

# MASS UNSPECIFIC SUPERVISED TAGGING (MUST) FOR BOOSTED JETS

**João Seabra**

Departamento de Física and CFTP, Instituto Superior Técnico, Lisboa



Phenomenology 2021 Symposium

24th of May 2021

Based on work made in collaboration with: **Juan A. Aguilar-Saavedra** and **Filipe R. Joaquim**

JHEP 03 (2021) 012

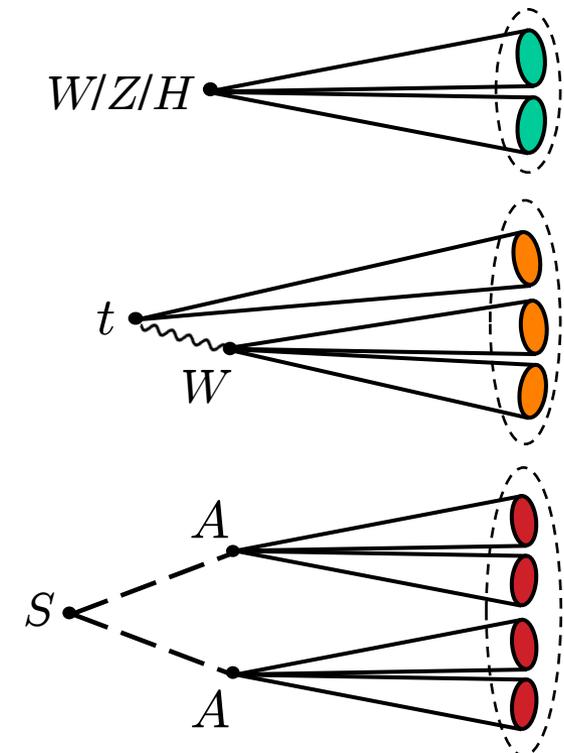
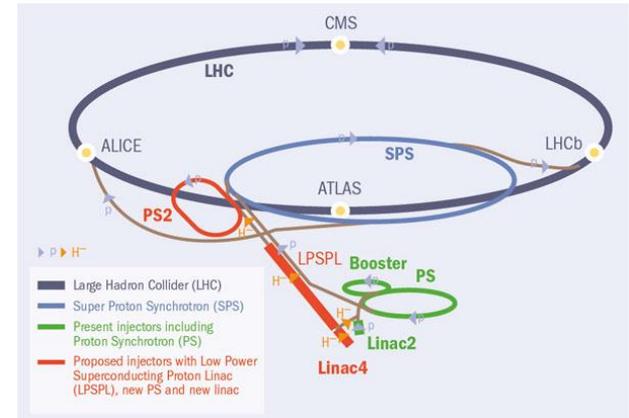


# MOTIVATION

- In the decades to come, the high-energy frontier of particle physics will continue to be explored at the **Large Hadron Collider (LHC)**;
- Most jets stem from **Quantum Chromodynamics (QCD)** processes...
- ...but when sufficiently boosted, the hadronic decays of Standard Model (SM) particles like the **W**, **Z** and **Higgs bosons** and the **top quark** also yield jets;
- **Hadronic decays of new particles** can produce jets too;
- A lot of theoretical frameworks beyond the SM predict multi-jet signals originated from direct or cascade decays of yet unseen particles.

e.g. J. A. Aguilar-Saavedra, F. R. Joaquim; JHEP 01 (2016) 183  
K. S. Agashe *et al.*; JHEP 05 (2017) 78

**Therefore...**



## Jet identification tools are crucial for new physics searches at the LHC.

### Searches for new gauge-bosons, scalars and spin-2 particles

A. M. Sirunyan *et al.* [CMS]; JHEP 08 (2017) 29  
A. M. Sirunyan *et al.* [CMS]; JHEP 09 (2018) 148  
M. Aaboud *et al.* [ATLAS]; Phys. Lett. B 781 (2018) 327  
M. Aaboud *et al.* [ATLAS]; Phys. Lett. B 788 (2019) 316  
M. Aaboud *et al.* [ATLAS]; Phys. Lett. B 783 (2018) 392  
M. Aaboud *et al.* [ATLAS]; Phys. Rev. D 98, 3 (2018) 32015  
A. M. Sirunyan *et al.* [CMS]; Phys. Rev. D 99, 1 (2019) 12005  
A. M. Sirunyan *et al.* [CMS]; Eur. Phys. J. C 80, 3 (2020) 237  
A. M. Sirunyan *et al.* [CMS]; Phys. Rev. D 100, 11 (2019) 112007  
G. Aad *et al.* [ATLAS]; Eur. Phys. J. C 80, 12 (2020) 1165

### Searches for vector-like quarks

A. M. Sirunyan *et al.* [CMS]; Phys. Lett. B 781 (2018) 574  
A. M. Sirunyan *et al.* [CMS]; Eur. Phys. J. C 79 (2019) 90  
M. Aaboud *et al.* [ATLAS]; JHEP 05 (2019) 41  
A. M. Sirunyan *et al.* [CMS]; Eur. Phys. J. C 79, 3 (2020) 36

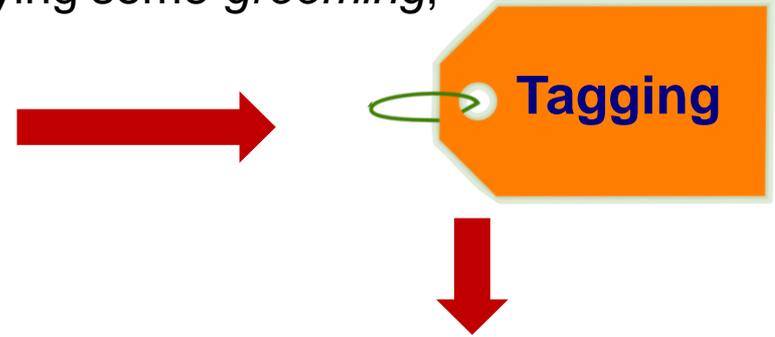
### Searches for dark-matter

A. M. Sirunyan *et al.* [CMS]; Eur. Phys. J. C 79, 3 (2019) 280

# JET IDENTIFICATION

It requires:

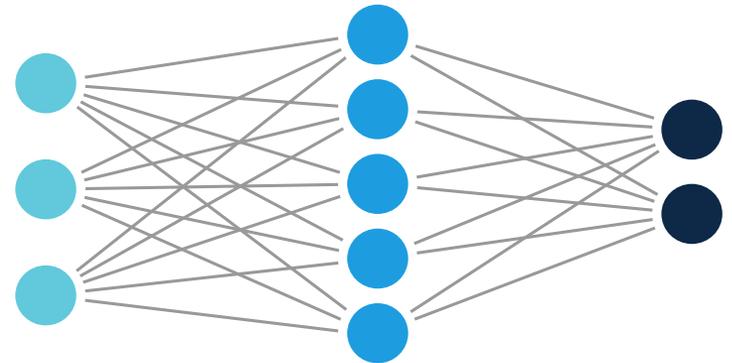
- Quantifying its mass, usually after applying some *grooming*;
- Inferring the number of quarks and gluons clustered inside it (*prongs*).



Examples:

Processes	Prongness	Classification
QCD	One-pronged (1P)	<b>Background</b>
$W/Z/H \rightarrow q\bar{q}$	Two-pronged (2P)	<b>Signal</b>
$t \rightarrow W^+b \rightarrow q\bar{q}b$	Three-pronged (3P)	
$S \rightarrow AA \rightarrow q\bar{q}q\bar{q}$	Four-pronged (4P)	

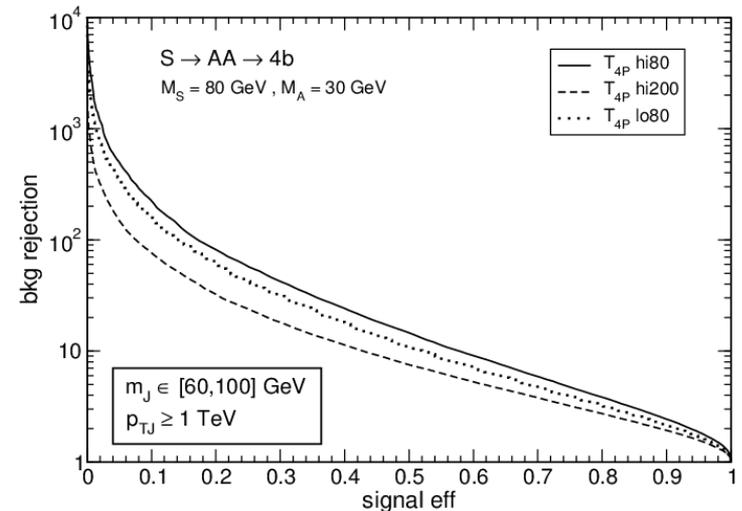
Here, this procedure relies on the training of **Neural Networks (NNs)**.



# MASS DECORRELATION

- The mass of a jet and the variables that encode its substructure are usually **correlated**;
- The mass decorrelation methods employed so far in supervised taggers leave a residual dependence of the results on the jet mass and transverse momentum training ranges. **Consequently...**

**Their performance drops when applied to kinematical regions different from those used to train them.**



J. A. Aguilar-Saavedra, B. Zaldívar, Eur. Phys. J. C 80, 6 (2020) 530

**Solution ??** 🤔

- **Build generic taggers, sensitive to any kind of jets**
- **Excellent performance for all jet masses**
- **Mass decorrelation**

# MASS UNSPECIFIC SUPERVISED TAGGING

Considering the **jet mass** and its **transverse momentum** varying over wide ranges, we make them input variables of a multivariate tool, together with **jet substructure observables**.

