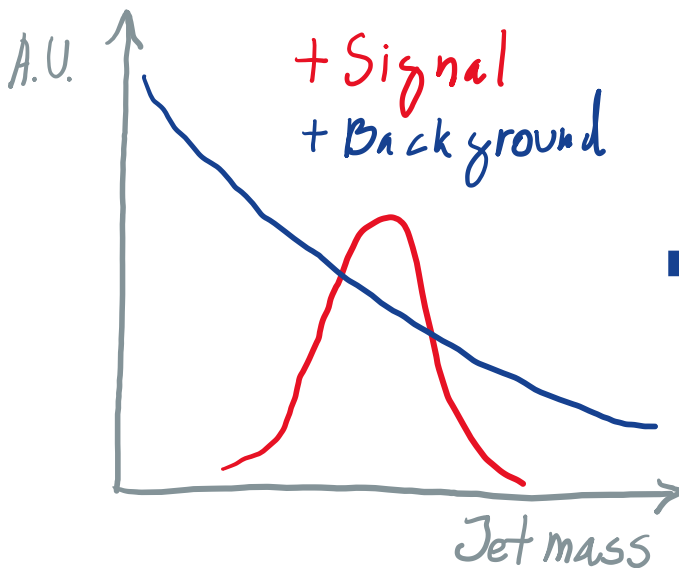
# Charm Tagging – Mass Decorrelation

Independence of the fitted variable on the selection is crucial

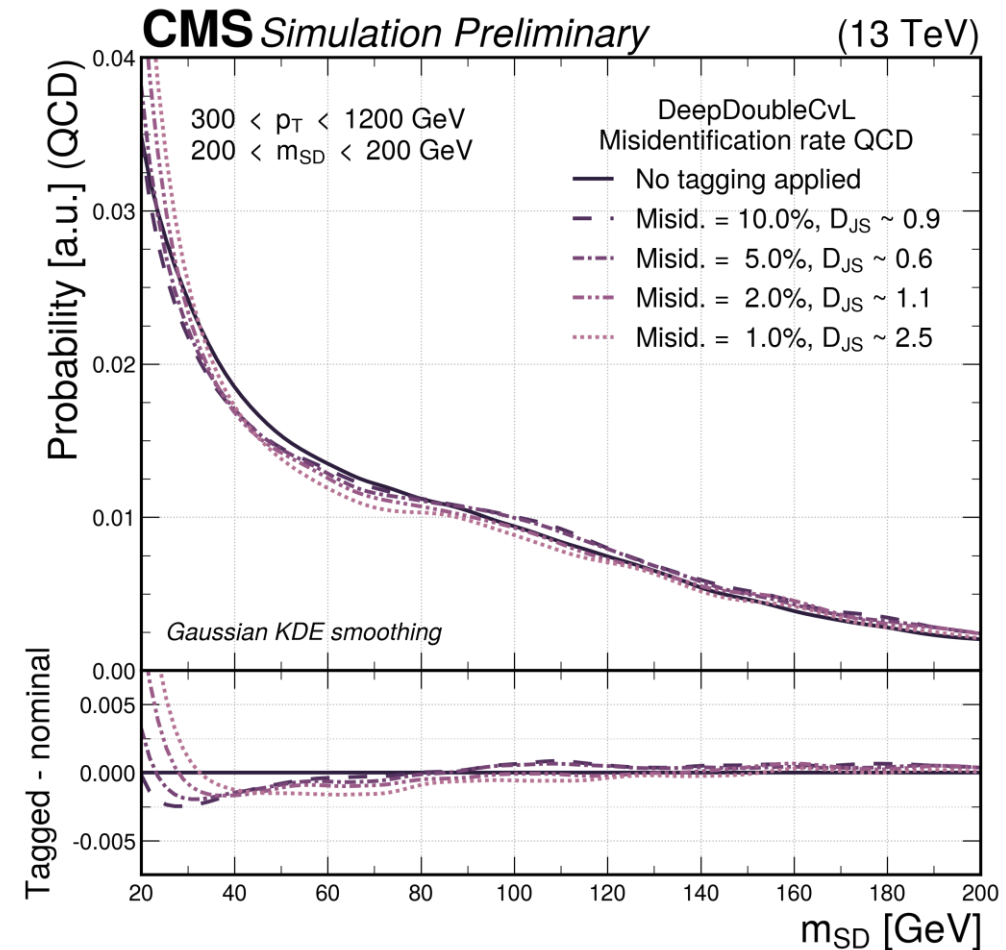
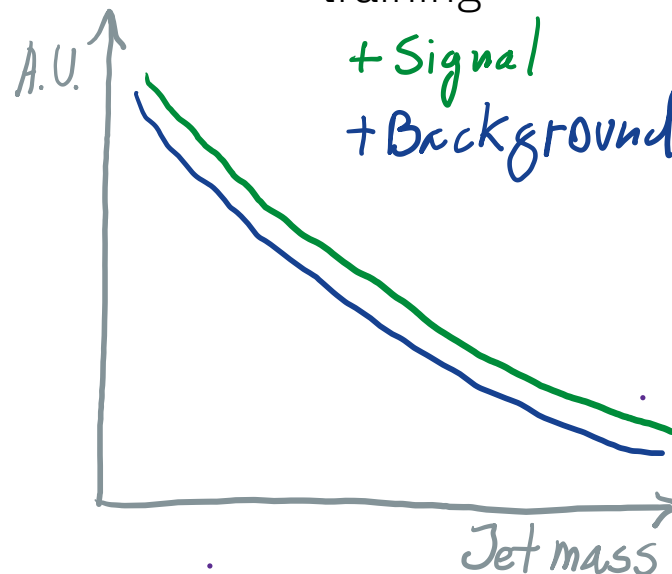
- In particular for data-driven background estimation

