

Interpretable Deep Learning for Two-Prong Jet Classification with Jet Spectra

Sung Hak Lim

Theory Center, KEK



BOOST 2019, MIT

July 2019

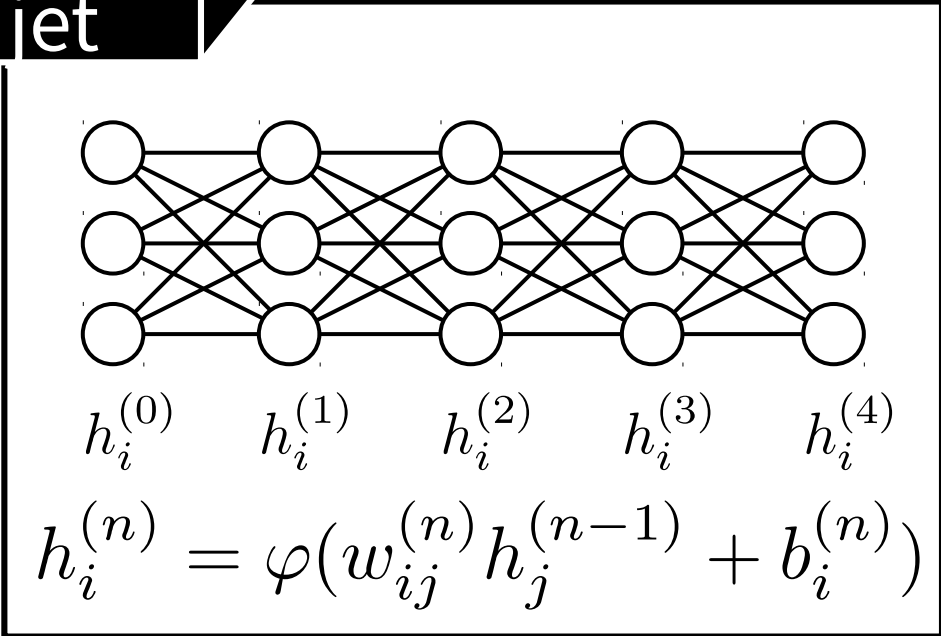
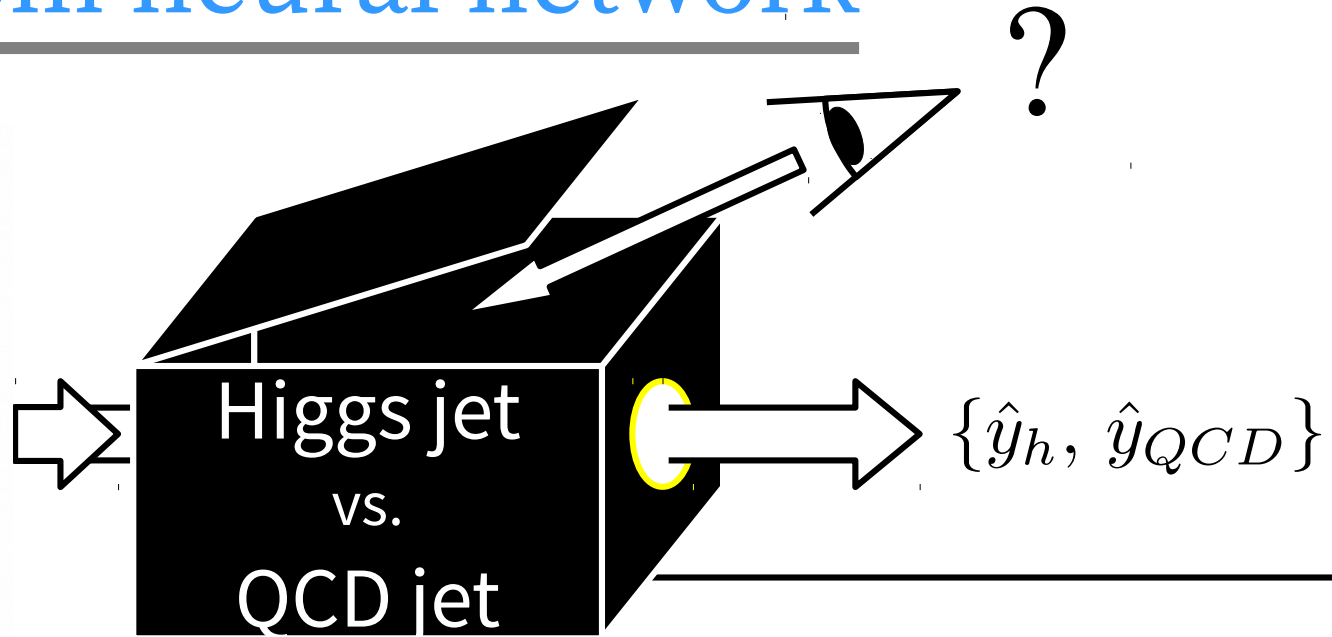
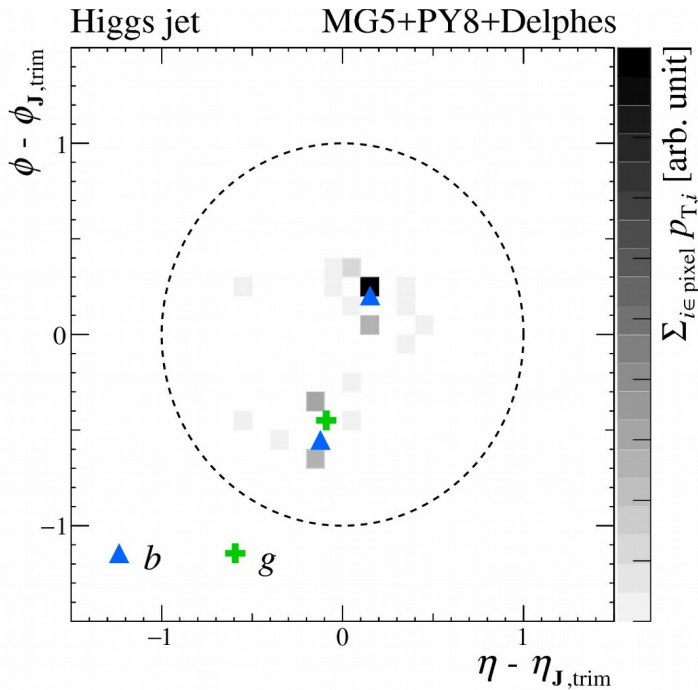