Mass Unspecific Supervised Tagging (MUST) for boosted jets

JHEP03(2021)012

J.A. Aguilar-Saavedra,<sup>a</sup> F.R. Joaquim<sup>b</sup> and J.F. Seabra<sup>b</sup>

<sup>a</sup>Departamento de Física Teórica y del Cosmos, Universidad de Granada, E-18071 Granada, Spain

<sup>b</sup>Departamento de Física and CFTP, Instituto Superior Técnico, Universidade de Lisboa, Av. Rovisco Pais 1, 1049-001 Lisboa, Portugal

Our MUST-inspired jet taggers have **19** input variables:

- **17 N-subjettiness observables** which characterise jet substructure,

$$\left\{ \tau_1^{(1/2)}, \tau_1^{(1)}, \tau_1^{(2)}, \dots, \tau_5^{(1/2)}, \tau_5^{(1)}, \tau_5^{(2)}, \tau_6^{(1)}, \tau_6^{(2)} \right\};$$

- **Jet's mass**,  $m_J$ ;
- **Jet's transverse momentum**,  $p_T$ .

**Note:**

Ratios of N-subjettiness variables will be denoted as

$$\tau_{mn} \equiv \frac{\tau_m^{(1)}}{\tau_n^{(1)}}.$$

All those variables should be **standardised** according to the SM background distributions.

# TRAINING SET GENERATION

	Background	Signal
Processes	$pp \rightarrow jj$	$pp \rightarrow ZS$ <ul style="list-style-type: none"> <li>All signal types: <math>Z \rightarrow \nu\nu</math></li> <li>2P: <math>S \rightarrow u\bar{u}, S \rightarrow b\bar{b}</math></li> <li>3P: <math>S \rightarrow F\nu; F \rightarrow udd, F \rightarrow udb</math></li> <li>4P: <math>S \rightarrow u\bar{u}u\bar{u}, S \rightarrow b\bar{b}b\bar{b}</math></li> </ul>
$p_T$ range	[200, 2200] GeV	
Mass ranges	$m_j \in [50, 250]$ GeV $M_{S,F} \in [30, 400]$ GeV $(M_S \leq p_T R / 2, R=0.8)$	

The decays of  $S$  and  $F$  are implemented with a flat matrix element (**to achieve generic taggers**).

# TAGGER PROPERTIES

Name	Types of events used in training	NN architecture	Output Layer
GenT	Background + 2P + 3P + 4P	2048 x 128	Sigmoid
GenT <sub>2P</sub>	Background + 2P	1028 x 64	Sigmoid
GenT <sub>3P</sub>	Background + 3P	1028 x 64	Sigmoid
GenT <sub>4P</sub>	Background + 4P	1028 x 64	Sigmoid
Prongness selection tagger	2P + 3P + 4P	2048 x 128	Softmax

To evaluate the performance of GenT, GenT<sub>2P</sub>, GenT<sub>3P</sub> and GenT<sub>4P</sub> we use the Area Under the ROC curve (AUC) whereas the Prongness selection tagger is evaluated by measuring its accuracy.

- All our NNs use the Rectified Linear Unit (ReLU) activation function;
- The optimisation of GenT, GenT<sub>2P</sub>, GenT<sub>3P</sub> and GenT<sub>4P</sub> (Prongness selection tagger) rely on the binary (categorical) cross-entropy;
- The Adam optimiser is applied to all NNs.

# BENEFIT OF MUST TAGGERS

**Non-MUST taggers:**

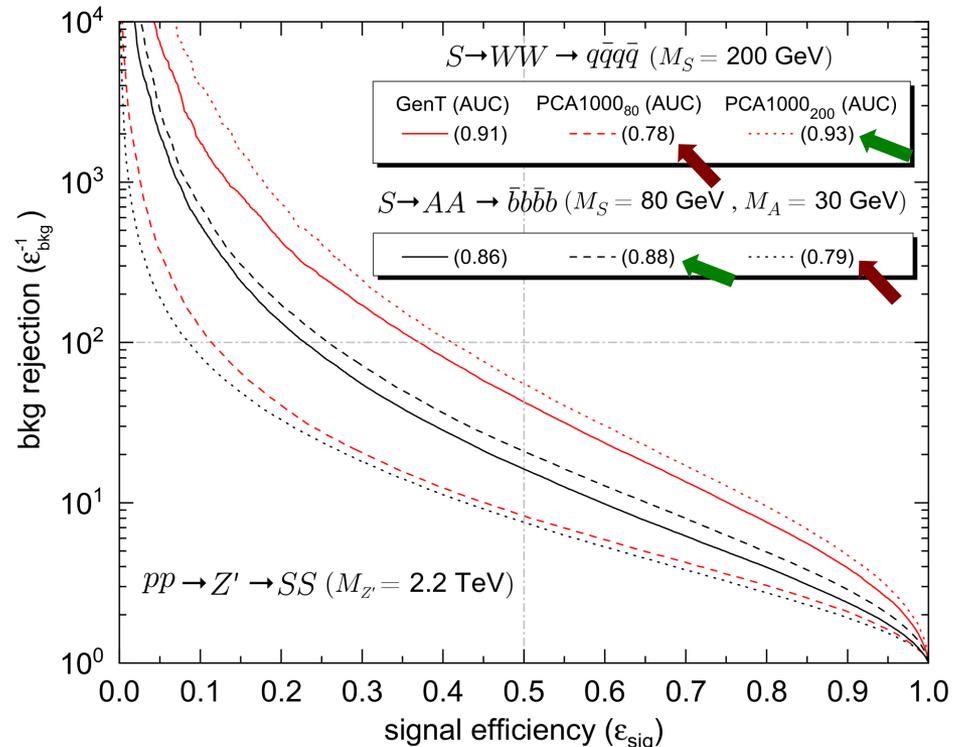
**PCA1000<sub>80</sub>**: Trained for  $p_T \geq 1.0$  TeV and on the mass interval  $m_J \in [60, 100]$  GeV.

**PCA1000<sub>200</sub>**: Trained on the same region of momentum but in a different mass interval,  $m_J \in [160, 240]$  GeV.

**Principal Component Analysis (PCA)** is used in both taggers to perform mass decorrelation.

✓ These taggers perform slightly better on a mass region close to the one where they were trained...

✗ ... but are much less efficient when applied to masses out of the training region.

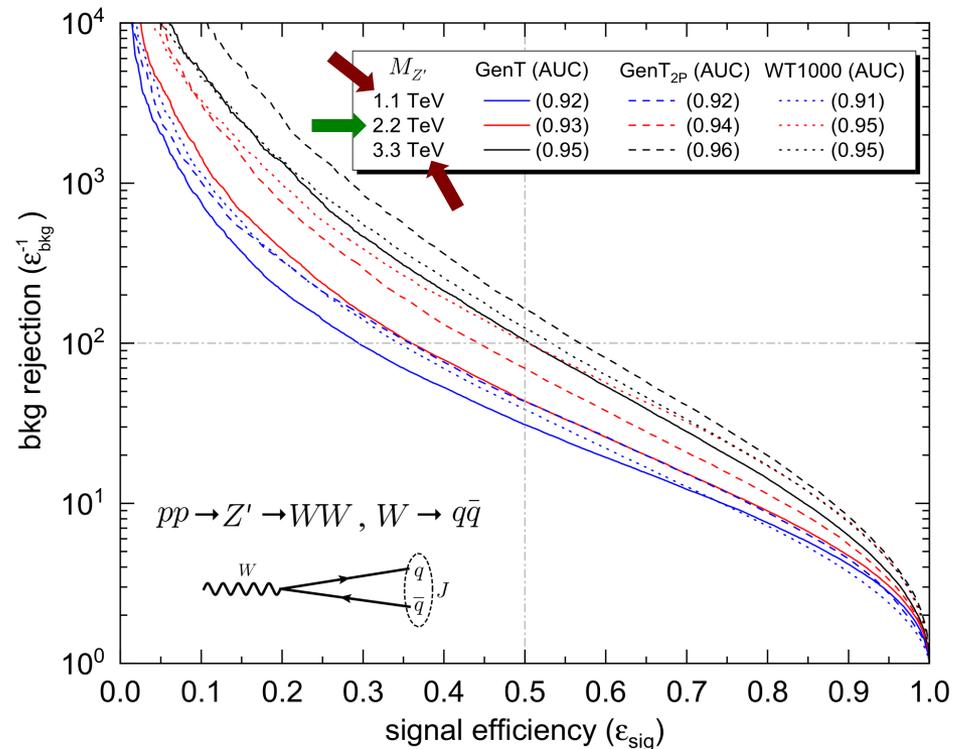


# BENEFIT OF MUST TAGGERS

Non-MUST tagger:

**WT1000:** PCA-decorrelated tagger trained with  $W$  jets obtained from  $Z' \rightarrow WW$  ( $M_{Z'} = 2.2$  TeV) and QCD jets with  $p_T \geq 1$  TeV and  $m_J \in [60, 100]$  GeV.

- ✓ It performs slightly better than  $\text{GenT}_{2P}$  for  $p_T \geq 1$  TeV and  $m_J \in [60, 100]$  GeV.
- ✗ It performs slightly worse than  $\text{GenT}_{2P}$  for  $p_T \geq 500$  GeV and  $p_T \geq 1.5$  TeV.

