



山东大学  
SHANDONG UNIVERSITY

# Measurement of boosted $VH(b\bar{b})$ using MVA method at ATLAS

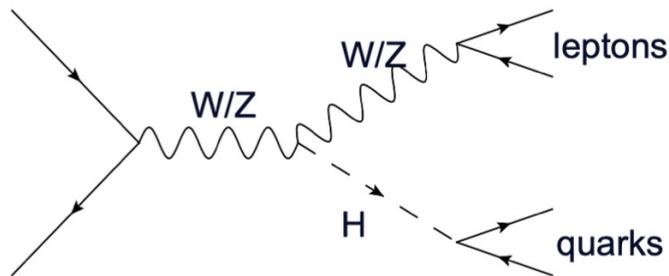
Jingyi Han

Shandong University

China LHC Physics Workshop, 27 Nov. 2021

➤ **Three  $VH(q\bar{q})$  stand-alone analyses** with full Run 2 dataset:

- $VH(b\bar{b})$  resolved - [Eur. Phys. J. C 81 \(2021\) 178](#)
- $VH(b\bar{b})$  boosted - [Phys. Lett. B 816 \(2021\) 136204](#)
- $VH(c\bar{c})$  - [ATLAS-CONF-2021-021](#), [paper draft](#)



V decays:

- $Z \rightarrow \nu\bar{\nu}$  : “0-lepton”;
- $W \rightarrow l\nu$  ( $l = e, \mu$ ) : “1-lepton”;
- $Z \rightarrow ll$  ( $l = e, \mu$ ) : “2-lepton”;

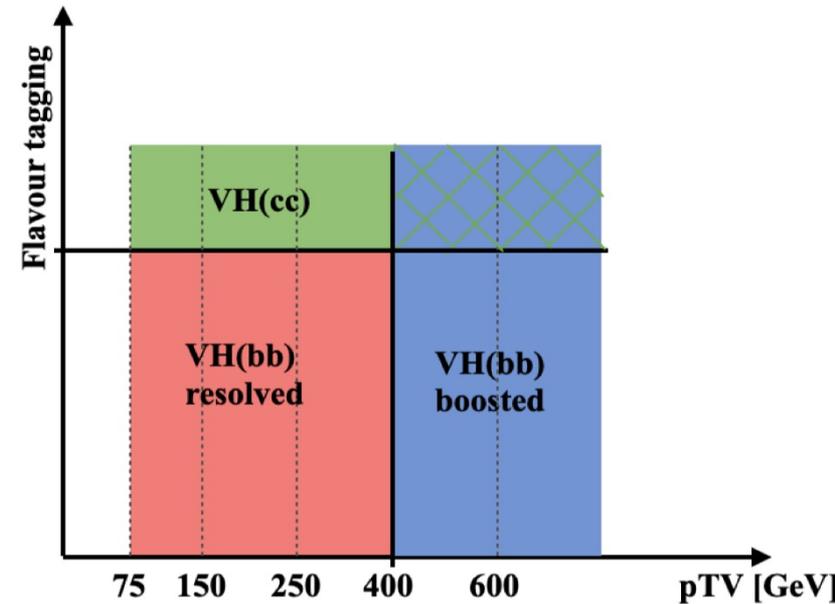
➤ **Goal of  $VH(b\bar{b}/c\bar{c})$  “legacy analysis”**: design one **coherent, harmonized** approach to extract  $VH(b\bar{b}/c\bar{c})$  results

# Analysis strategy

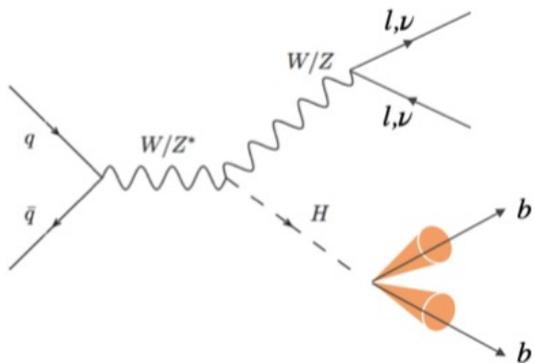
- No overlap between **VH(bb) resolved** and **VH(cc)** analysis
- Possible overlap between **VH(bb) boosted** and **VH(cc)** analysis
- Simple combination strategy between **VH(bb) resolved** and **VH(bb) boosted**

tested at the moment:

- Use **resolved analysis** in  $p_T^V < 400$  GeV region
- Use **boosted analysis** in  $p_T^V > 400$  GeV region
  - Add a split at **600 GeV**

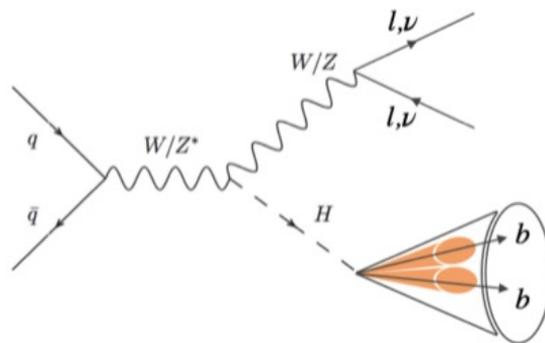


# Event selection



## VH( $b\bar{b}$ ) resolved

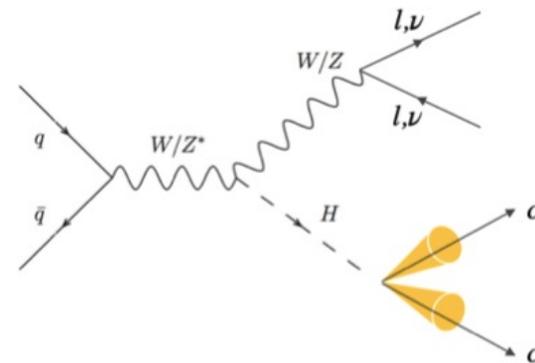
- $\geq 2$  pFlow jets with exactly **2 b-tags**;
- $150 \text{ GeV} < P_T^V < 400 \text{ GeV}$   
( $75 \text{ GeV} < P_T^V < 400 \text{ GeV}$  in 2L)



## VH( $b\bar{b}$ ) boosted

- $\geq 1$  large-R calo jet;
- Leading large-R associated with at least 2 VR jets;
- 2 b-tags in the leading 3 associated VR jets;
- $P_T^V > 400 \text{ GeV}$

**Final state discussed  
in following slides.**



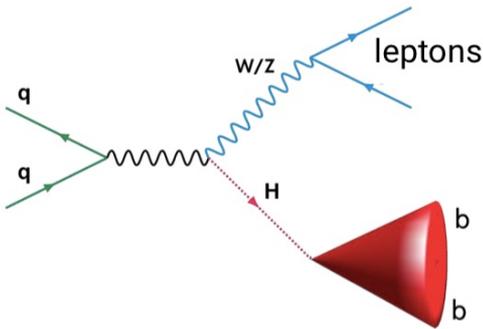
## VH( $c\bar{c}$ ) analysis

- $\geq 2$  pFlow jets;
- 2 leading jets **c-tagged + not b-tagged**
- $P_T^V > 150 \text{ GeV}$   
( $P_T^V > 75 \text{ GeV}$  in 2L)

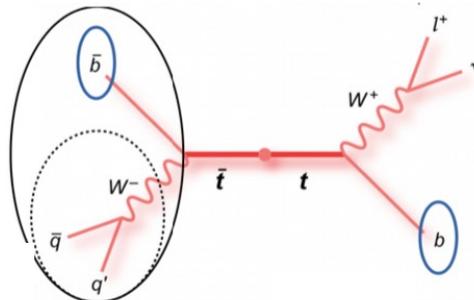
# Event categorization

Channel	Categories					
	$400 < p_T^V < 600 \text{ GeV}$			$p_T^V > 600 \text{ GeV}$		
	0 additional b-tagged track jets: <b>Signal region</b>		$\geq 1$ add. b-tagged track jets	0 additional b-tagged track jets: <b>Signal region</b>		$\geq 1$ add. b-tagged track jets
	0 add. small-R jets	$\geq 1$ add. small-R jets	<b>top control region</b>	0 add. small-R jets	$\geq 1$ add. small-R jets	<b>top control region</b>
0-lepton	high purity	low purity		high purity	low purity	
1-lepton	high purity	low purity		high purity	low purity	
2-lepton	SR			SR		

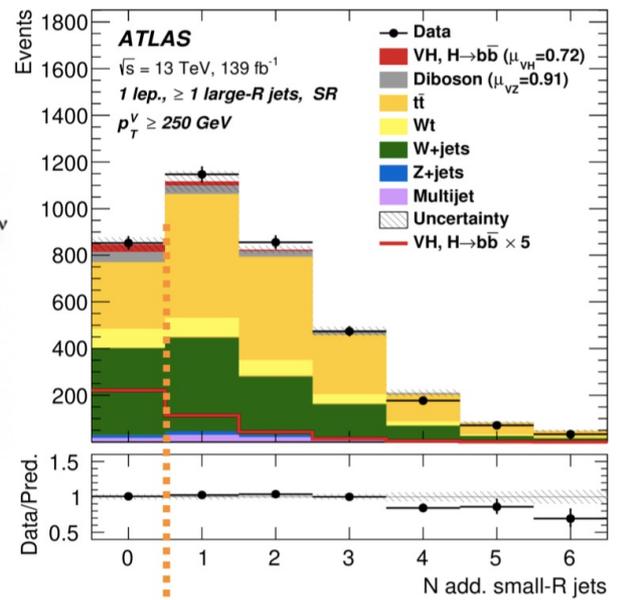
Signal:



ttbar:



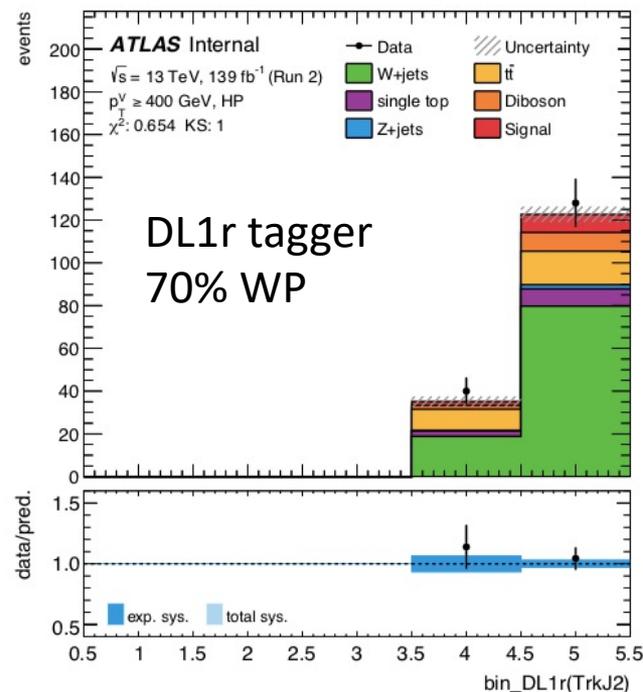
(Feynman diagrams at leading order)