Train classifiers with flat mass spectrum signal

Plain training: no mass decorrelation



Mass-decorrelated training



# Analysis In a Nutshell

## Event Selection

- Lepton veto  
Top veto
- $E_T^{miss} < 140$
  - b-tag veto

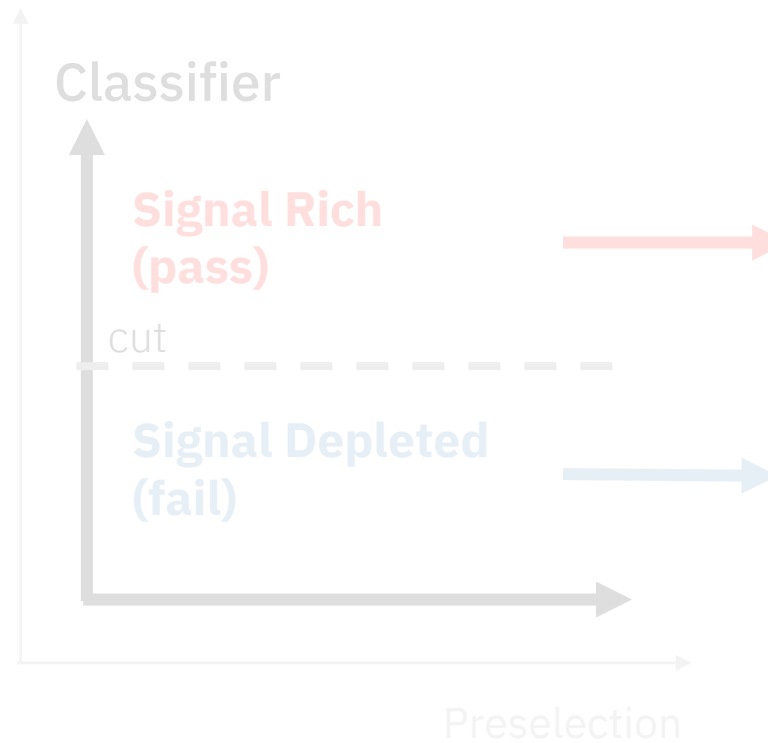


\*Mix of Fat Jet Triggers

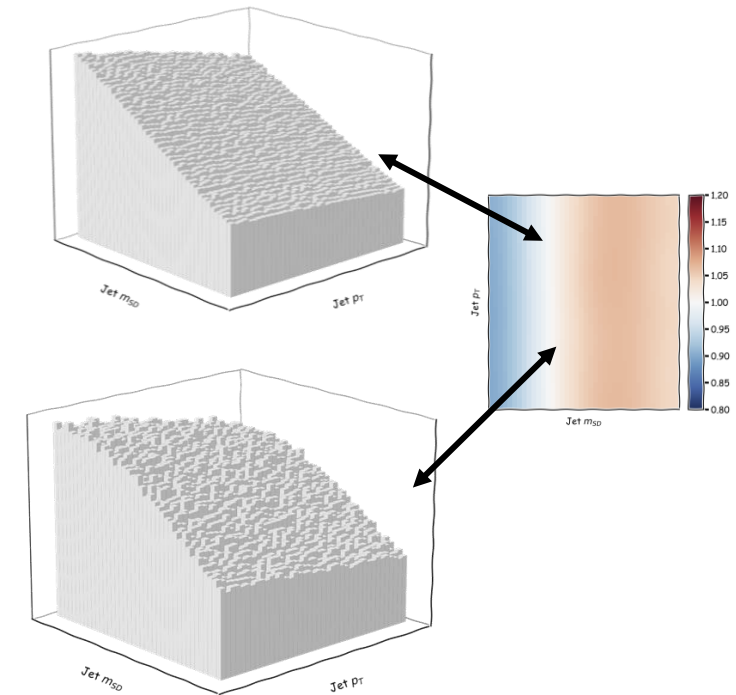
### AK8 Jet (Higgs Cand.)

- $m_{SD}, p_T$  cuts
- $N_2$  substructure cut
- Charm vs. bottom cut

## Signal Region Definition



## Simultaneous Fit



**Fit jet soft-drop mass in bins of  $p_T$**   
**Constrain QCD from data *in-situ***

# Transfer Factor Fit

## Rhalphabet

- Differential **Alphabet** (ABCD)

## Transfer Factor – residual correction

- Accounting for different tagger response
- Flat TF** → **Regular ABCD method**

TF defined (constrained) bin **barycenters**

- 6 pT x 23 ρ (mass) bins
- Surface parametrized in **Bernstein basis**

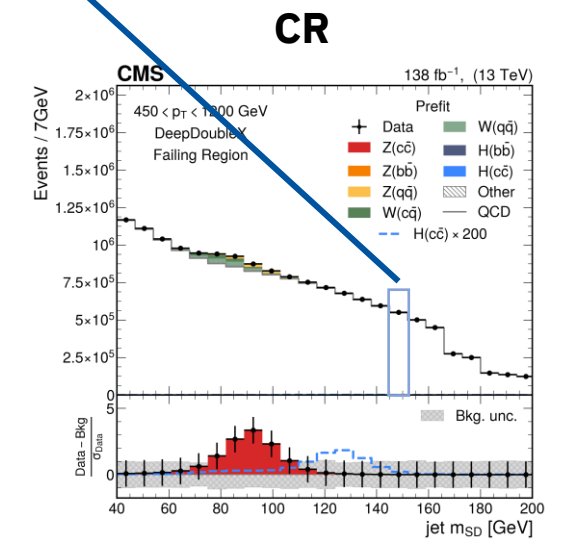
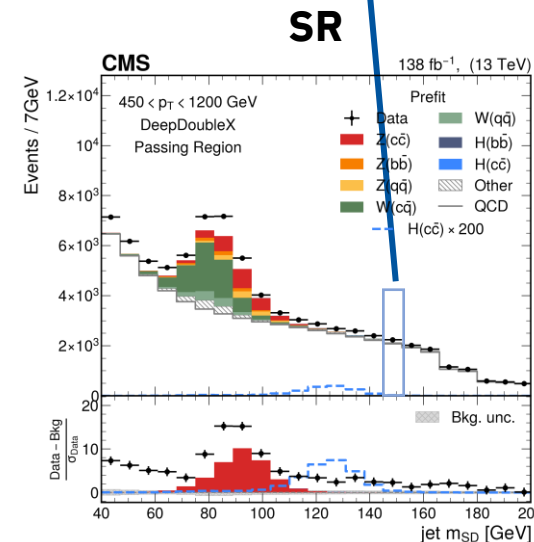
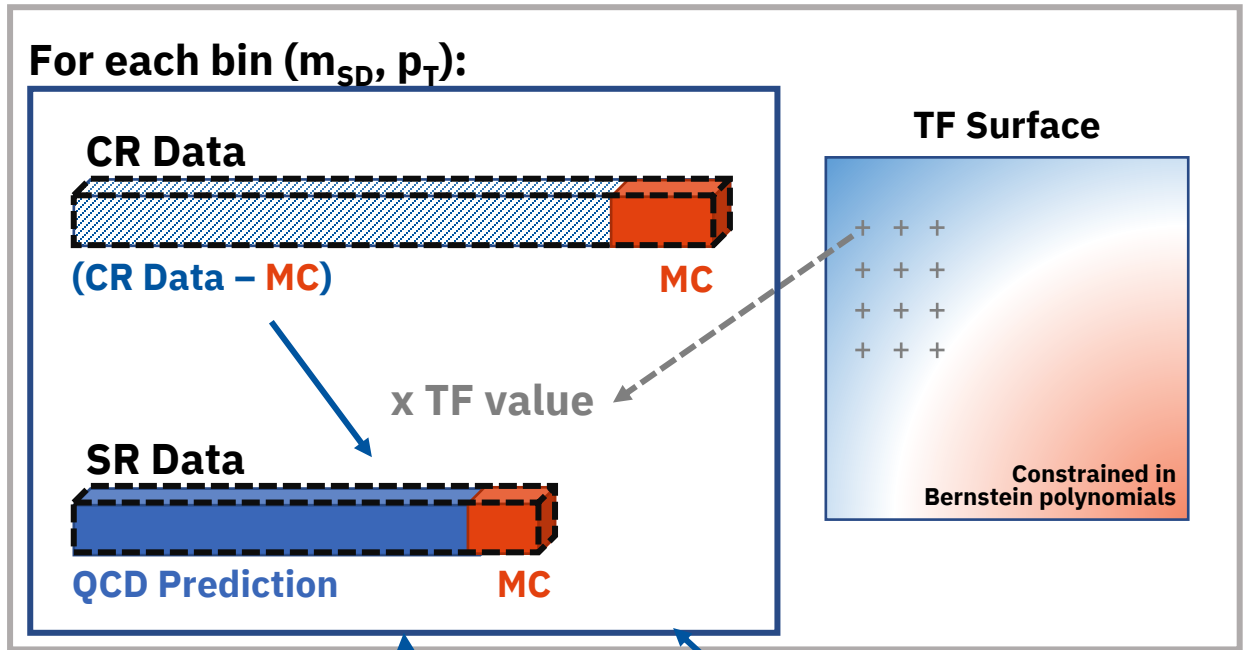
$$R(\rho, p_T) = \sum_{k=0}^{n_\rho} \sum_{\ell=0}^{n_{p_T}} a_{k,\ell} b_{k,n_\rho}(\rho) b_{\ell,n_{p_T}}(p_T) \epsilon^{\text{QCD}}(\rho, p_T)$$