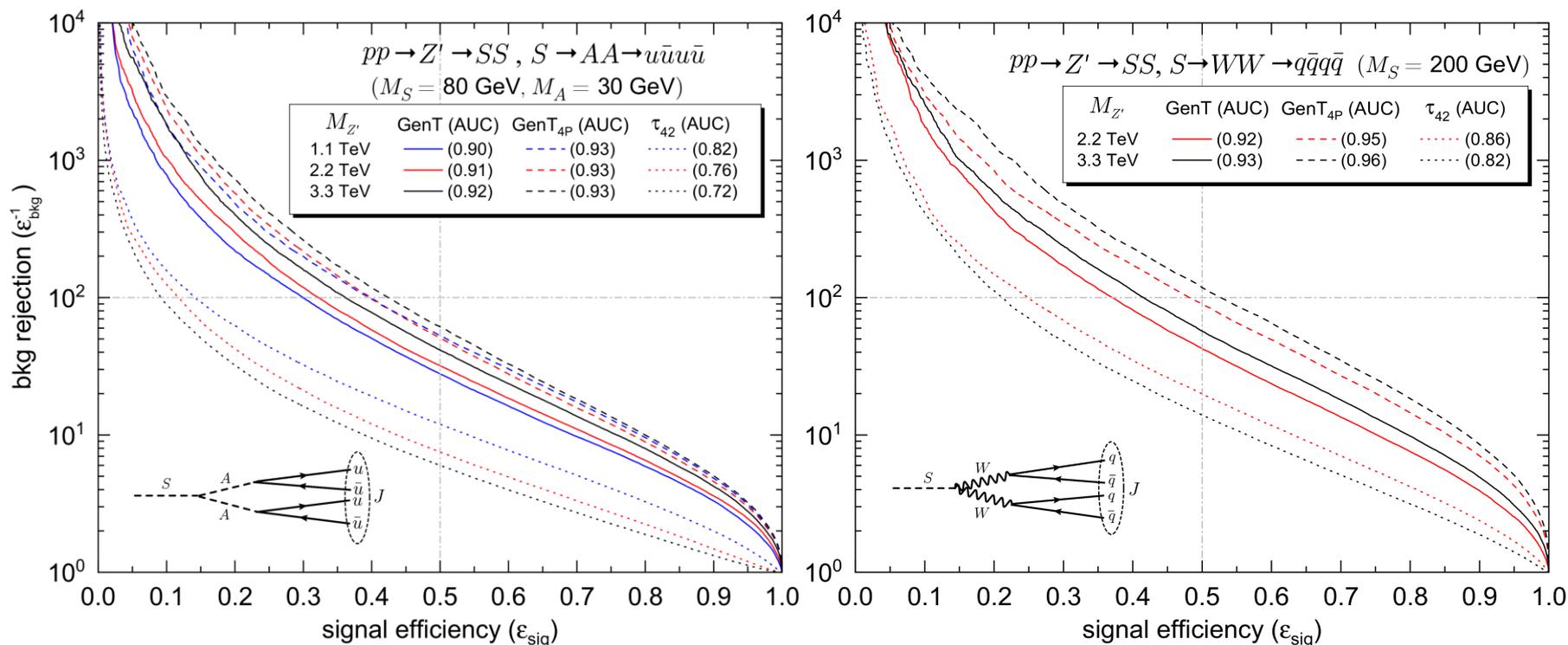


**$\text{GenT}_{2P}$  is nearly optimal for  $W$  jets**

# TAGGER PERFORMANCE (4P SIGNALS)

Background: Quark and gluon jets generated in  $pp \rightarrow Zq$ ,  $pp \rightarrow Zg$ , with  $Z \rightarrow \nu\nu$

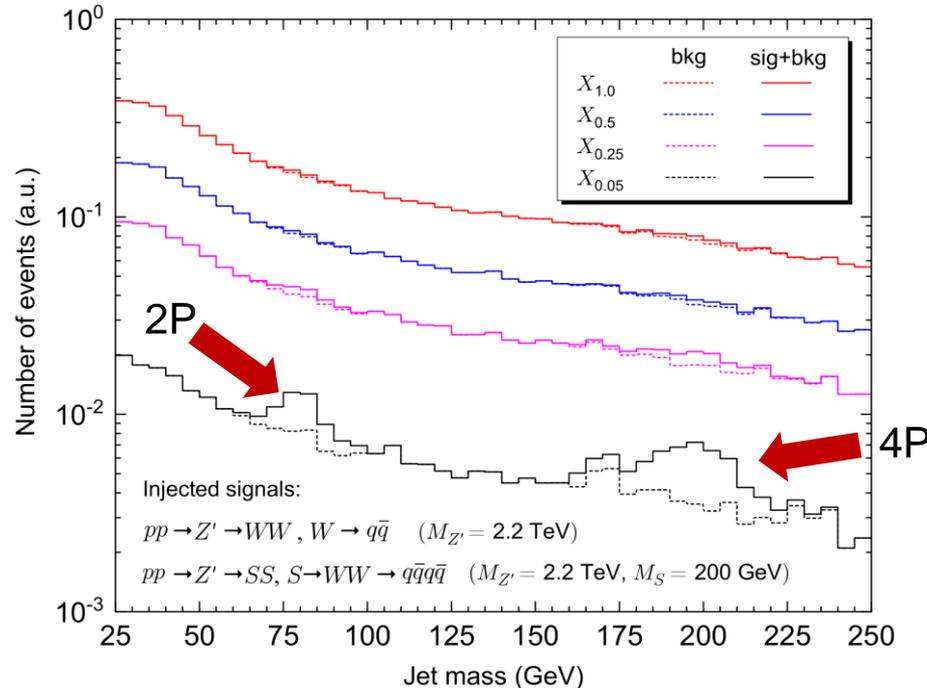
- $p_T \geq 0.5, 1.0, 1.5$  TeV for  $M_{Z'} = 1.1, 2.2, 3.3$  TeV, respectively;



- The performance of GenT and GenT<sub>4P</sub> is **significantly better** than that of  $\tau_{42}$ .
- The performance **improves** as  $M_{Z'}$  increases

# MASS DECORRELATION

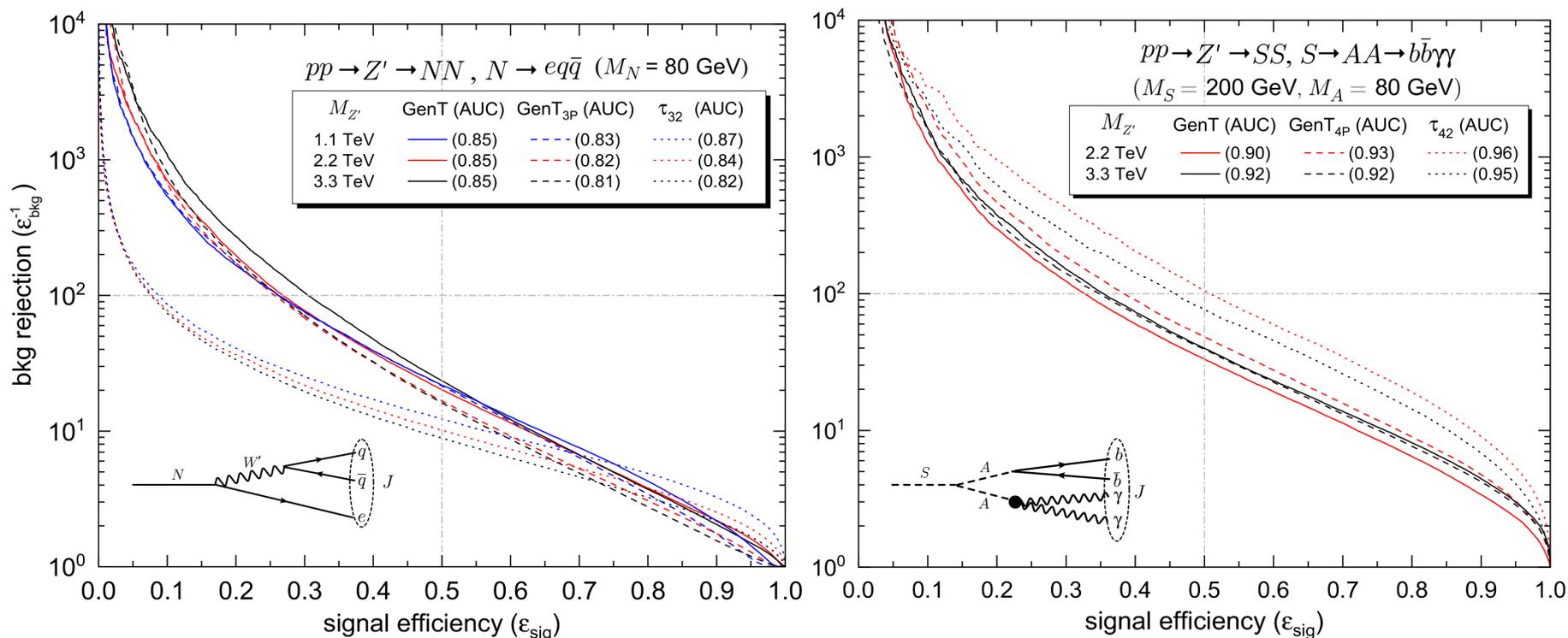
- Defining  $\rho = 2 \log(m_J / p_T)$ , we compute at each bin of a 2D grid  $(\rho, p_T)$  the 5%, 25% and 50% percentiles of the NN score ( $X_{0.05}$ ,  $X_{0.25}$  and  $X_{0.5}$  respectively);



- This varying threshold preserves the SM background distribution and the injected signals show up when the cut is sufficiently tight.

**Our generic taggers also provide a perfect solution to the mass correlation problem.**

# JETS NOT USED TO TRAIN MUST TAGGERS



- **MUST taggers can detect unseen signals with good efficiency;**
- AUC is not good to evaluate the performance of taggers for neutrino jets;
- Simpler multivariate methods like **logistic regression** may achieve better performance for stealth boson jets with two photons in the final state.