# Multivariate analysis

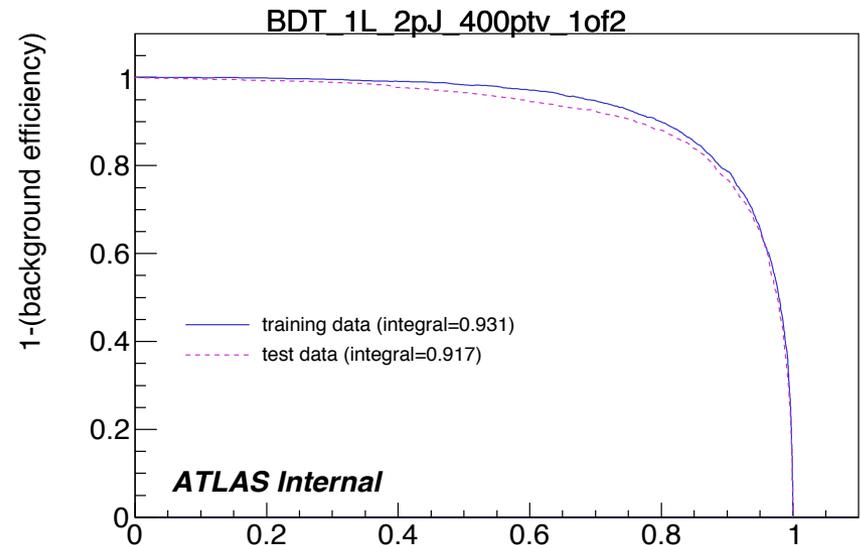
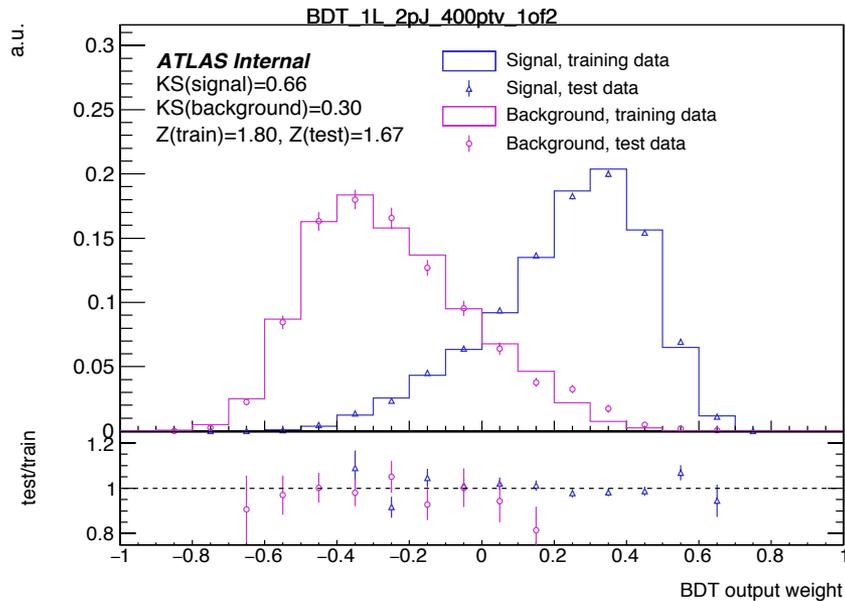
- **Multivariate analysis** used since long time in  $VH(b\bar{b})$  resolved analysis
  - **New: Introduce also in  $VH(b\bar{b})$  boosted analysis and in  $VH(c\bar{c})$  analysis**
- **Baseline: Boosted Decision Tree (BDT)**
- Input variables in  $VH(b\bar{b})$  boosted are inspired by information used in resolved MVA

	0L	1L	2L
mJ			
dR(TrkJ1,TrkJ2)			
pT(TrkJ1)			
pT(TrkJ2)			
pTV			
dPhi(V,J)			
bin_DL1r(TrkJ1)			
bin_DL1r(TrkJ2)			
MEff			
dY(V,J)	n/a		
Lepton pT imbalance / cosTheta(l,Z)	n/a	n/a	
N(add. calo jets)			
N(trk jets in large-R jet)			



# MVA baseline settings: 1L as an example

	$p_T^V > 400 \text{ GeV}$
MiniNodeSize	5
MaxDepth	3
NTrees	200
nCuts	60
AdaBoostedBeta	0.35



Training and test agreement looks ok.

# Boosted MVA fit

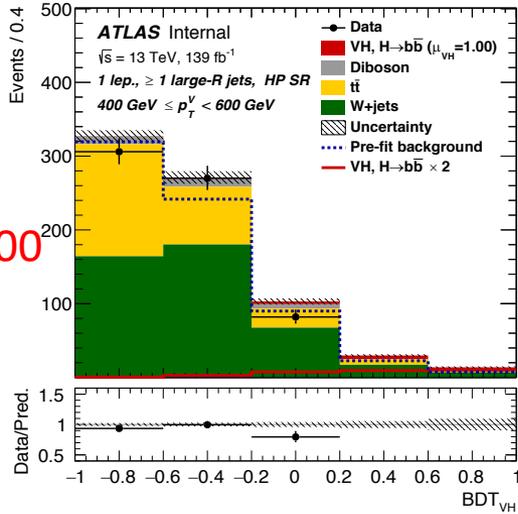
- First results of the MVA fit in the boosted regime
- Fit performed in  $p_T^V > 400$  GeV region
  - Split at  $p_T^V = 600$  GeV
  - HP SR and LP SR merged in  $p_T^V > 600$  GeV
  - MVA training performed in  $p_T^V > 400$  GeV, but evaluated in  $400 \text{ GeV} < p_T^V < 600 \text{ GeV}$  and  $p_T^V > 600 \text{ GeV}$
- Preliminary studies considering the **statistic only**

		0L		1L		2L	
		# add small-R jets		# add small-R jets		# add small-R jets	
		0	>=1	0	>=1	0	>=1
400-600 GeV	no add b-tagged jets	HP SR	LP SR	HP SR	LP SR	SR	
	add b-tagged jets	CR		CR			
>600 GeV	no add b-tagged jets	(HP + LP) SR		(HP + LP) SR		SR	
	add b-tagged jets	CR		CR			

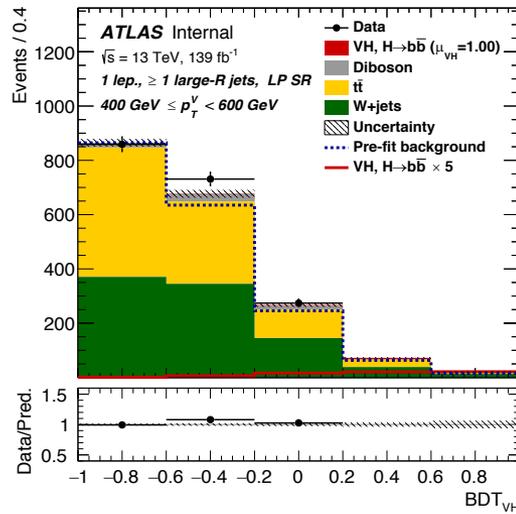
Fit in 8 SRs and 4 CRs

# Post-fit plots: 1L

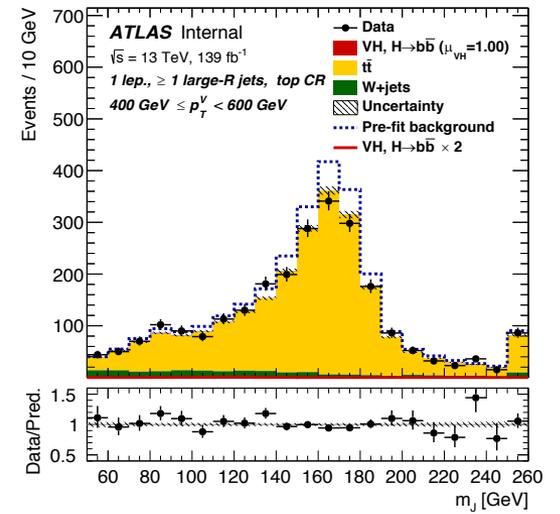
HP SR



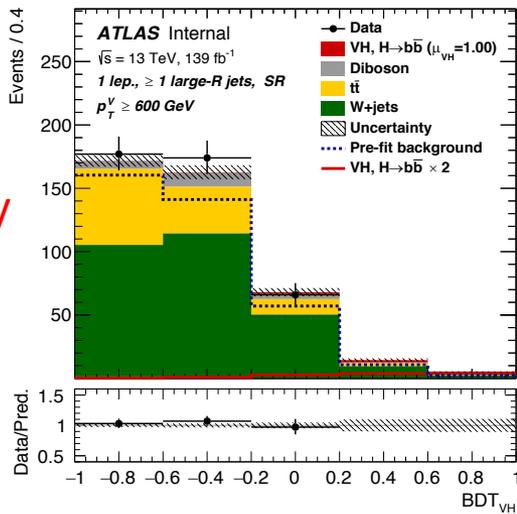
LP SR



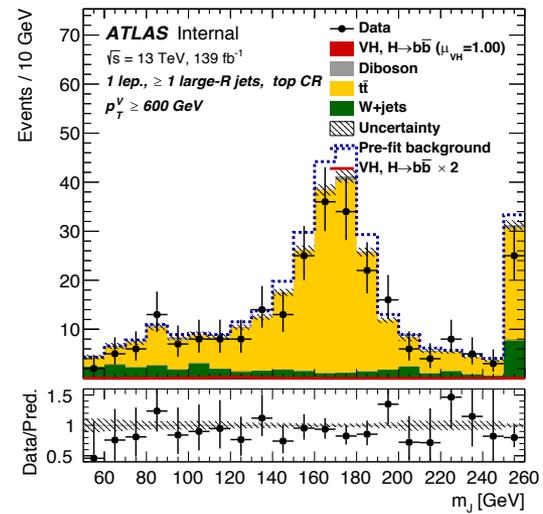
CR



SR



CR



# Cut-based analysis

- Cut based strategy used in the standard alone boosted  $VH(b\bar{b})$  analysis