$$b_{\nu,n}(x) = \binom{n}{\nu} x^\nu (1-x)^{n-\nu} \quad \rho = \log \frac{p_T^2}{m_{SD}^2}$$

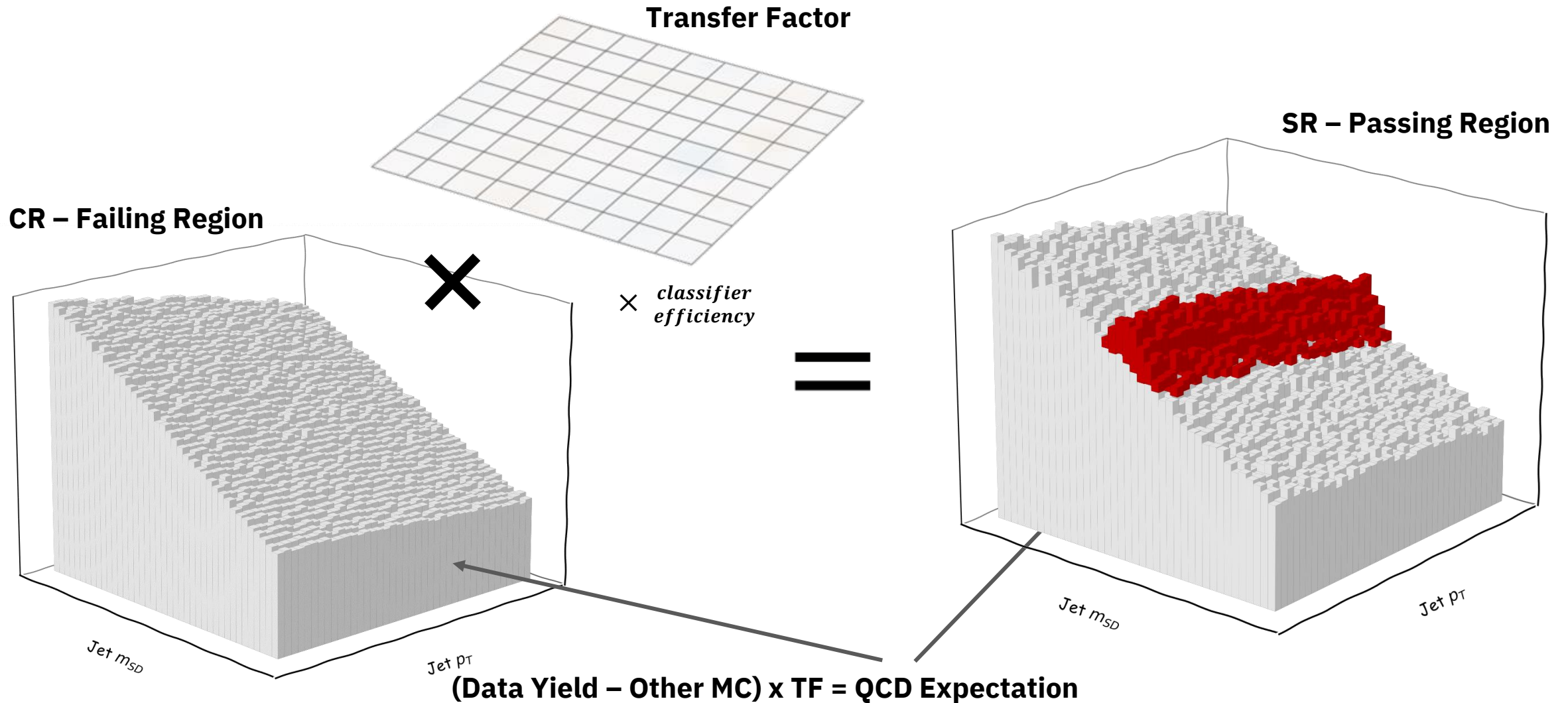
## Order of polynomials is arbitrary

- Determine optimal configuration based from goodness of fit (F-test)