J. A. Aguilar-Saavedra, B. Zaldivar; Eur. Phys. J. C 80, 6 (2020) 530

# IDENTIFICATION OF NEW PHYSICS SIGNALS

Using the [prongness selection tagger](#), we apply the following classification criteria in the four benchmark examples below:

$$\begin{cases} 2P, & \text{if } P_{2P} \geq 0.5 \\ 3P, & \text{if } P_{3P} \geq 0.5 \\ 4P, & \text{if } P_{4P} \geq 0.5 \\ \text{Undefined,} & \text{otherwise} \end{cases}$$

## Benchmark 1 (4P)

$$\begin{aligned} Z' &\rightarrow SS, \\ S &\rightarrow AA \rightarrow b\bar{b}b\bar{b}, \\ M_{Z'} &= 2.2 \text{ TeV}, \\ M_S &= 80 \text{ GeV}, \\ M_A &= 30 \text{ GeV} \end{aligned}$$

## Benchmark 2 (2P)

$$\begin{aligned} Z' &\rightarrow AA, \\ A &\rightarrow b\bar{b}, \\ M_{Z'} &= 2.2 \text{ TeV}, \\ M_A &= 80 \text{ GeV} \end{aligned}$$

## Benchmark 3 (4P)

$$\begin{aligned} Z' &\rightarrow SS, \\ S &\rightarrow WW \rightarrow q\bar{q}q\bar{q}, \\ M_{Z'} &= 3.3 \text{ TeV}, \\ M_S &= 200 \text{ GeV} \end{aligned}$$

## Benchmark 4 (2P)

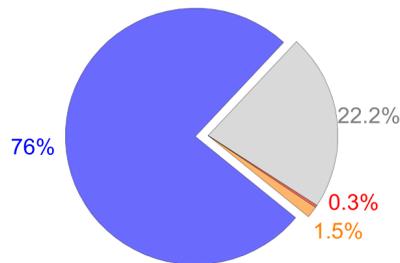
$$\begin{aligned} Z' &\rightarrow AA, \\ A &\rightarrow u\bar{u}, \\ M_{Z'} &= 3.3 \text{ TeV}, \\ M_A &= 200 \text{ GeV} \end{aligned}$$

# IDENTIFICATION OF NEW PHYSICS SIGNALS

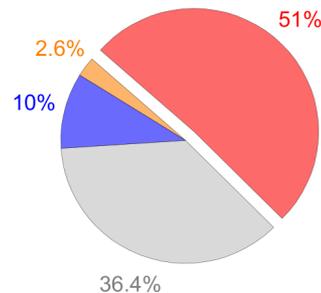
Using the **prongness selection tagger**, we apply the following classification criteria in the four benchmark examples below:

$$\begin{cases} 2P, & \text{if } P_{2P} \geq 0.5 \\ 3P, & \text{if } P_{3P} \geq 0.5 \\ 4P, & \text{if } P_{4P} \geq 0.5 \\ \text{Undefined,} & \text{otherwise} \end{cases}$$

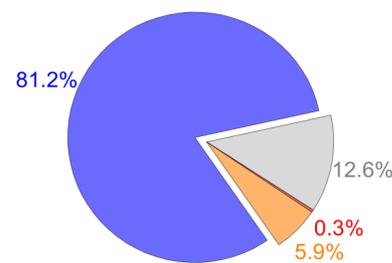
Benchmark 1 (4P)



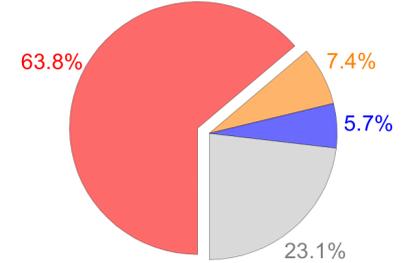
Benchmark 2 (2P)



Benchmark 3 (4P)



Benchmark 4 (2P)



Classification: ■ 2P ■ 3P ■ 4P ■ Undefined

- The fraction of correctly identified jets is several times larger than that of misidentified ones;
- Mistag rates can be further reduced by raising the value of the threshold that separates undefined jets from the classified ones.

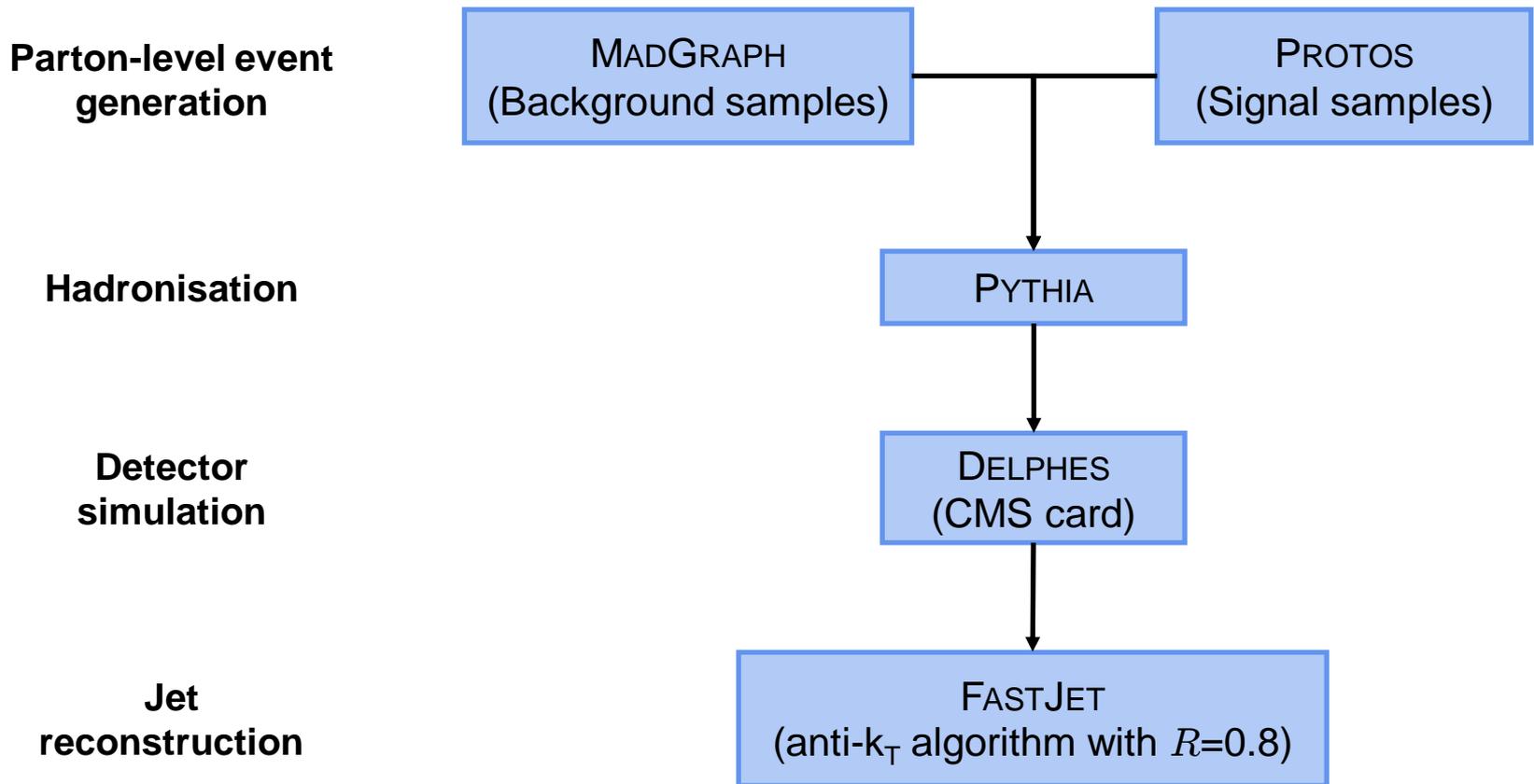
# CONCLUDING REMARKS

- We introduced the method of **MUST** for multi-pronged jets;
- Taggers built upon MUST keep an excellent performance across a very wide  $m_J$  and  $p_T$  range;
- Our taggers are sensitive to any kind of multi-pronged jets, outperforming simple variables;
- Mass decorrelation can easily be implemented using the varying threshold method;
- MUST taggers can achieve good performances on signals for which they were not trained;
- The MUST concept can also be applied to selection taggers that can determine the prongness of signal jets.

**Thank you!**

# Backup slides

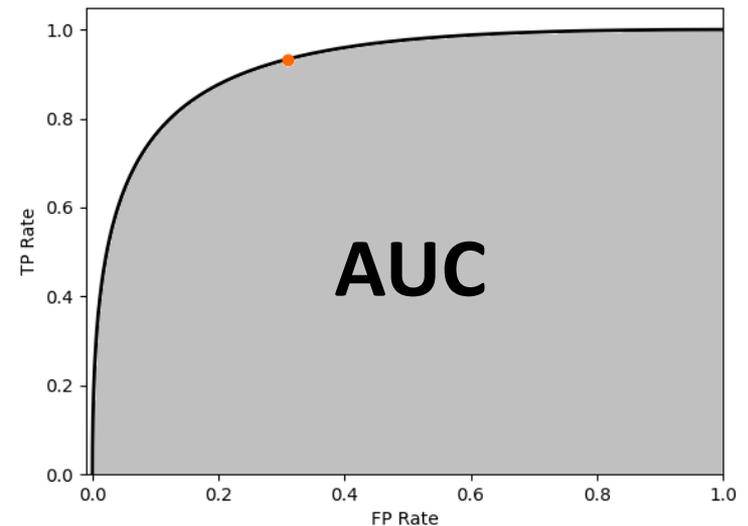
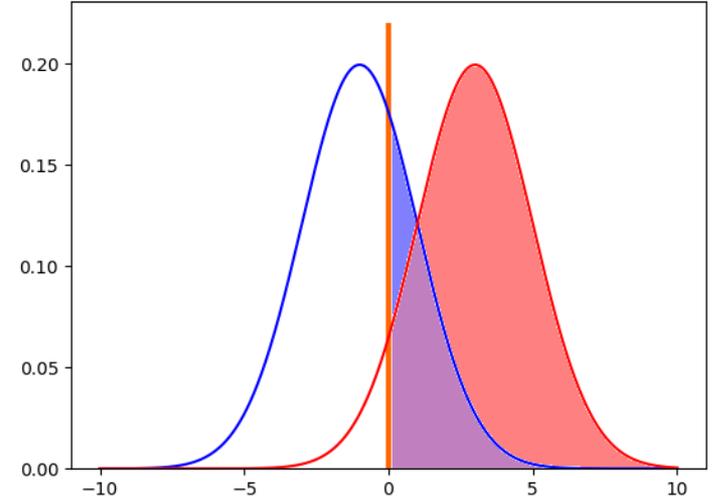
# EVENT SIMULATION AND RECONSTRUCTION



# ROC CURVES

- In a binary classification task, there is always a **threshold** separating the two classes;
- The fraction of **True Positives (TP)** and **False Positives (FP)** for all possible thresholds defines the classifier's **ROC curve**;
- The **Area Under the ROC curve (AUC)** is often used to evaluate the performance of the classifier (it assumes a value between 0 and 1).