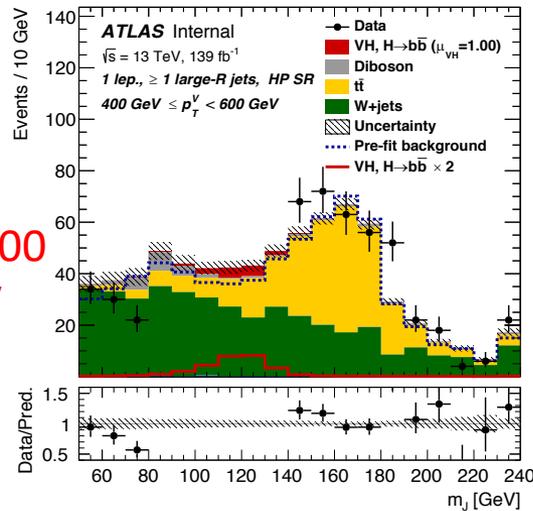
Selection	0 lepton channel	1 lepton channel		2 lepton channel	
		$e$ sub-channel	$\mu$ sub-channel	$e$ sub-channel	$\mu$ sub-channel
Trigger	$E_T^{\text{miss}}$	Single lepton	$E_T^{\text{miss}}$	Single lepton	$E_T^{\text{miss}}$
Leptons	0 VH-loose lepton	1 WH-signal lepton no second VH-loose lepton		$\geq 1$ ZH-signal lepton 2 VH-loose leptons	
$E_T^{\text{miss}}$	$> 250$ GeV	$> 50$ GeV	-	-	
$p_T^V$	$p_T^V > 250$ GeV				
Large-R jet	at least one large-R jet, $p_T > 250$ GeV, $ \eta  < 2$				
Track-Jets	at least two track-jets, $p_T > 10$ GeV, $ \eta  < 2.5$ , matched to the leading large-R jet				
$b$ -jets	leading two track-jets matched to the leading large-R must be $b$ -tagged				
$m_J$	$> 50$ GeV				
$\min[\Delta\phi(E_T^{\text{miss}}, \text{jets})]$	$> 30^\circ$	-			
$\Delta\phi(E_T^{\text{miss}}, H_{\text{cand}})$	$> 120^\circ$	-			
$\Delta\phi(E_T^{\text{miss}}, E_{T, \text{trk}}^{\text{miss}})$	$< 90^\circ$	-			
$ \Delta Y(V, H) $	-	$ \Delta Y(V, H)  < 1.4$			
$m_{ll}$	-	-		$66 \text{ GeV} < m_{ll} < 116 \text{ GeV}$	
lepton $p_T$ imbalance	-	-		$(p_T^{l_1} - p_T^{l_2})/p_T^Z < 0.8$	
lepton flavor	-	-		two lepton same flavour	
lepton charge	-	-		opposite sign muons	

Same large-R jet cuts as MVA analysis

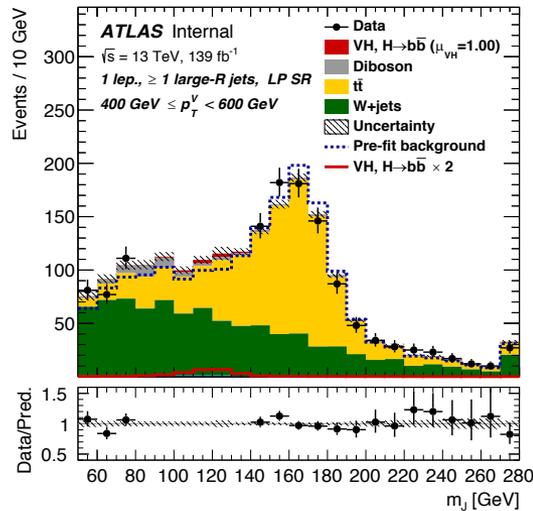
- Final discriminating variable is **mJ**, the leading large-R jet mass

# post-fit plots: 1L (cut based)

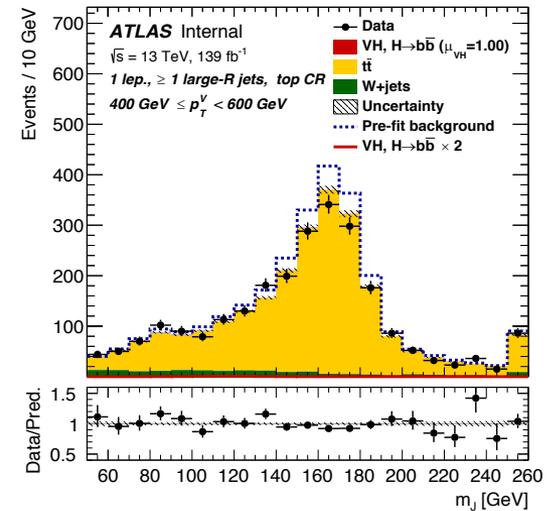
HP SR



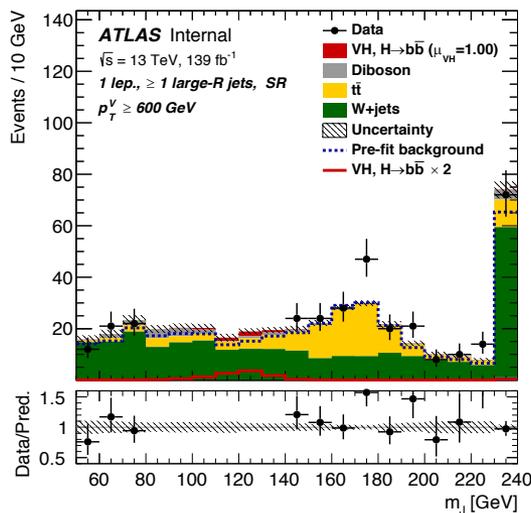
LP SR



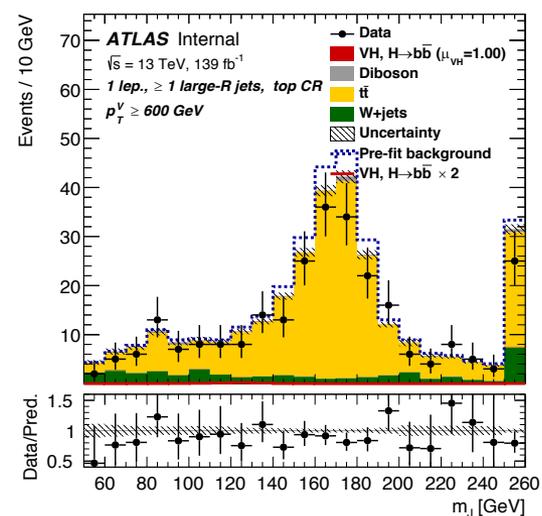
CR



SR



CR



# Significance: MVA vs. cut based

	<b>cut based</b>	<b>MVA</b>	<b>MVA improvement</b>
0-lepton fit	$1.04\sigma$	$2.10\sigma$	+102%
1-lepton fit	$1.06\sigma$	$2.35\sigma$	+122%
2-lepton fit	$0.70\sigma$	$1.46\sigma$	+109%
combined fit	$1.64\sigma$	$3.48\sigma$	+112%

(Significance are from statistic only maximum-profile-likelihood fit)

- Introduced the  $VH(bb^-/cc^-)$  “legacy” analysis
- Showed preliminary results for MVA analysis in  $VH(bb)$  boosted regime
- MVA analysis significance improves 112% compared with cut based analysis
- Target date for publication: **Moriond 2023**



**Back up**

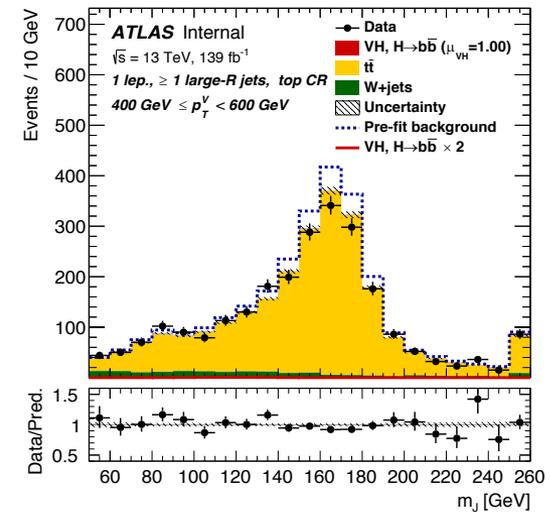
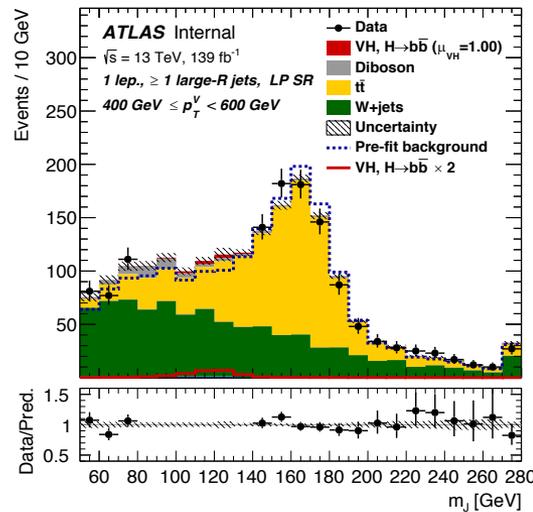
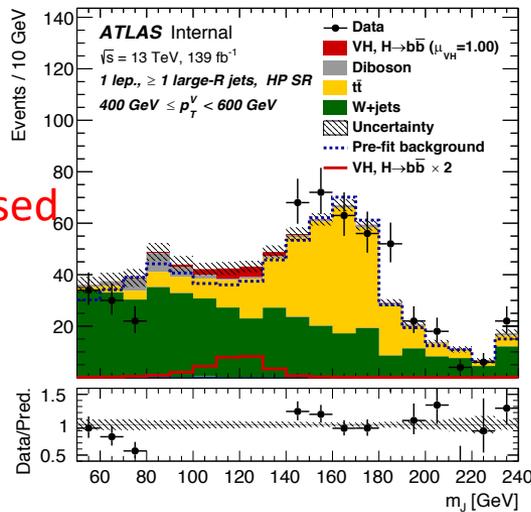
# Post-fit plots: mJ vs. BDT (400\_600GeV)