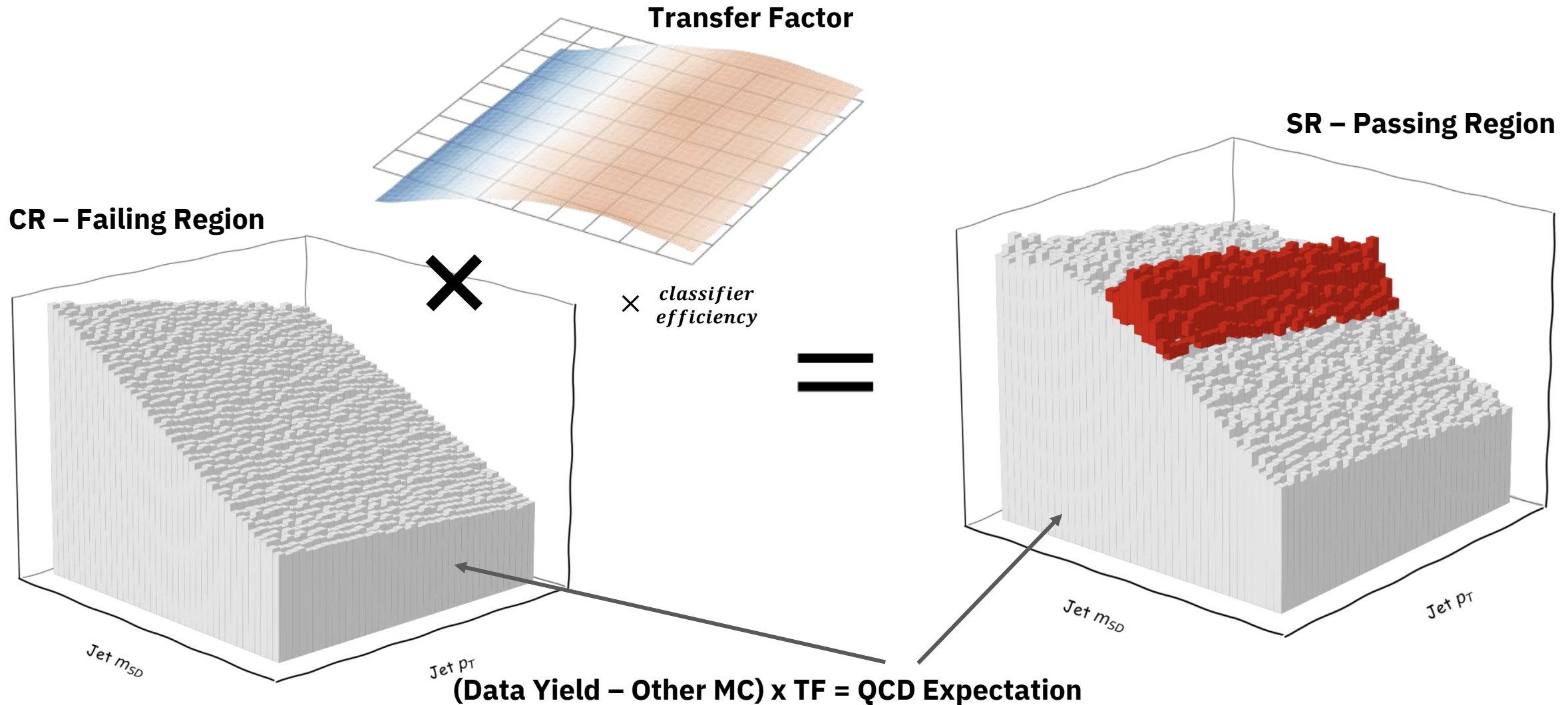
## Simultaneous Pass and Fail Fit



# Transfer Factor Fit – Optimal Case (ABCD)



# Transfer Factor Fit – Realistic Example



# Transfer Factor Fit – In Practice

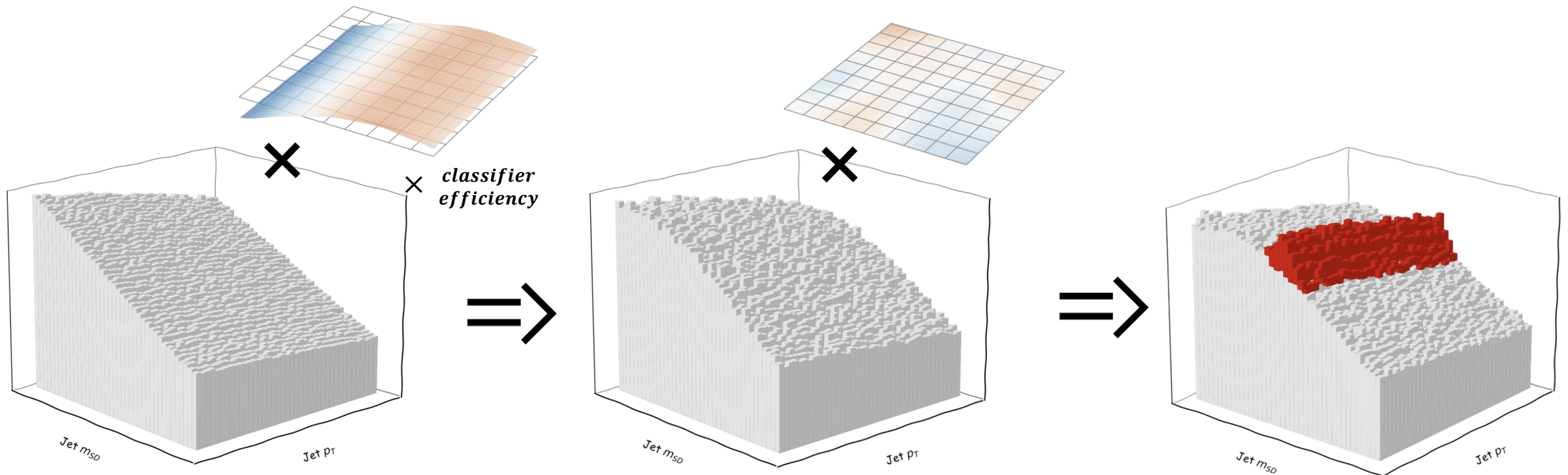
**2 effects** in to parametrize

- Classifier **Mass Sculpting** (MC only fit)
- **Discrepancies** due to **Data/MC** mismodelling