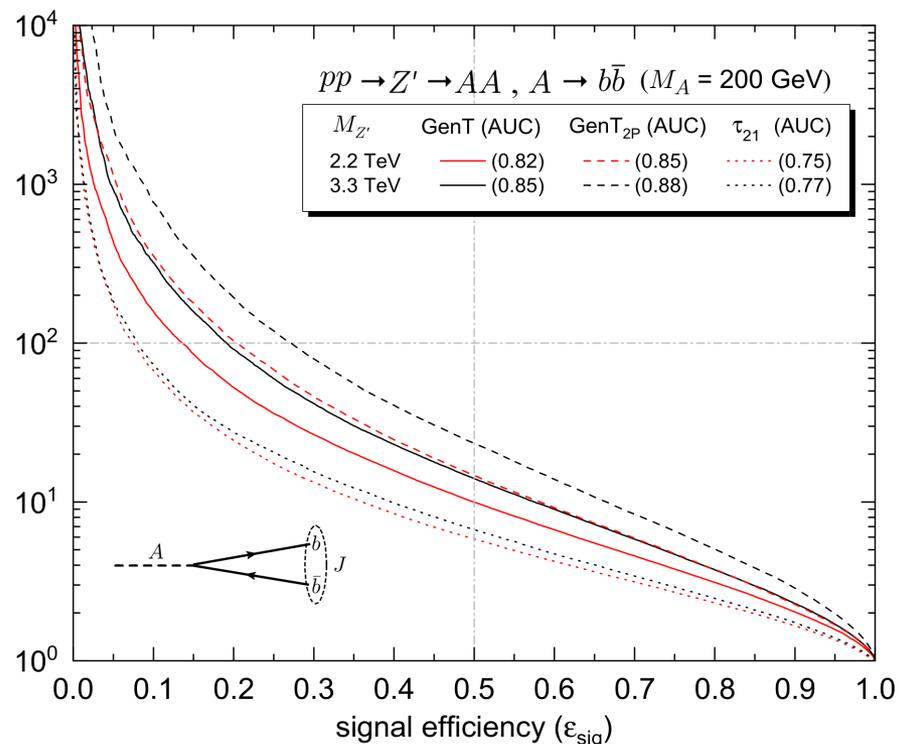
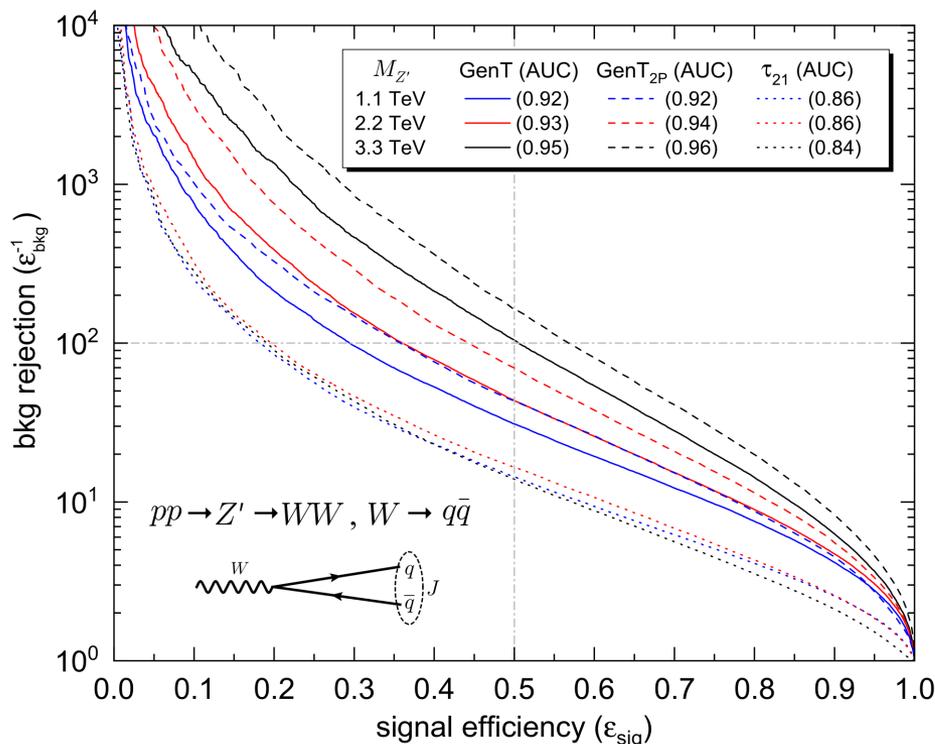
**Note:** In our results, we represent the ROC curves on the plane  $(\epsilon_{\text{sig}}, \epsilon_{\text{bkg}}^{-1})$ . Considering Background and Signal events as being Negative and Positive, respectively,  $\epsilon_{\text{sig}} = \text{TP Rate}$  and  $\epsilon_{\text{bkg}} = \text{FP Rate}$ .



# TAGGER PERFORMANCE (2P SIGNALS)

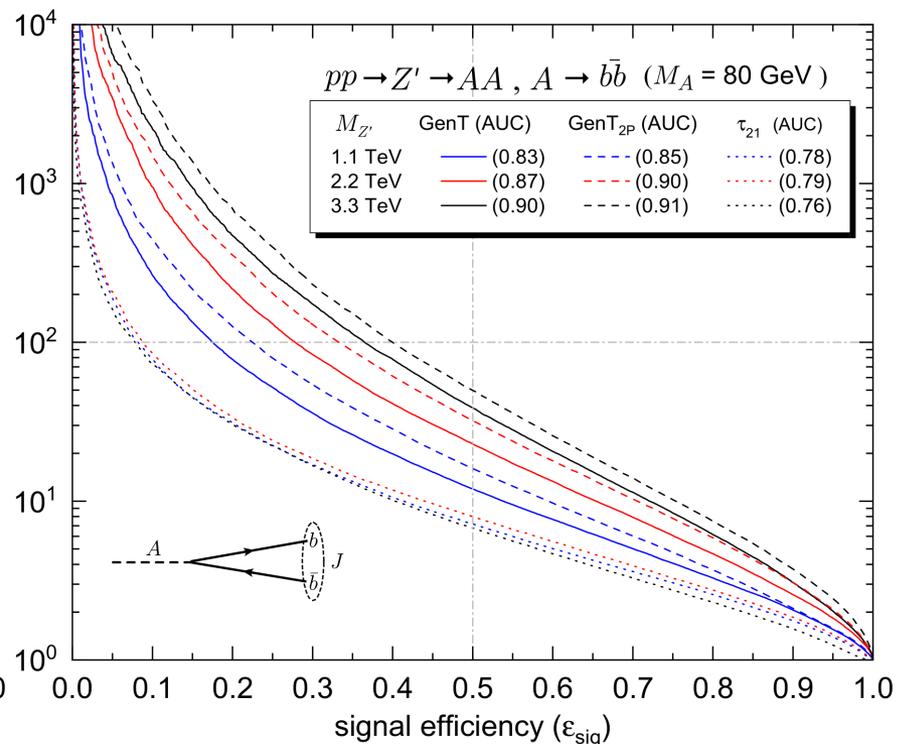
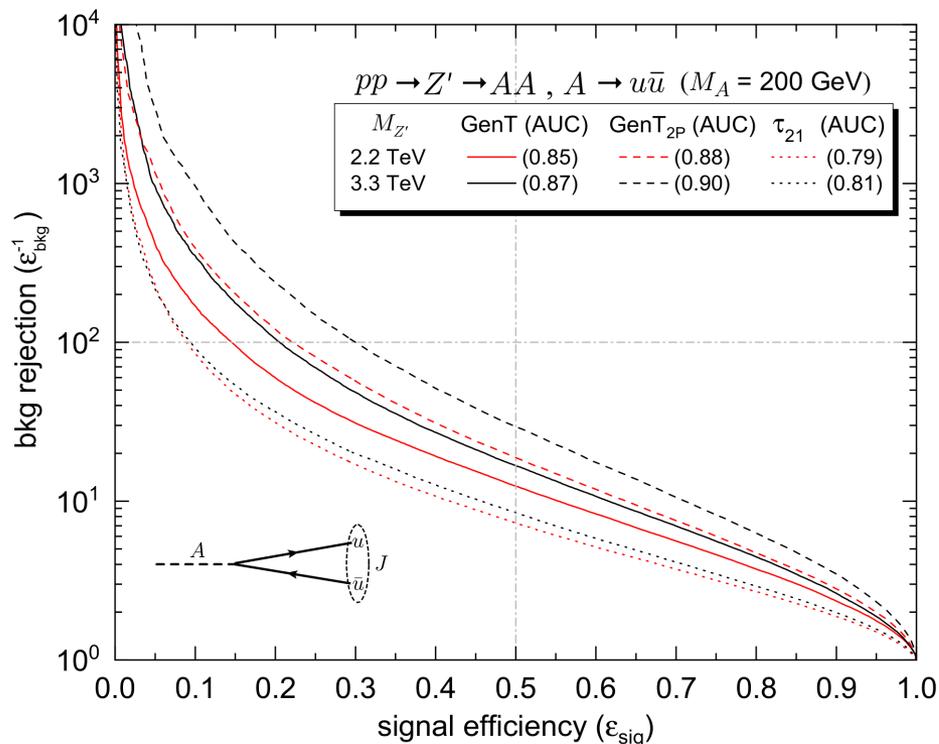
Background: Quark and gluon jets generated in  $pp \rightarrow Zq$ ,  $pp \rightarrow Zg$ , with  $Z \rightarrow \nu\nu$