HP SR

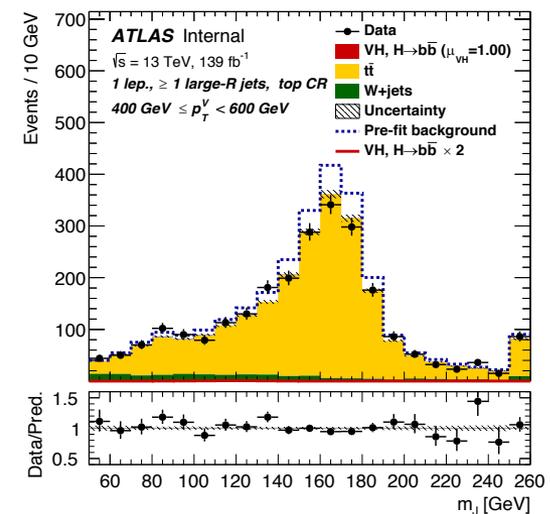
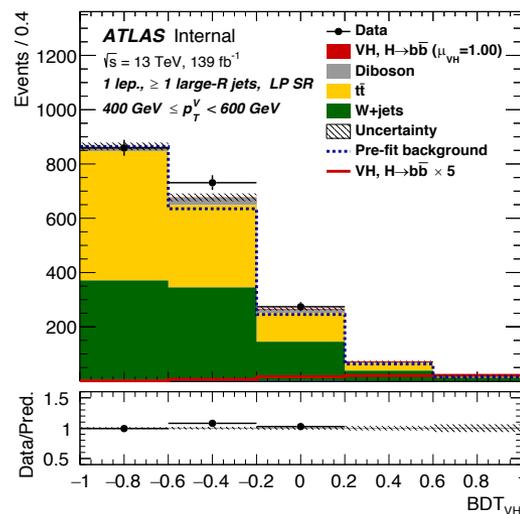
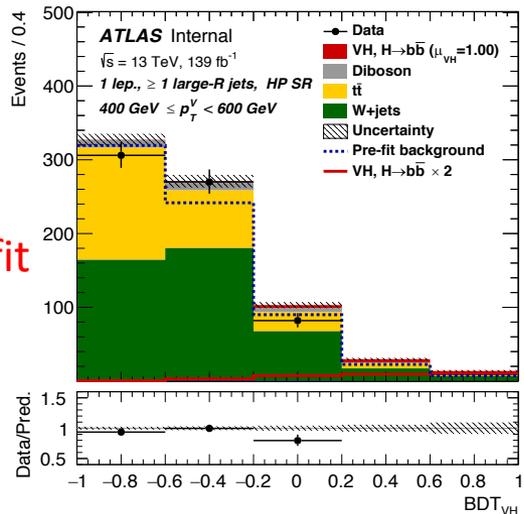
LP SR

CR

Cut-based fit



MVA fit

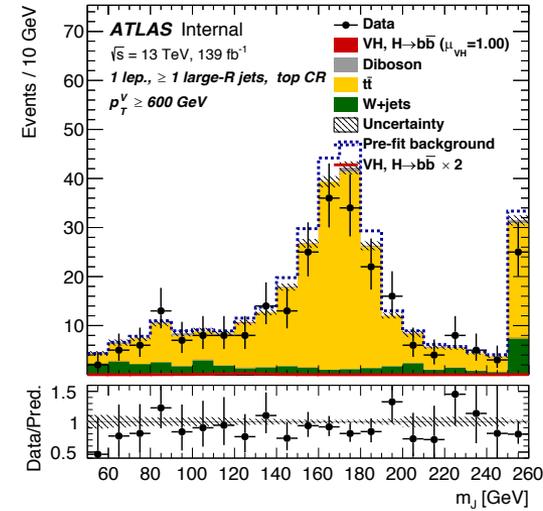
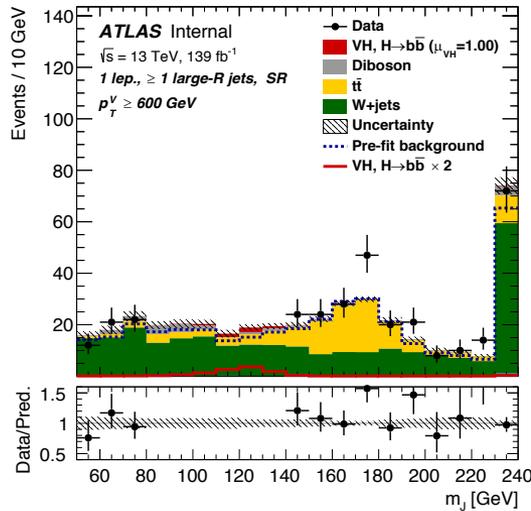


# Post-fit plots: mJ vs. BDT ( $p_T^V > 600\text{GeV}$ )

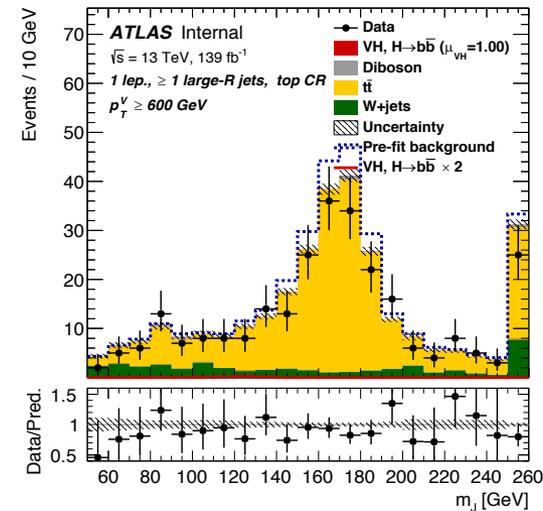
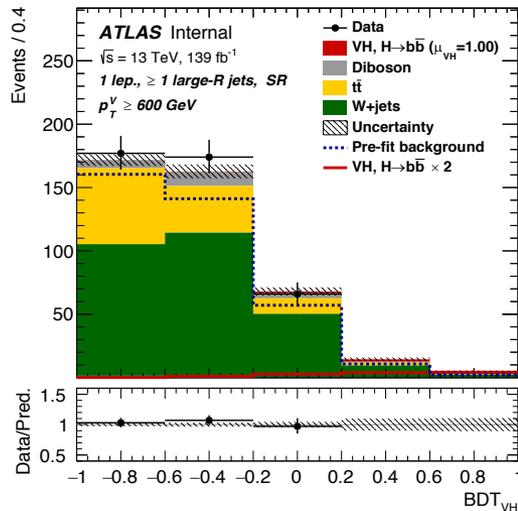
SR