F-Tests to optimize #dof

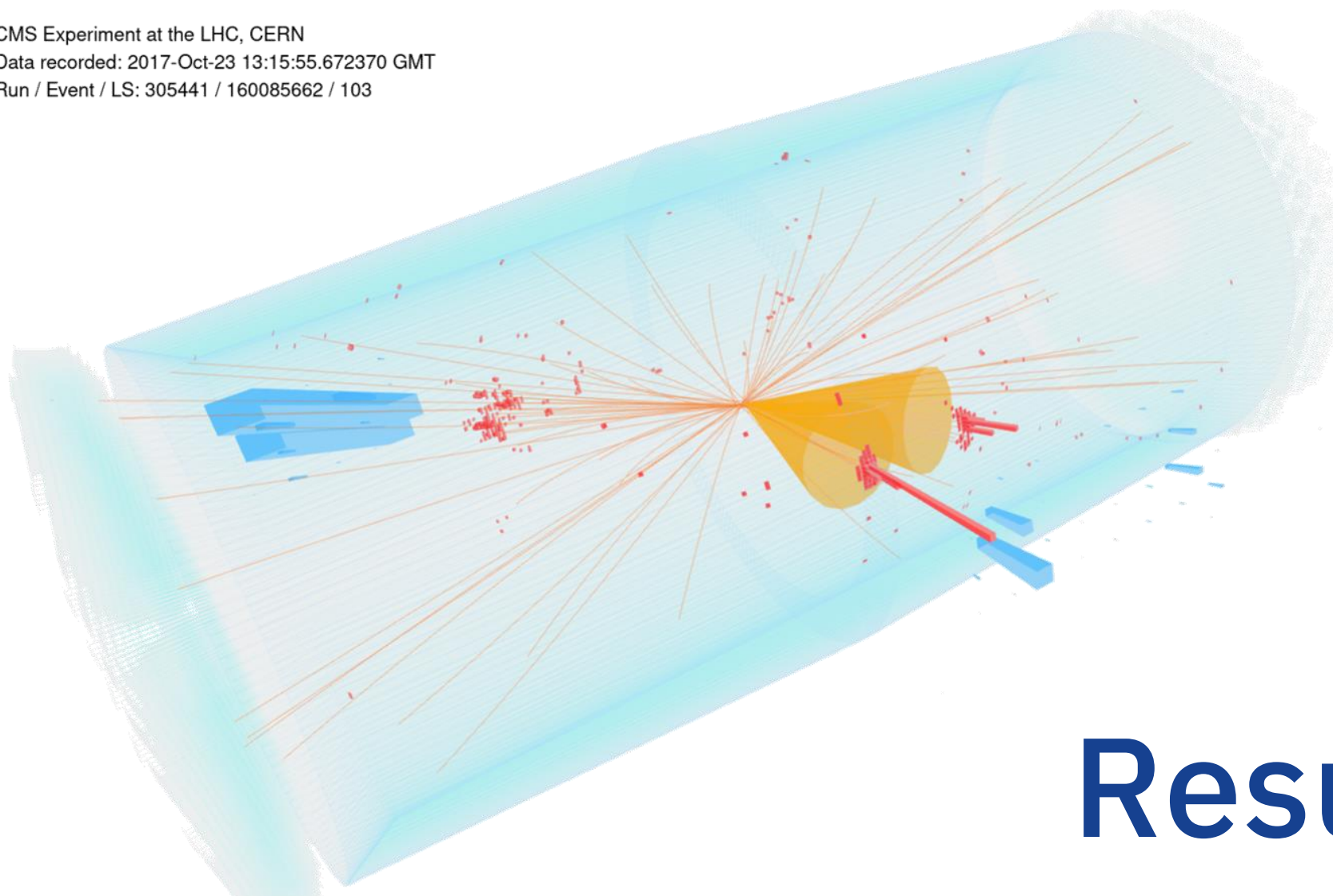
Method tested against bias

- Spurious “peaky” signal