- $p_T \geq 0.5, 1.0, 1.5$  TeV for  $M_{Z'} = 1.1, 2.2, 3.3$  TeV, respectively;

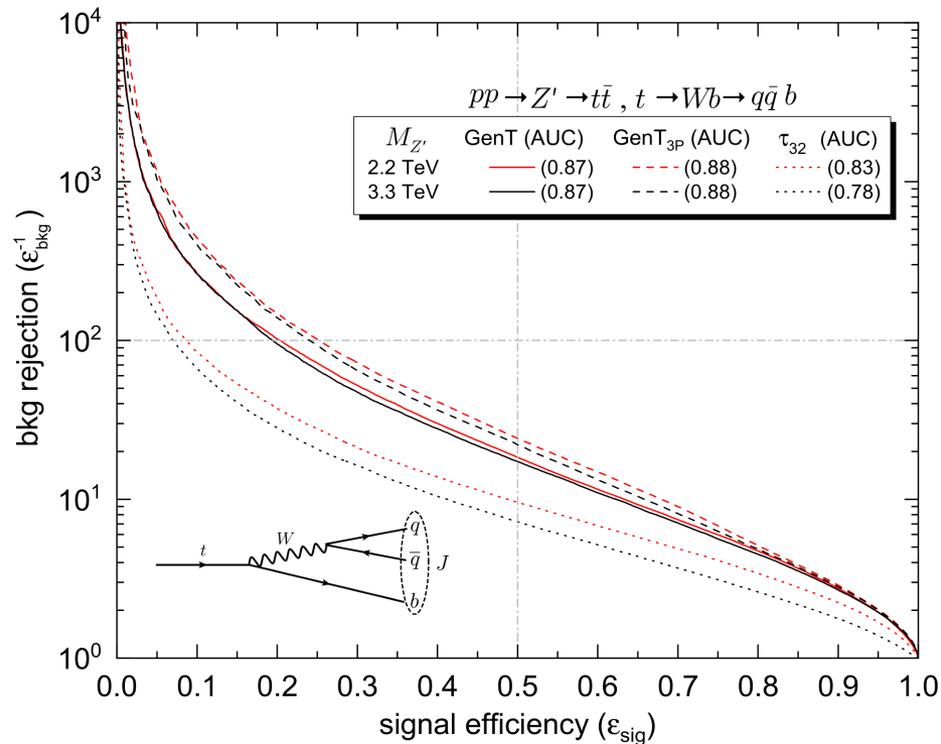


- In general, GenT and GenT<sub>2P</sub> perform better than the commonly used ratio  $\tau_{21}$ ;
- The performance improves as  $M_{Z'}$  increases.

# TAGGER PERFORMANCE (2P SIGNALS)



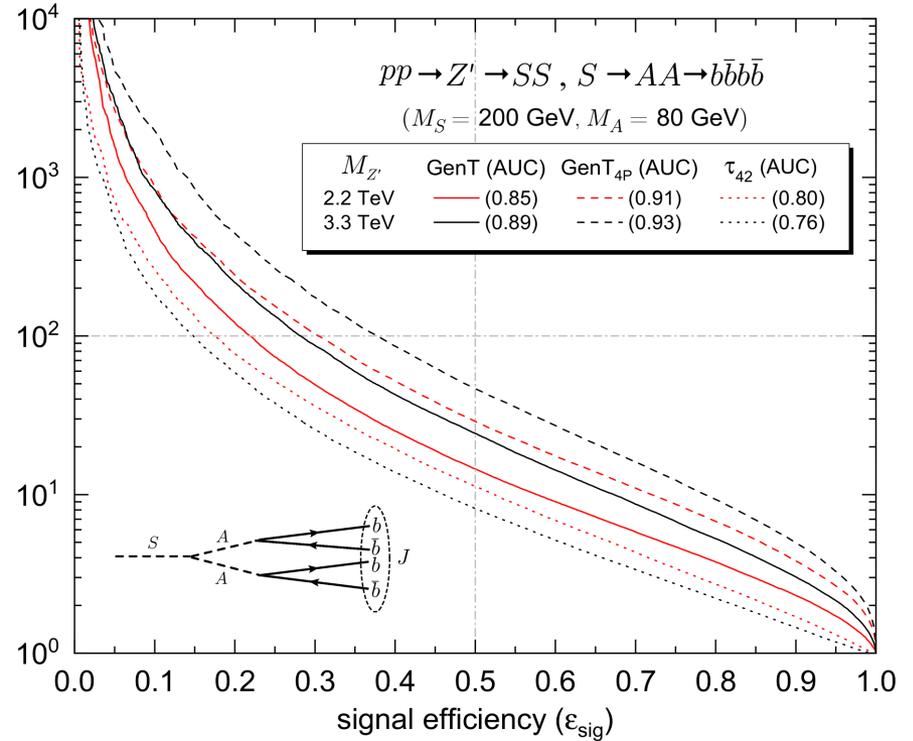
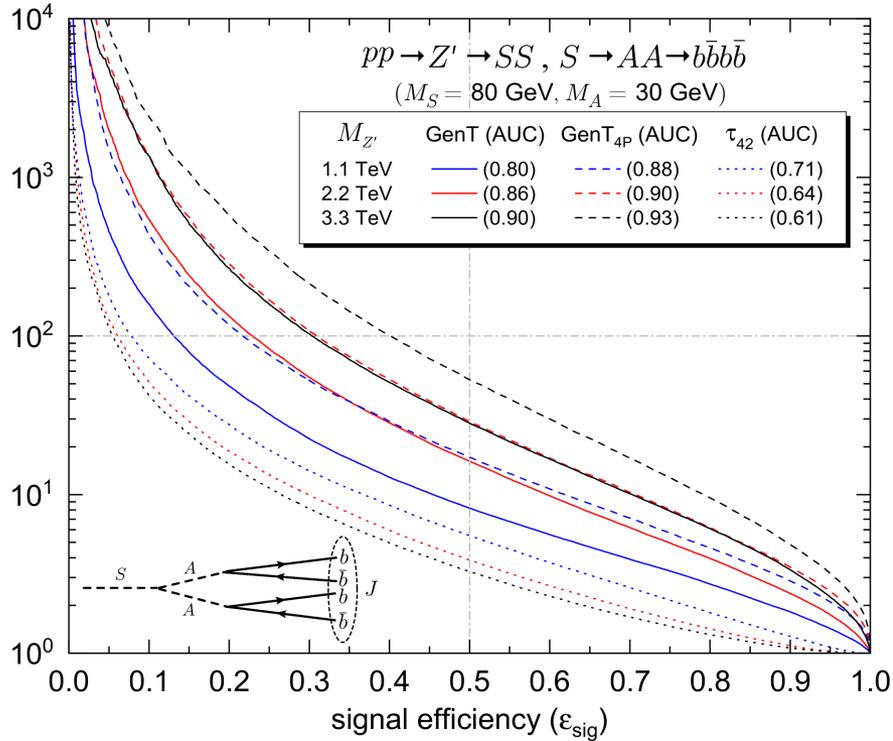
# TAGGER PERFORMANCE (3P SIGNAL)



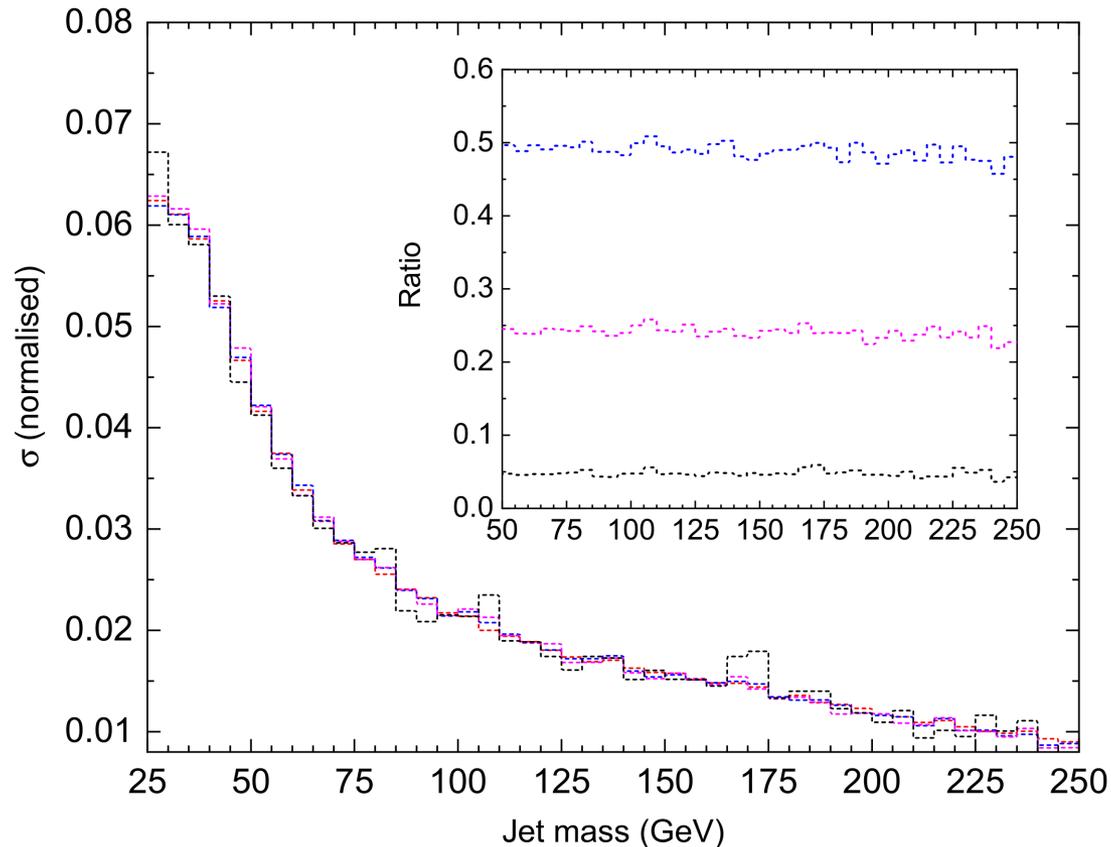
- Although GenT and GenT<sub>3P</sub> perform well on top quark jets and **would not miss those signals**, fully-dedicated top taggers perform better.