CR

Cut-based fit

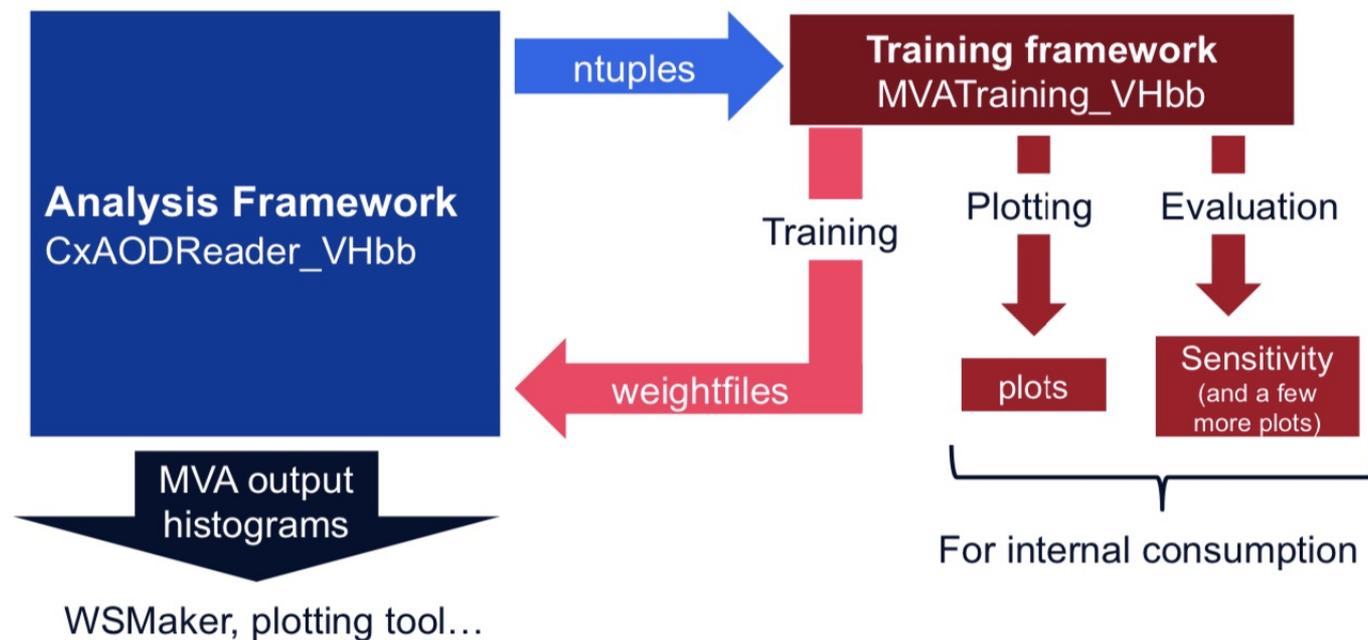


MVA fit



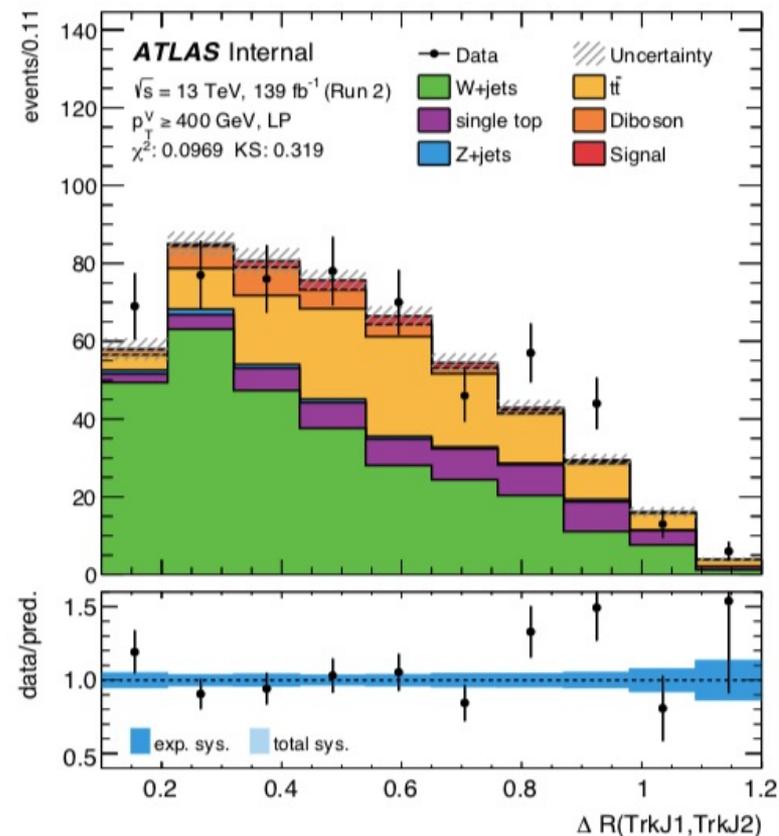
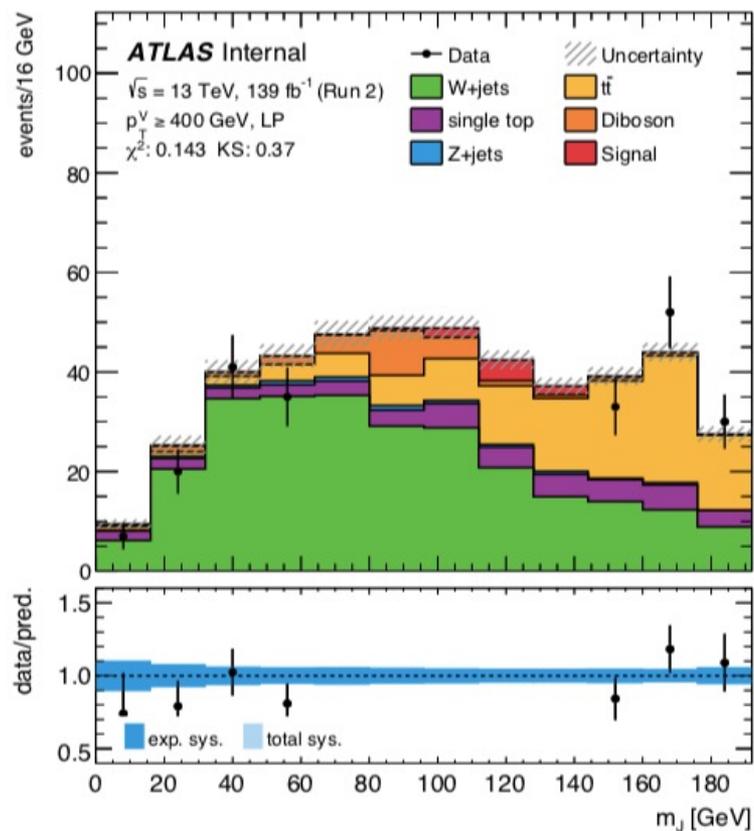
# Multivariate analysis

- **Multivariate analysis** used since long time in  $VH(bb\bar{)}$  resolved analysis
  - **New:** Introduce also in  $VH(bb\bar{)}$  boosted analysis and in  $VH(cc\bar{)}$  analysis



- **Baseline: Boosted Decision Tree (BDT)**

# Input variables data/MC (blinded)



$m_j$  and  $\Delta R(\mathbf{b1}, \mathbf{b2})$  are the most important two variables in the training.

More plots here: [link](#)

# Input variable importance ranking

```
-----  
Rank : Variable                : Variable Importance  
-----  
 1 : mJ                        : 1.789e-01  
 2 : deltaRbTrkJbTrkJ         : 1.447e-01  
 3 : pTbTrkJ2                 : 1.007e-01  
 4 : deltaYVJ                  : 9.876e-02  
 5 : NAdditionalCaloJets       : 9.132e-02  
 6 : pTbTrkJ1                 : 8.874e-02  
 7 : pTV                       : 7.915e-02  
 8 : absdeltaPhiVJ            : 6.894e-02  
 9 : NMatchedTrackJetLeadFatJet : 6.516e-02  
10 : bin_bTagBTrkJ2           : 4.597e-02  
11 : bin_bTagBTrkJ1           : 3.766e-02  
-----
```

# MVA fit: Transformation parameter optimization

Transformation of the BDT output distribution: optimise the expected sensitivity and reduce the number of bins

$$Z = z_s \frac{n_s}{N_s} + z_b \frac{n_b}{N_b}$$

Default :  
 $Z_s=3, Z_b=2$

- $N_s$  ( $N_b$ ): total number of signal (background) events in the histogram;
- $n_s$  ( $n_b$ ): number of signal (background) events in interval  $I[k, l]$  (interval between bin  $k$  and  $l$ );

The starting point for the transformation is the BDT output distribution with 500 equidistant bins between -1 and 1. The re-binning is then conducted using the following procedure:

1. Starting from the last bin on the right of the original histogram, increase the range of the interval  $I(k, last)$  by adding one after the other, the bins from the right to the left;
2. Calculate the value of  $Z$  at each step;
3. Once  $Z(I[k_0, last]) > 1$ , rebin all the bins in the interval  $I(k_0, last)$  into a single bin;
4. Repeat steps 1-3, starting this time from the last bin on the right, not included in the previous remap (the new last is  $k_0 - 1$ ), until  $k_0$  in the first bin.

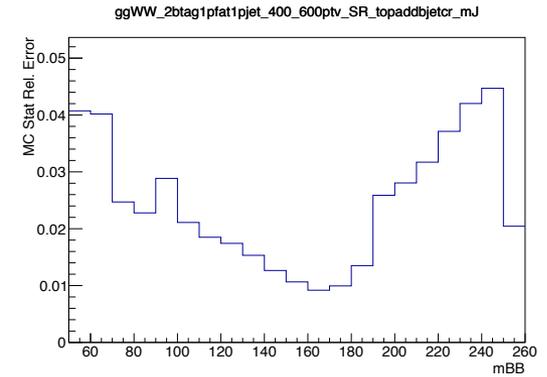
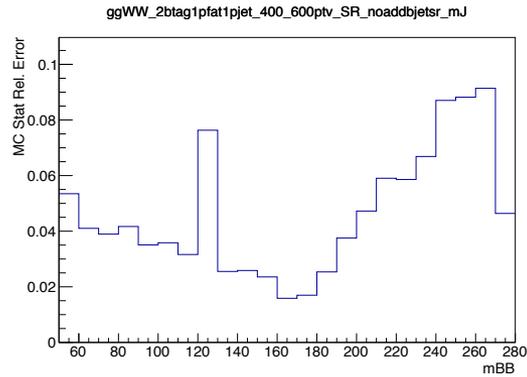
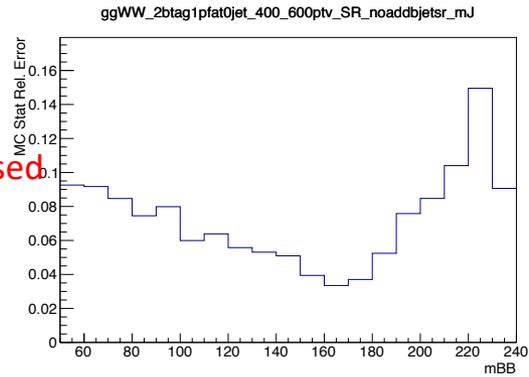
If the statistical uncertainty of the newly formed bin is larger than 20% step 2 is extended until the statistical uncertainty is below 20%.

# mJ fit vs. MVA fit: MC stat relative error: 400\_600 GeV (1L)

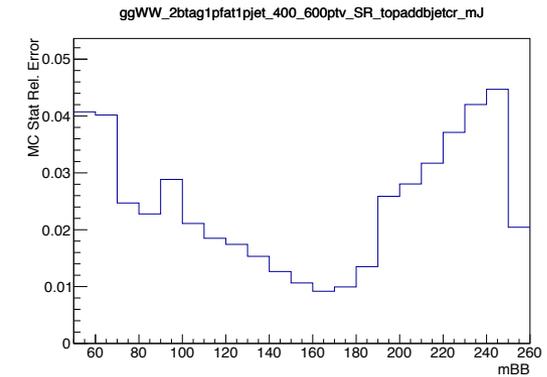
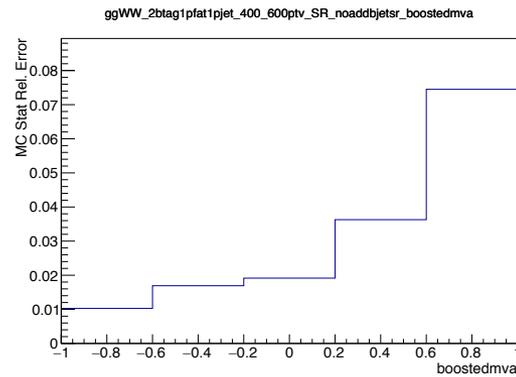
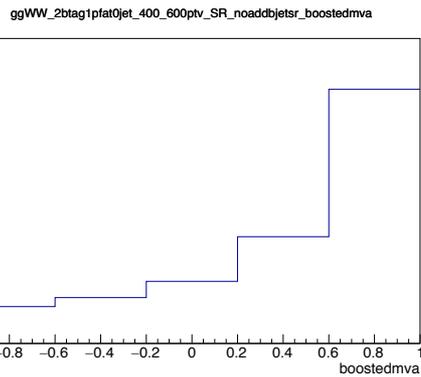
HP SR

LP SR

CR



Cut-based  
fit

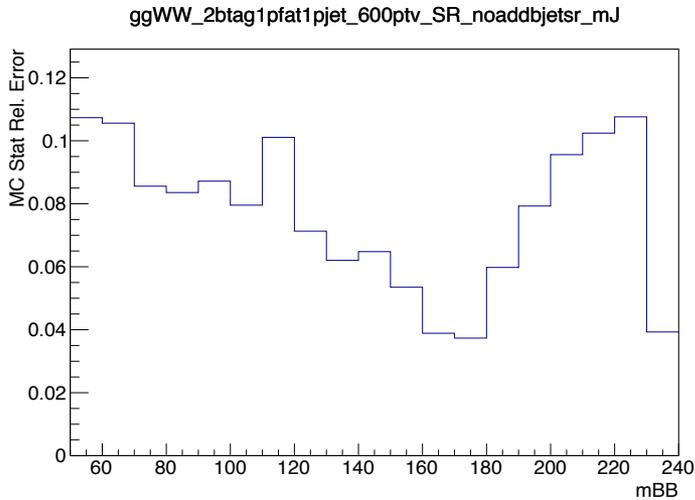


MVA

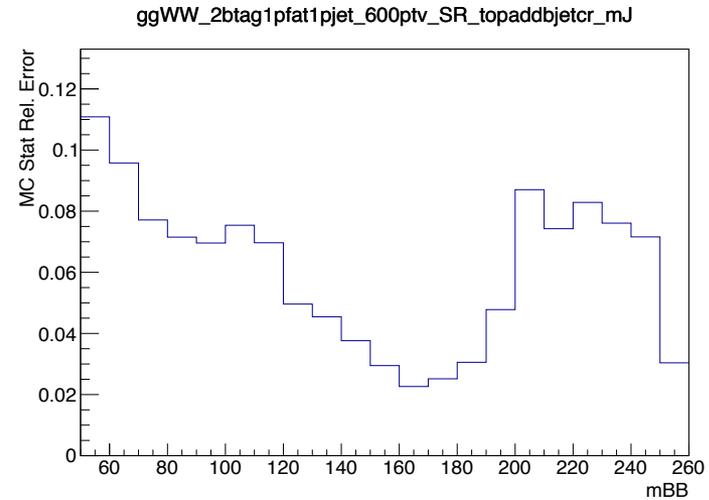
# mJ fit vs. MVA fit: MC stat relative error: $p_T^V > 600$ GeV (1L)