- Different TF parametrizations



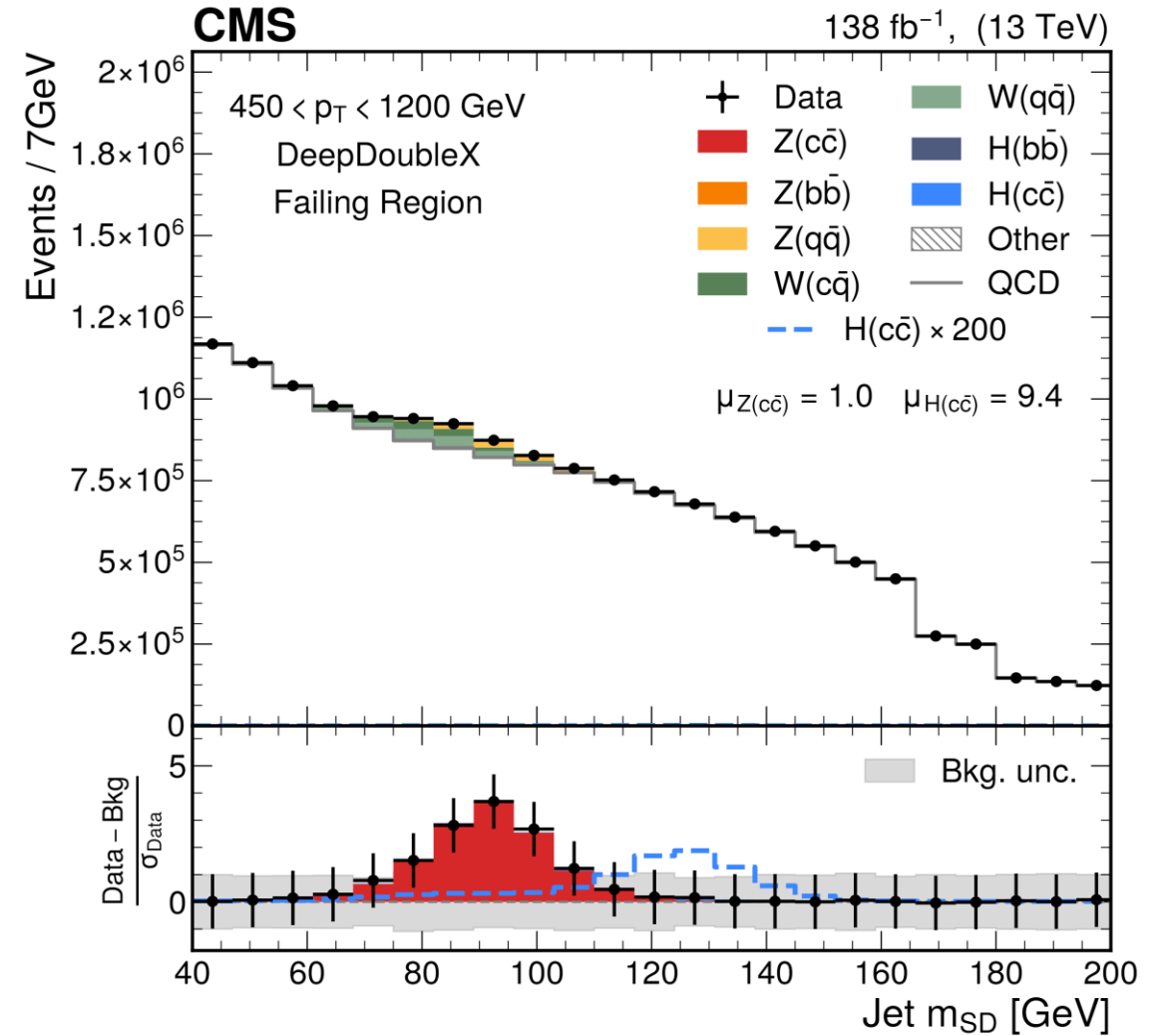
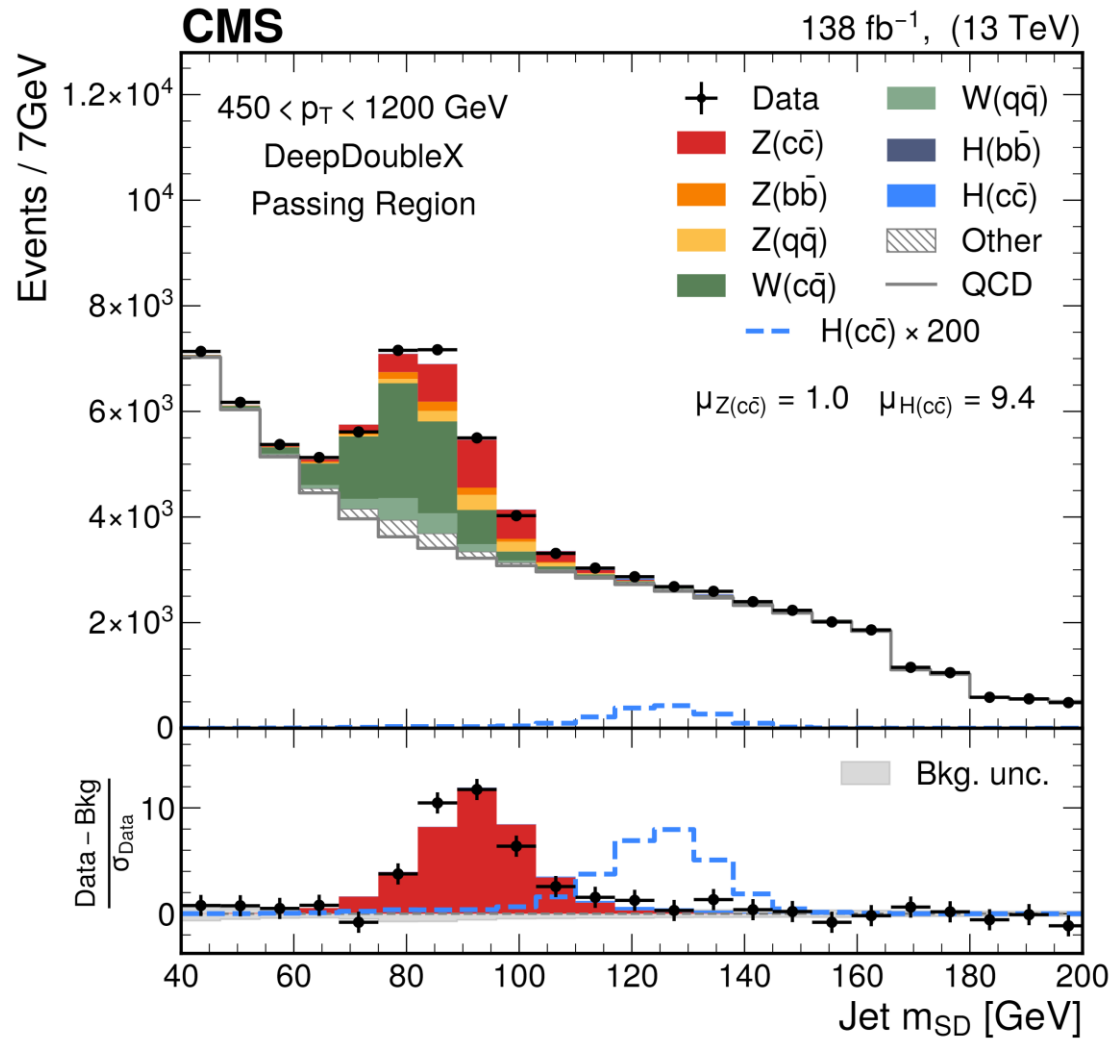