e.g. S. Macaluso, D. Shih; JHEP 10 (2018) 121

# TAGGER PERFORMANCE (4P SIGNALS)



# ANOTHER PLOT FOR MASS DECORRELATION



- Main plot – Normalised background distributions before and after cuts
- Inner plot – Ratios of distributions after/before cuts

# MORE ABOUT STEALTH BOSONS

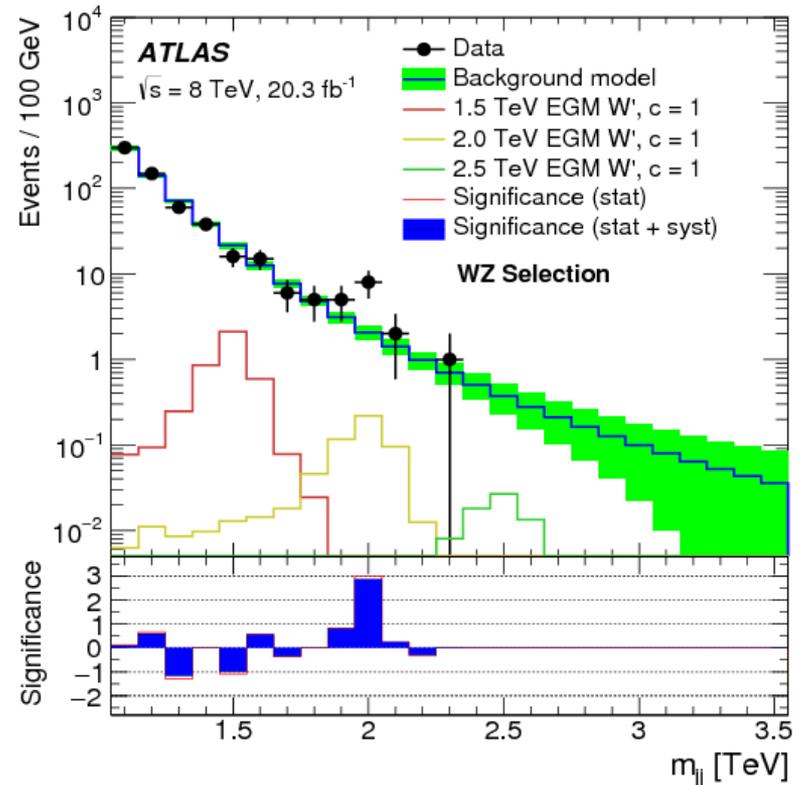
- Stealth bosons are relatively light boosted particles with a cascade decay:

$$S \rightarrow AA \rightarrow q\bar{q}q\bar{q}$$



A particles can be weak bosons W, Z, a Higgs boson or new relatively light (pseudo-)scalars.

- Heavy resonances decaying into two such stealth bosons, or one plus a W/Z boson, may offer an explanation for **small excesses** found in hadronic diboson resonance searches near an invariant mass of 2 TeV (example on the right).



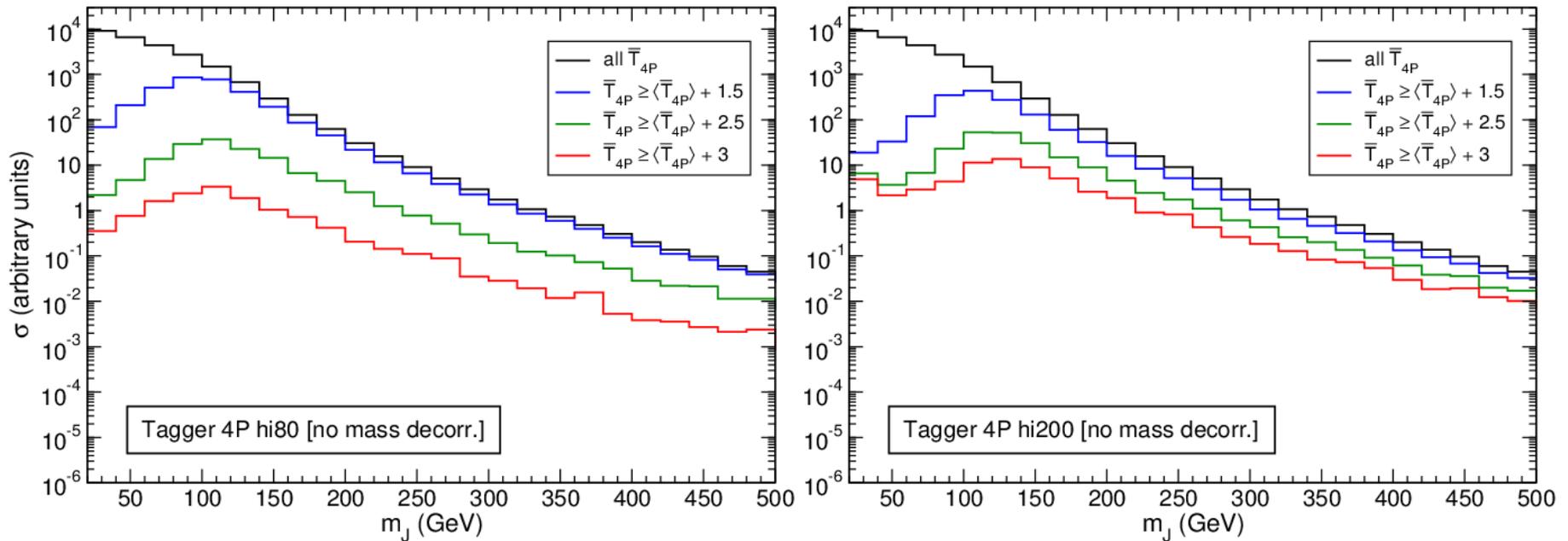
(For more information about this topic, see:

J. A. Aguilar-Saavedra; Eur. Phys. J. C 77, 10 (2017) 703)

G. Aad *et al.* [ATLAS]; JHEP 12 (2015) 55

# TAGGERS WITHOUT MASS DECORRELATION

Jet mass spectrum for QCD background produced by 4P taggers with no prior mass decorrelation



J. A. Aguilar-Saavedra, B. Zaldívar, Eur. Phys. J. C 80, 6 (2020) 530

- The peak-like structure produced near 100 GeV is not in any case related to the design mass interval.

# N-SUBJETTINESS OBSERVABLES

- The N-subjettiness observable  $\tau_N^{(\beta)}$  is a measure of the radiation about  $N$  axes in the jet, specified by an angular exponent  $\beta > 0$ ,

$$\tau_N^{(\beta)} = \frac{1}{p_T} \sum_{i \in \text{jet}} p_{Ti} \min \left\{ R_{1i}^\beta, R_{2i}^\beta, \dots, R_{Ni}^\beta \right\}$$

Transverse momentum  
of particle  $i$  in the jet

Angular distance between  
particle  $i$  and axis  $N$  in the jet

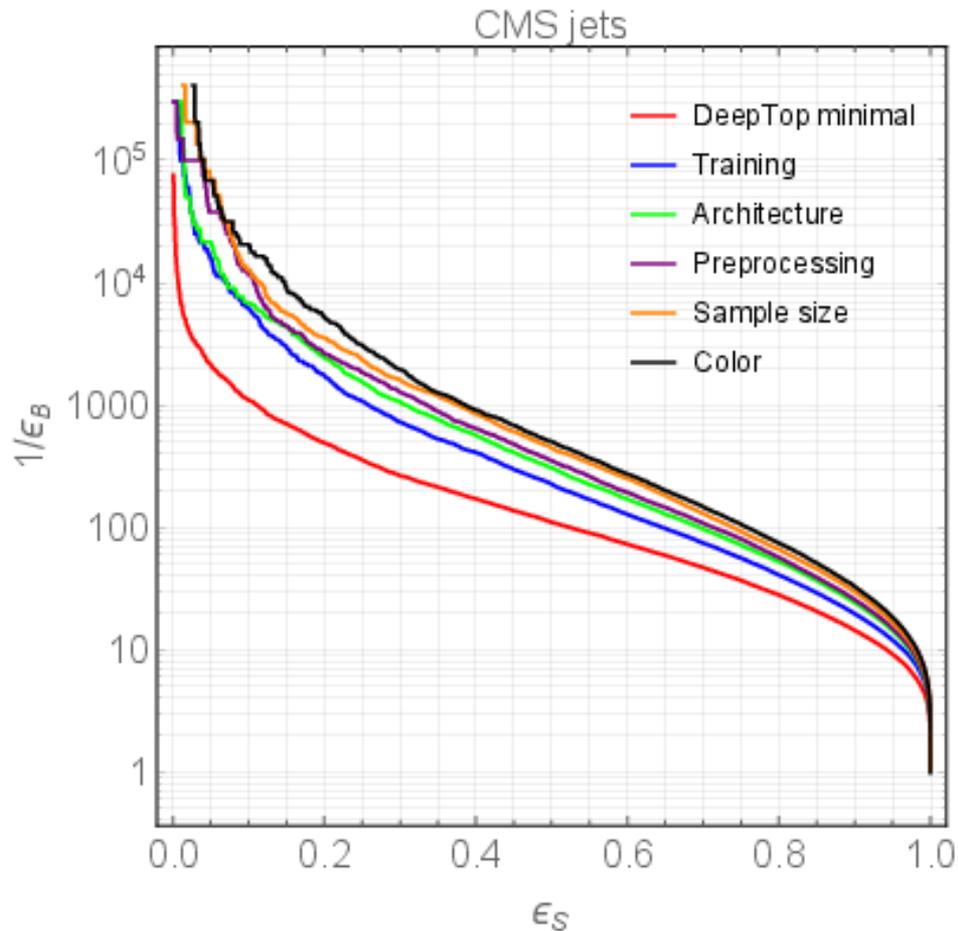
- The coordinates of the  $M$ -Body phase space can be defined by  $(M - 1)$  transverse momentum fractions and  $(2M - 3)$  angles, so we need  **$(3M - 4)$**  N-subjettiness observables to completely specify the coordinates of that space:

$$\left\{ \tau_1^{(0.5)}, \tau_1^{(1)}, \tau_1^{(2)}, \tau_2^{(0.5)}, \tau_2^{(1)}, \tau_2^{(2)}, \dots, \tau_{M-2}^{(0.5)}, \tau_{M-2}^{(1)}, \tau_{M-2}^{(2)}, \tau_{M-1}^{(1)}, \tau_{M-1}^{(2)} \right\}$$

(For more information about this topic, see [K. Datta, A. Larkoski; JHEP 06 \(2017\) 73](#))

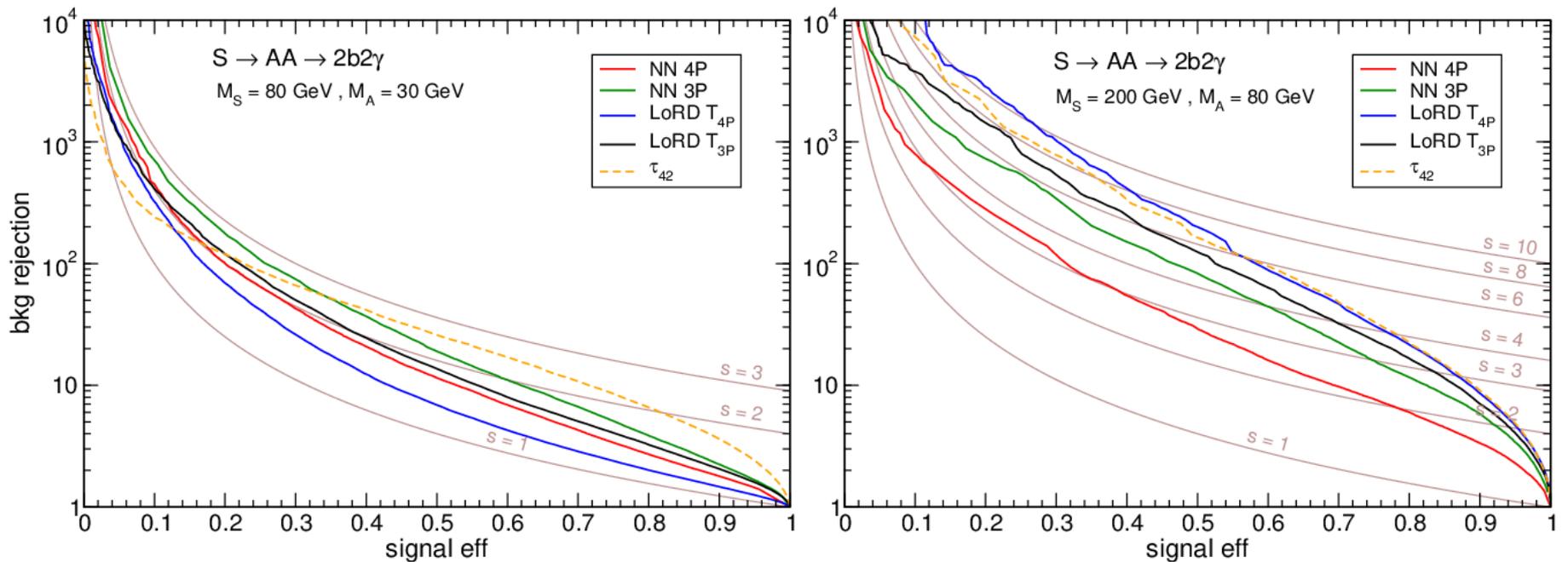
# TOPDEEP TAGGER

Performance of DeepTop tagger (after several improvements)  
discriminating top quark jets



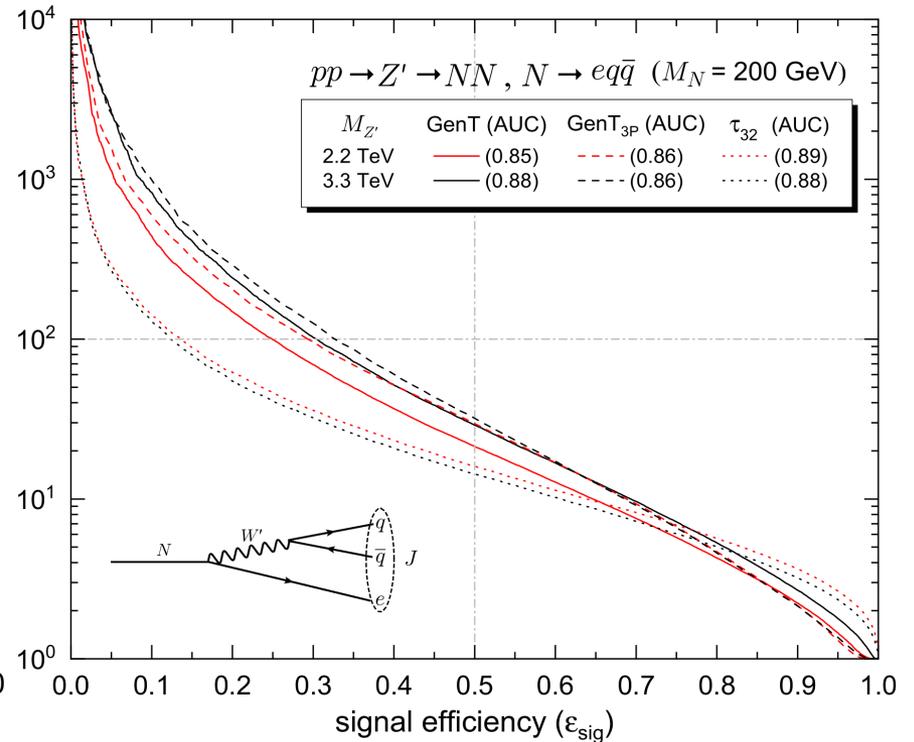
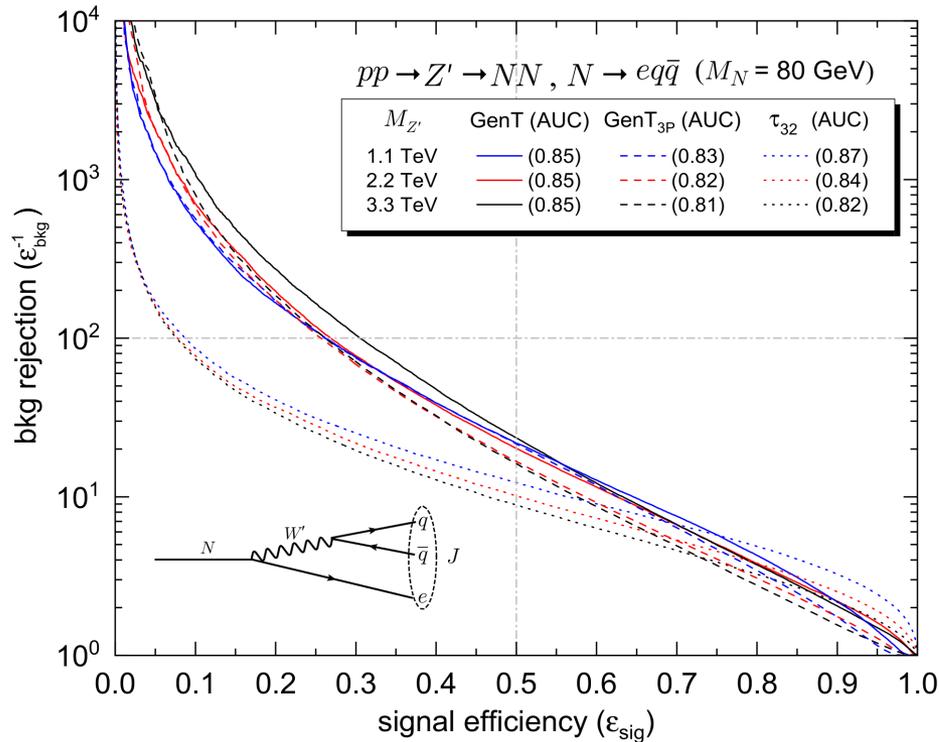
S. Macaluso, D. Shih; JHEP 10 (2018) 121

## Performance of Logistic Regression Design (LoRD) classifying jets with two hard photons



J. A. Aguilar-Saavedra, B. Zaldivar, Eur. Phys. J. C 80, 6 (2020) 530

# JETS NOT USED TO TRAIN MUST-TAGGERS



# JETS NOT USED TO TRAIN MUST-TAGGERS