mJ fit

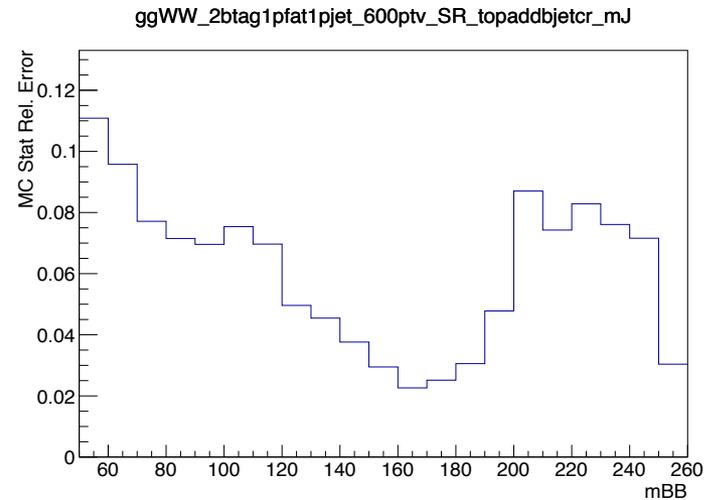
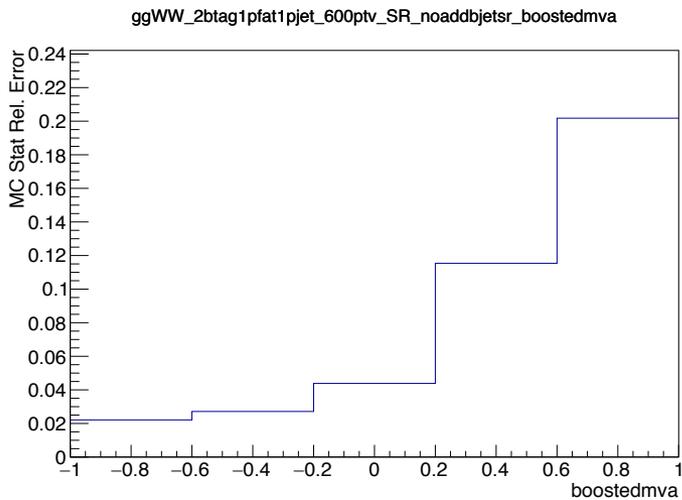
SR



CR



MVA fit



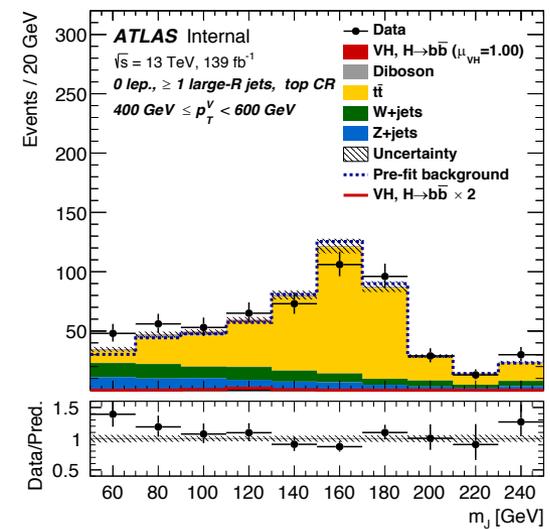
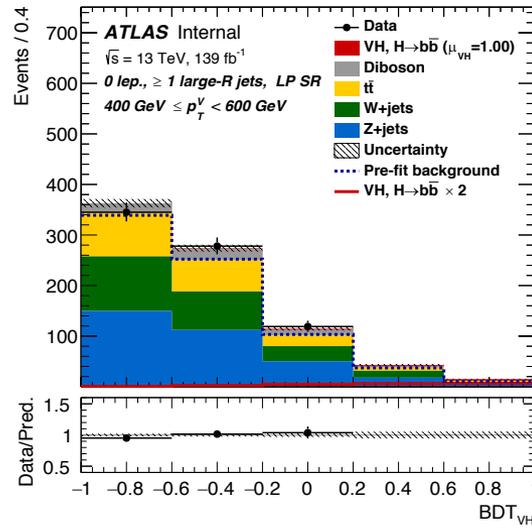
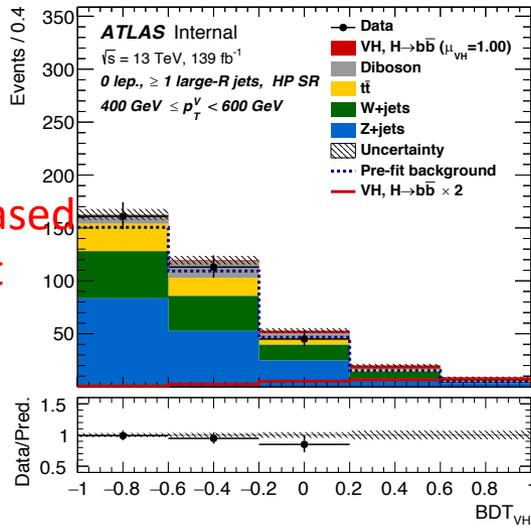
# 0L post-fit plots: mJ vs. BDT (400\_600GeV)

HP SR

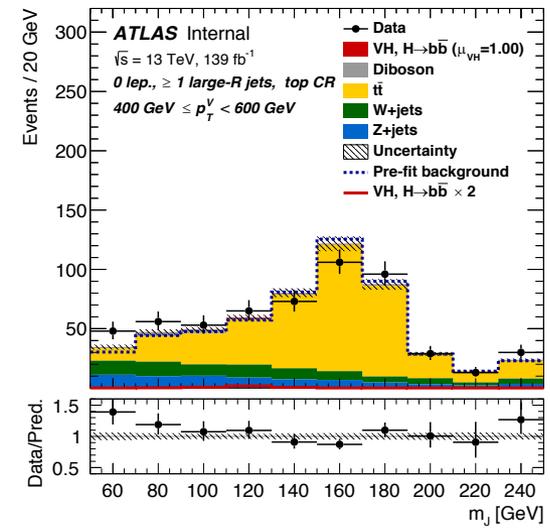
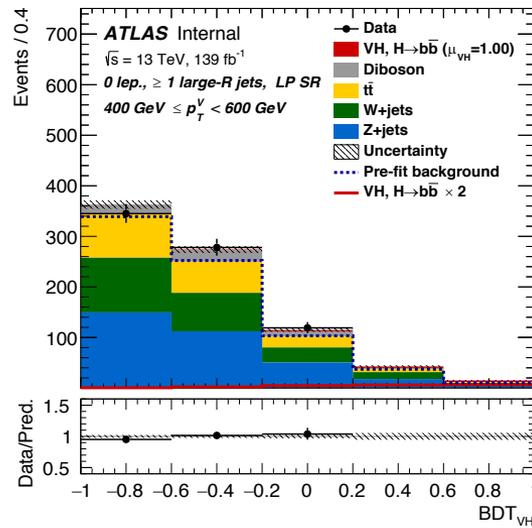
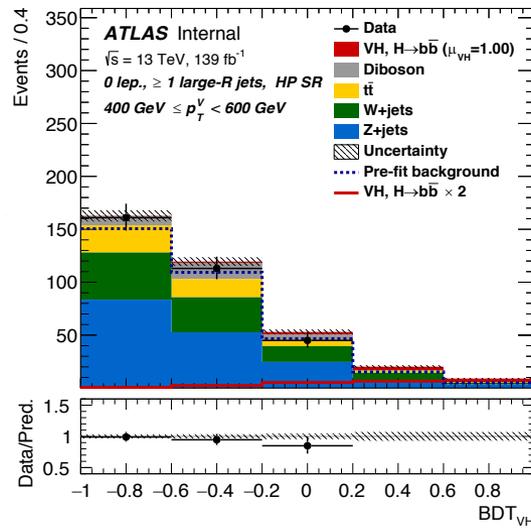
LP SR