CMS Experiment at the LHC, CERN  
Data recorded: 2017-Oct-23 13:15:55.672370 GMT  
Run / Event / LS: 305441 / 160085662 / 103



# Results

# Results – $Z \rightarrow cc$ Observation

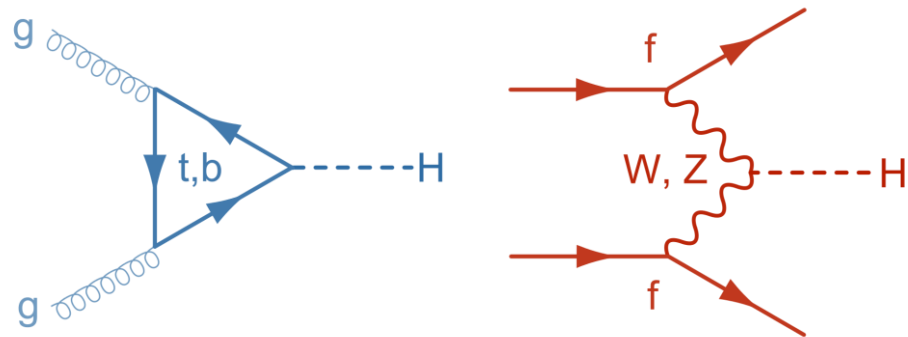




# ggH $\rightarrow$ cc Limit

Inclusive in H production mode

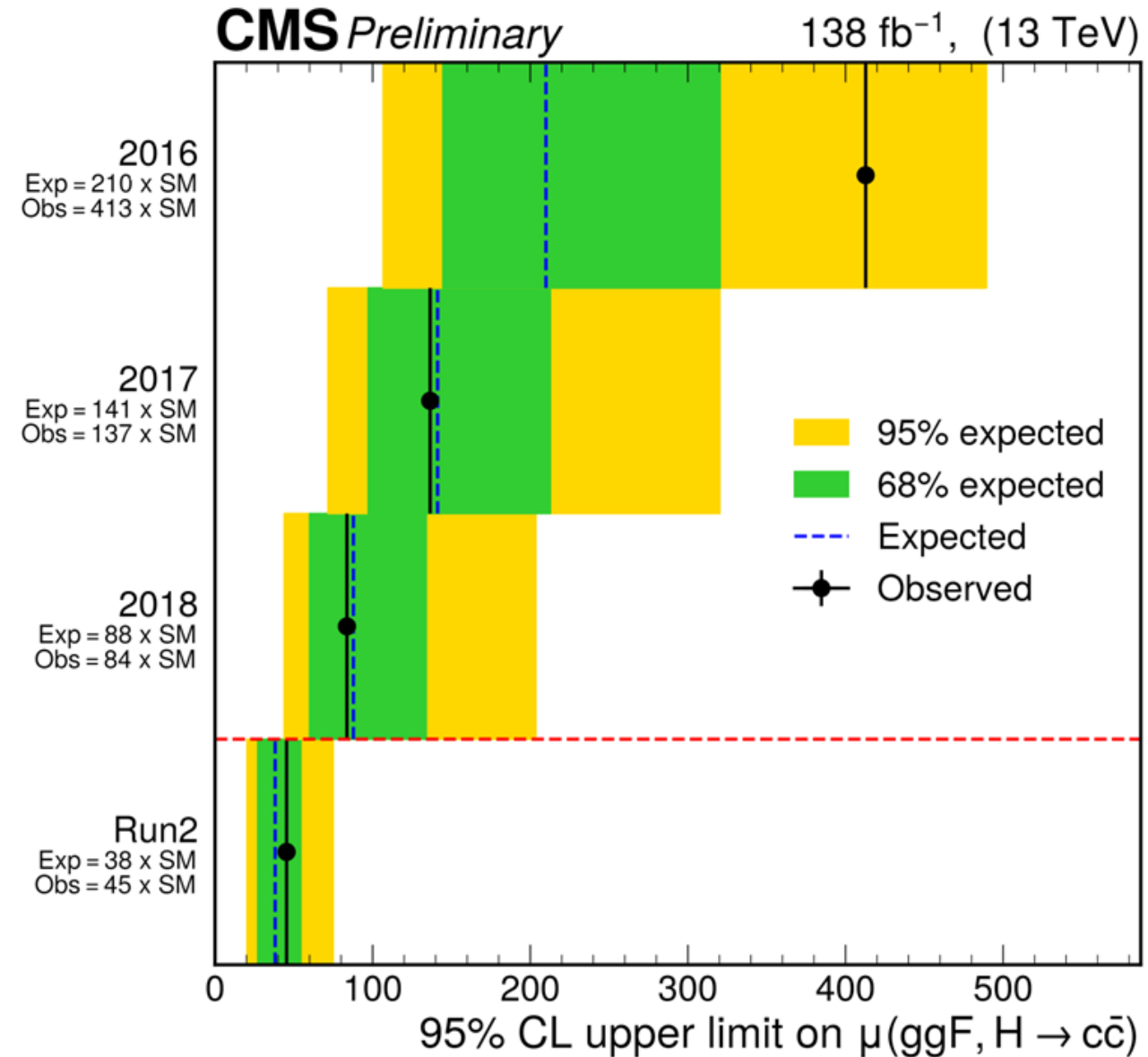
- $\sim 50\%$  ggF,  $\sim 30\%$  VBF



$\mu_{H \rightarrow cc} < 47$  (39)  $\times$  SM@95%

- **Orthogonal to  $VH \rightarrow cc$**
- Higher  $p_T$  regime

Statistical uncertainties dominant



# Summary

Full **Run2 ggH → cc Analysis** in a boosted regime, submitted to PRL

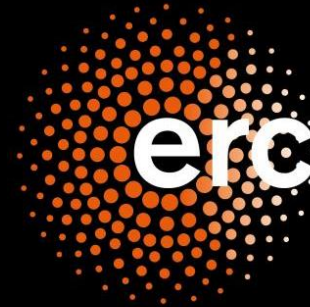
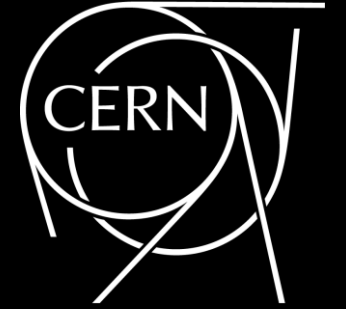
- Largely possible due to **DeepDoubleX** tagger
- First time this channel is explored at the LHC

**Observation** of **Z → cc** validating analysis method

- Significance  $\gg 5\sigma$
- Strongest constraint at the LHC yet
- First measurement in this production mode

Observed (expected) **H → cc Limit** < **47 (39) x SM expectation**

- Entirely orthogonal configuration to previous searches



**European Research Council**  
Established by the European Commission

# Thank You

This project has received funding from the European Research Council under the European Unions' Horizon 2020 research and innovation programme. [ID: 724704]

# Technical Aspects



Implementation is *nearly* ROOT independent

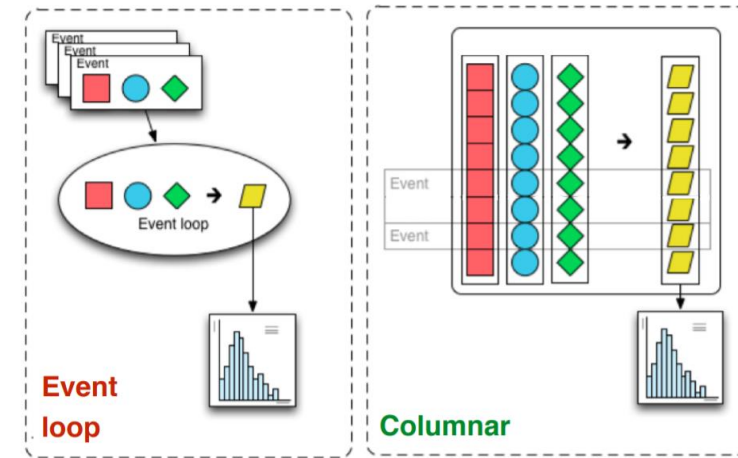
Using [Coffea](#) framework with Scikit-HEP Software stack

[awkward](#) – handling jagged data

[uproot](#) – ROOT file reading

(Fitting still dependent on RooFit/CMS Combine)

[mplhep](#) – plotting



Easy scale out

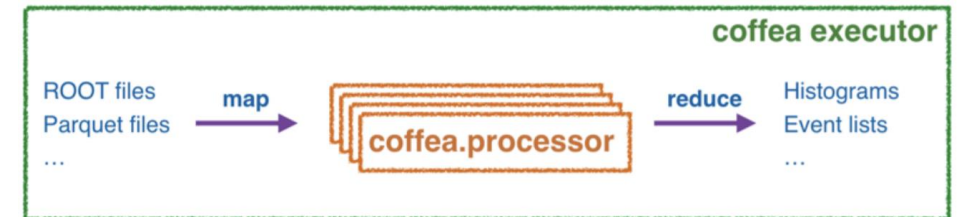


Full Run2 analysis ~ **30TB data+sim**, ~**3TB branches accessed**

Incl. systematic variations (JES/JER...)