CR

Cut-based fit

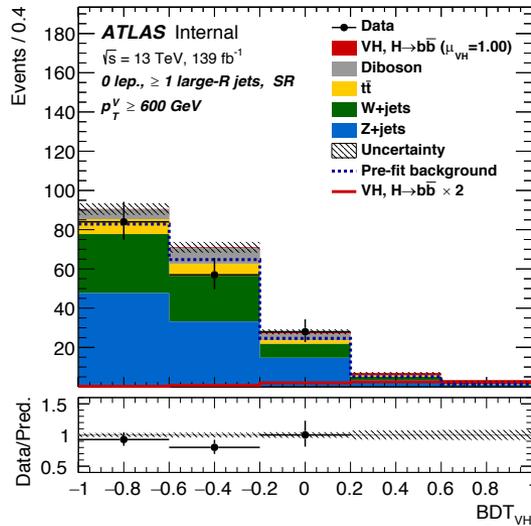


MVA

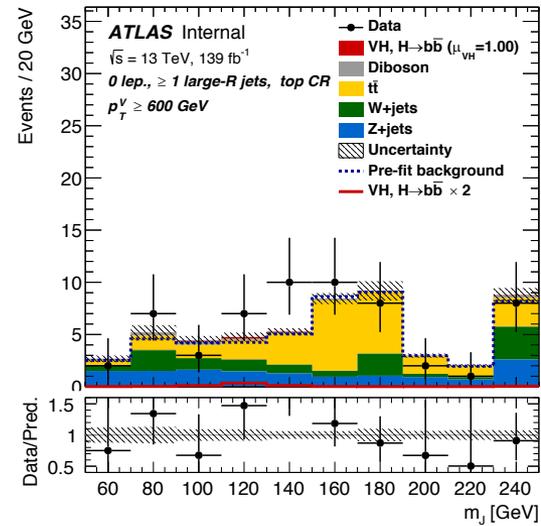


# 0L post-fit plots: mJ vs. BDT ( $p_T^V > 600\text{GeV}$ )

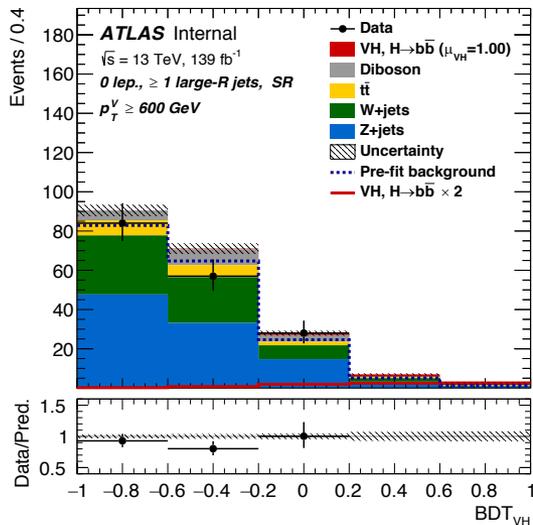
SR



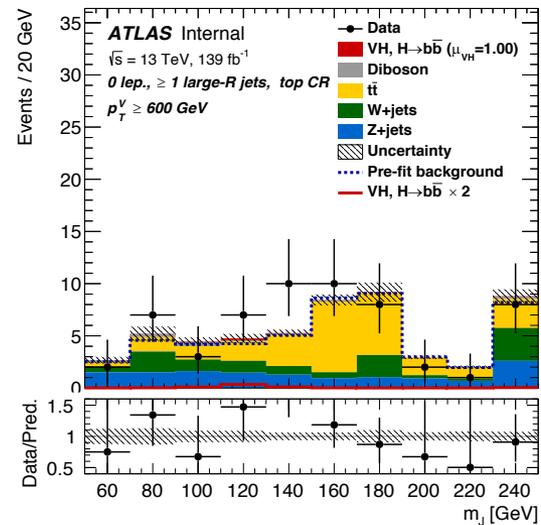
CR



Cut-based fit



MVA fit



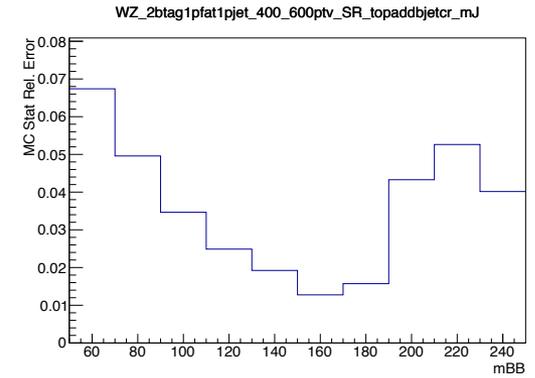
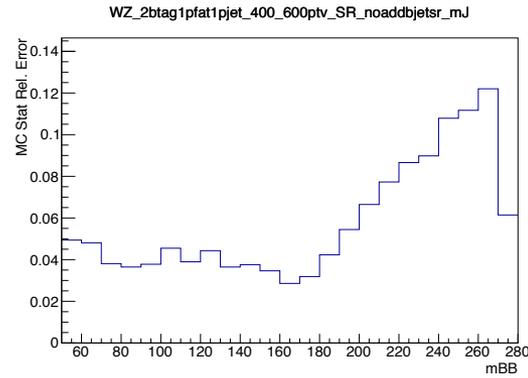
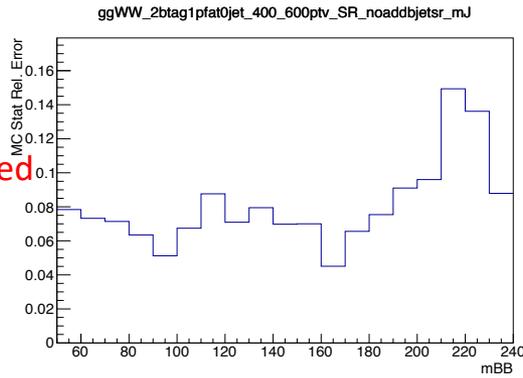
# mJ fit vs. MVA fit: MC stat relative error: 400\_600 GeV (0L)

HP SR

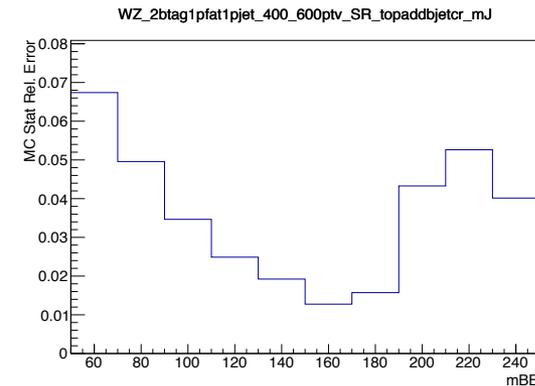
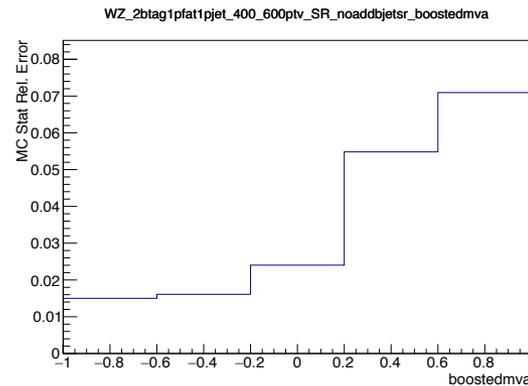
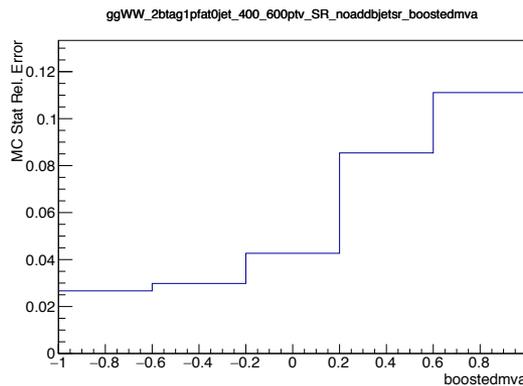
LP SR

CR

Cut-based  
fit



MVA

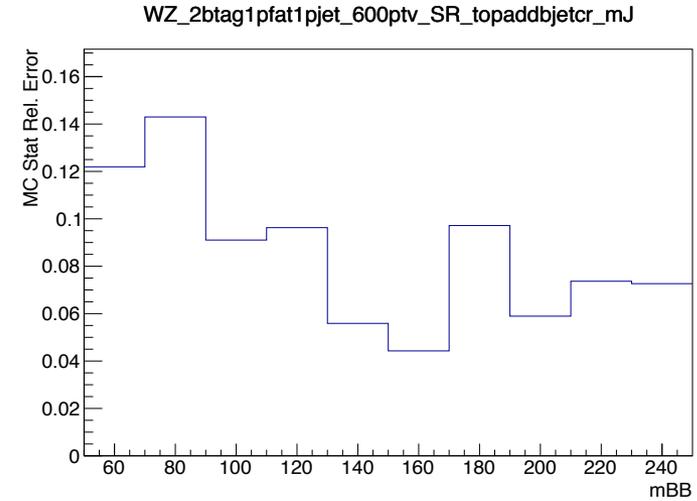
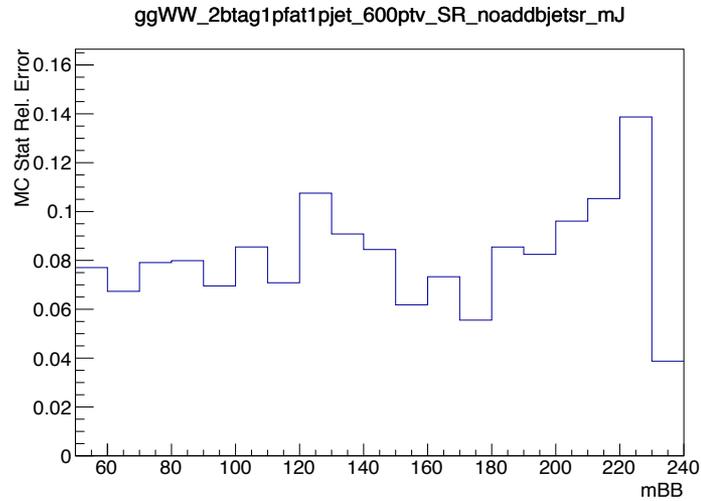


# mJ fit vs. MVA fit: MC stat relative error: $p_T^V > 600$ GeV (0L)

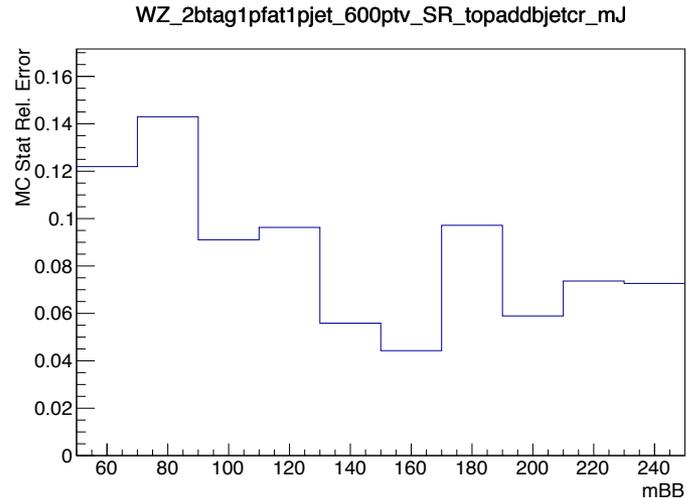
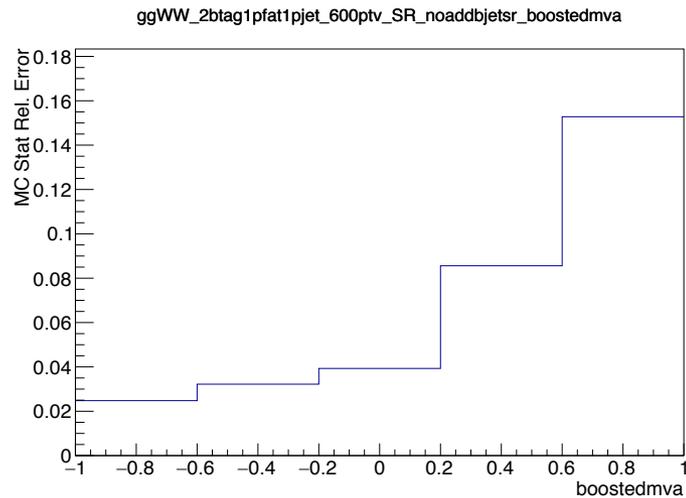
SR

CR

Cut-based  
fit



MVA fit

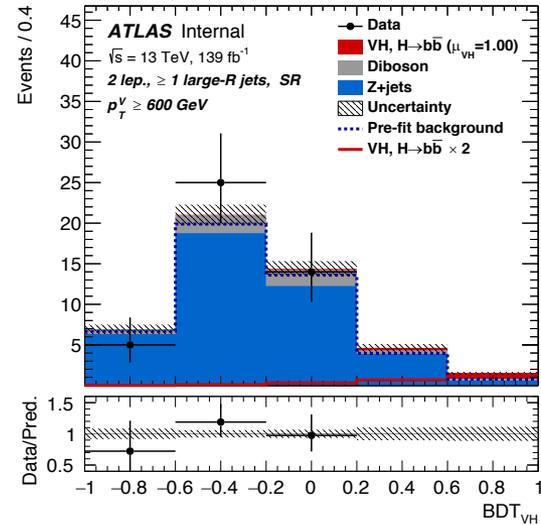
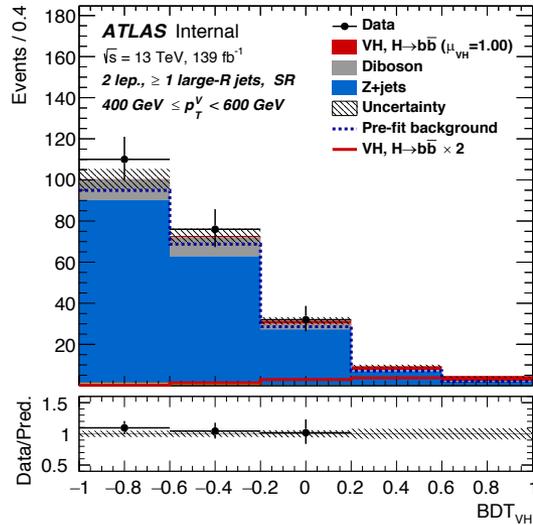


# 2L post-fit plots: mJ vs. BDT ( $p_T^V > 600\text{GeV}$ )

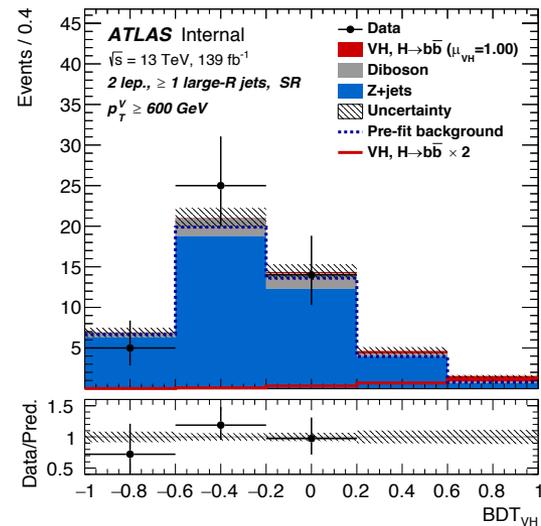
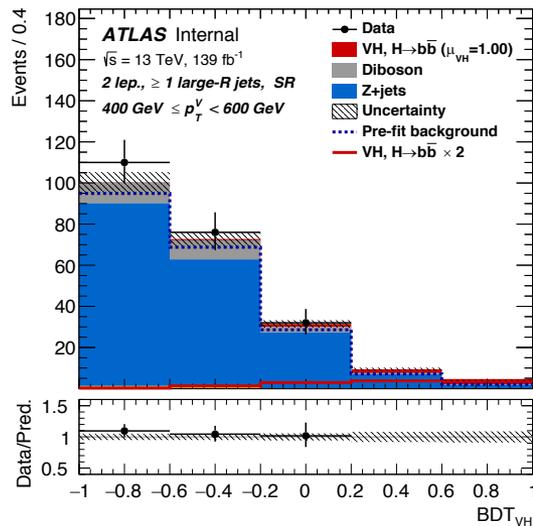
400\_600 GeV

$p_T^V > 600\text{ GeV}$

Cut-based fit



MVA fit

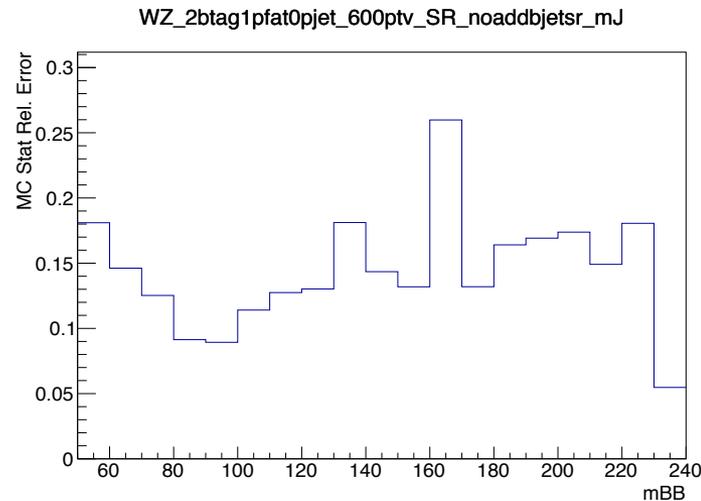
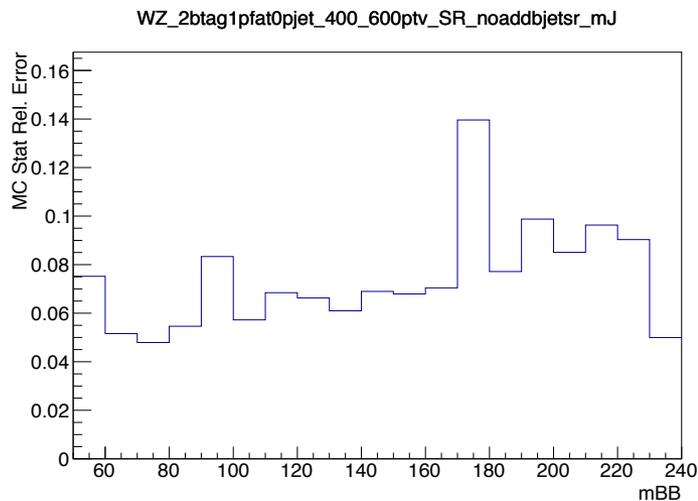


# mJ fit vs. MVA fit: MC stat relative error (2L)

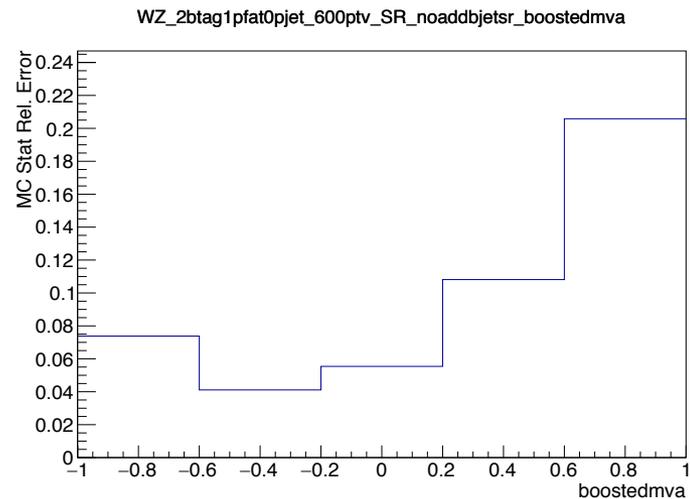
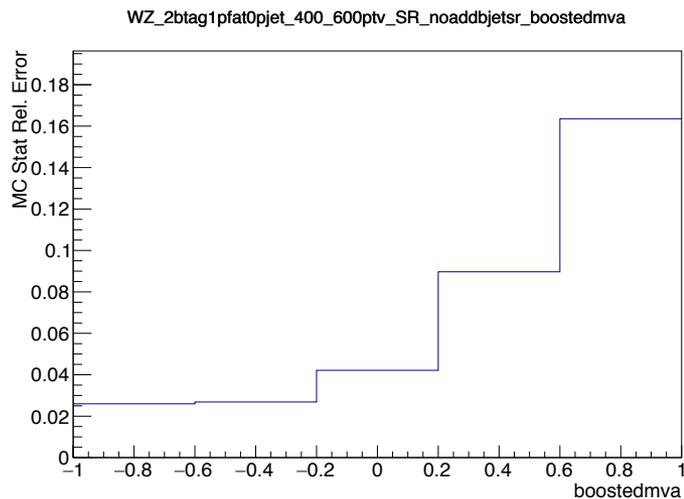
400\_600 GeV

$p_T^V > 600$  GeV

Cut-based  
fit



MVA fit



# Breakdown

Set of nuisance parameters	Impact on error			
Total	+ 0.628	- 0.615	+ - 0.621	
DataStat	+ 0.601	- 0.585	+ - 0.593	
FullSyst	+ 0.183	- 0.189	+ - 0.186	
Data stat only	+ 0.505	- 0.490	+ - 0.498	
Floating normalizations	+ 0.338	- 0.336	+ - 0.337	
Modelling: VH	+ 0.000	- 0.000	+ - 0.000	
Modelling: Background	+ 0.183	- 0.189	+ - 0.186	
Multi Jet	+ 0.000	- 0.000	+ - 0.000	
Modelling: single top	+ 0.000	- 0.000	+ - 0.000	
Modelling: ttbar	+ 0.000	- 0.000	+ - 0.000	
Modelling: W+jets	+ 0.000	- 0.000	+ - 0.000	
Modelling: Z+jets	+ 0.000	- 0.000	+ - 0.000	
Modelling: Diboson	+ 0.000	- 0.000	+ - 0.000	
MC stat	+ 0.183	- 0.189	+ - 0.186	
Experimental Syst	+ 0.000	- 0.000	+ - 0.000	
Detector: lepton	+ 0.000	- 0.000	+ - 0.000	
Detector: MET	+ 0.000	- 0.000	+ - 0.000	
Detector: JET	+ 0.000	- 0.000	+ - 0.000	
Detector: FATJET	+ 0.000	- 0.000	+ - 0.000	
Detector: FTAG (b-jet)	+ 0.000	- 0.000	+ - 0.000	
Detector: FTAG (c-jet)	+ 0.000	- 0.000	+ - 0.000	
Detector: FTAG (l-jet)	+ 0.000	- 0.000	+ - 0.000	
Detector: FTAG (extrap)	+ 0.000	- 0.000	+ - 0.000	
Detector: PU	+ 0.000	- 0.000	+ - 0.000	
Lumi	+ 0.000	- 0.000	+ - 0.000	

# mJ fit vs. MVA fit: Floating normalization