10-20 nodes of 40-120 threads depending on availability

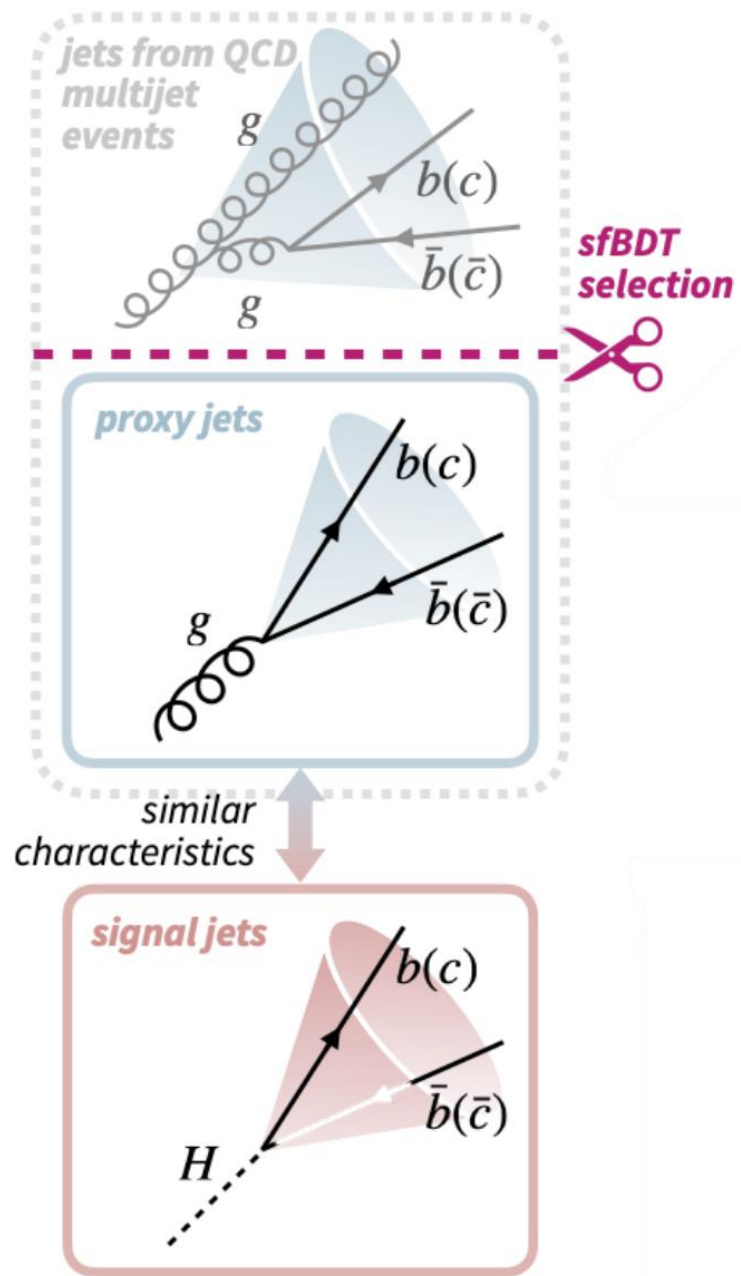
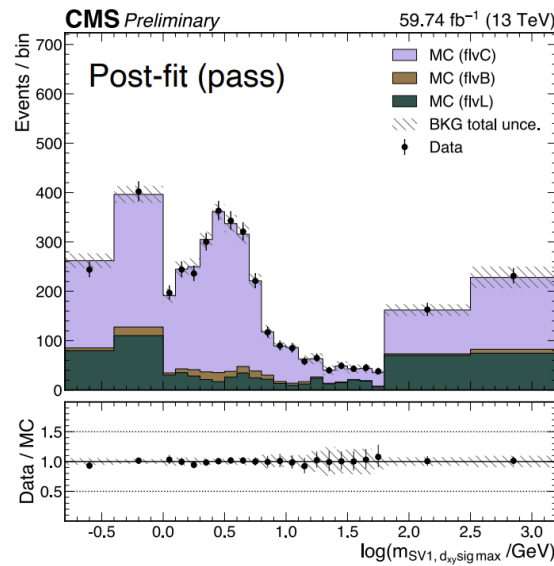
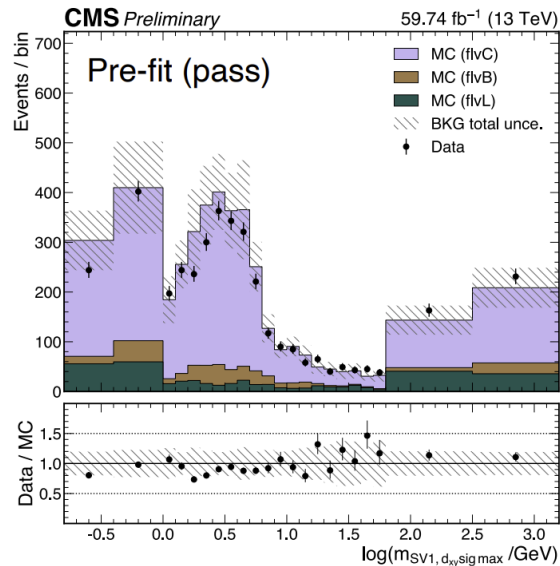
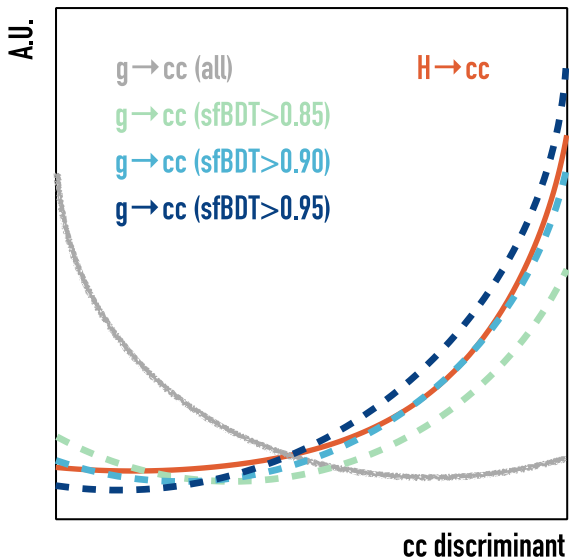
< **40 minutes total runtime** (I/O limited) in optimal conditions



# Tagger Calibration

- No existing pure  $H \rightarrow cc$  or  $Z \rightarrow cc$  region exists
- Use  $g \rightarrow cc$  as a proxy (not entirely signal-like)
- Two ways of “forcing” a good proxy – compatible results
  1. Use soft-muon presence to select Higgs-like  $cc$  jets
  2. Train a BDT to select Higgs-like  $cc$  jets

Effects of the BDT



# Background Parametrization Alternatives

- Bernstein – baseline
  - Previously used, easy to fit
  - Polynomial space

1	-2	1
	2	-2
		1

**Bernstein**

1		
-1	2	
1	-8	8

**Chebyshev**

1		
	1	
		1

**Polynomial**

$$R_{p/f}(\rho, p_T) = \sum_{k=0}^{n_\rho} \sum_{\ell=0}^{n_{p_T}} a_{k,\ell} \left[ b_{k,n_\rho}(\rho) b_{\ell,n_{p_T}}(p_T) \right]$$

- Chebyshev
  - More sensitive to initial values
  - Polynomial space
  - Replace  $b_{k,l}$  terms with  $c_{k,l}$

- Exponential transform of Bernstein
  - Less stable
  - Independent parameter space

$$R_{p/f}(\rho, p_T) = \sum_{k=0}^{n_\rho} \sum_{\ell=0}^{n_{p_T}} a_{k,\ell} \exp \left[ b_{k,n_\rho}(\rho) b_{\ell,n_{p_T}}(p_T) \right]$$