S. H. Lim, M. M. Nojiri, arXiv:1807.03312, JHEP10(2018)181.

A. Chakraborty, **S. H. Lim**, M. M. Nojiri, arXiv:1904.02092, JHEP07(2019)135.

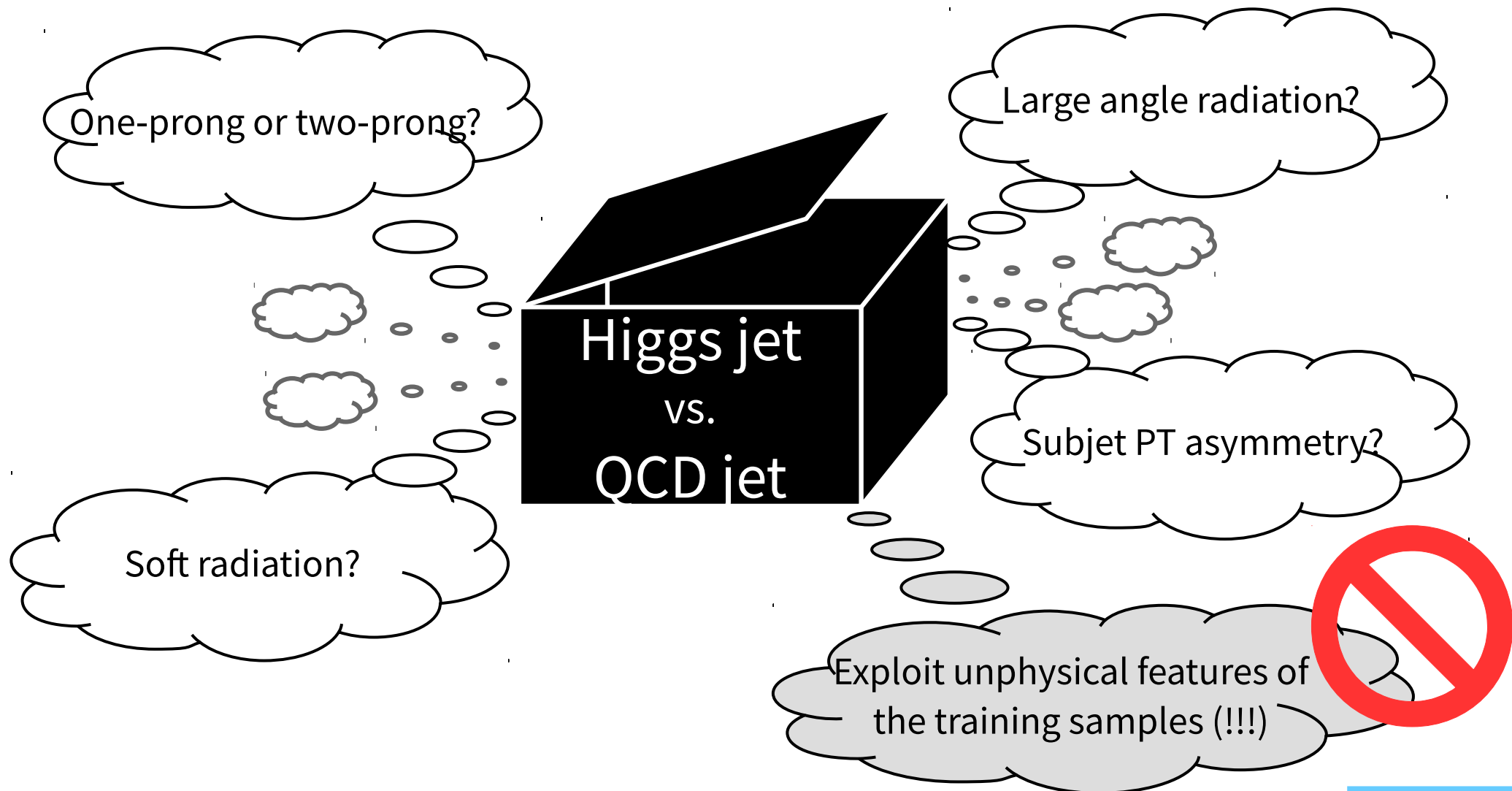
A. Chakraborty, **S. H. Lim**, M. M. Nojiri, M. Takeuchi, in preparation.

Difficulties on understanding the results from neural network



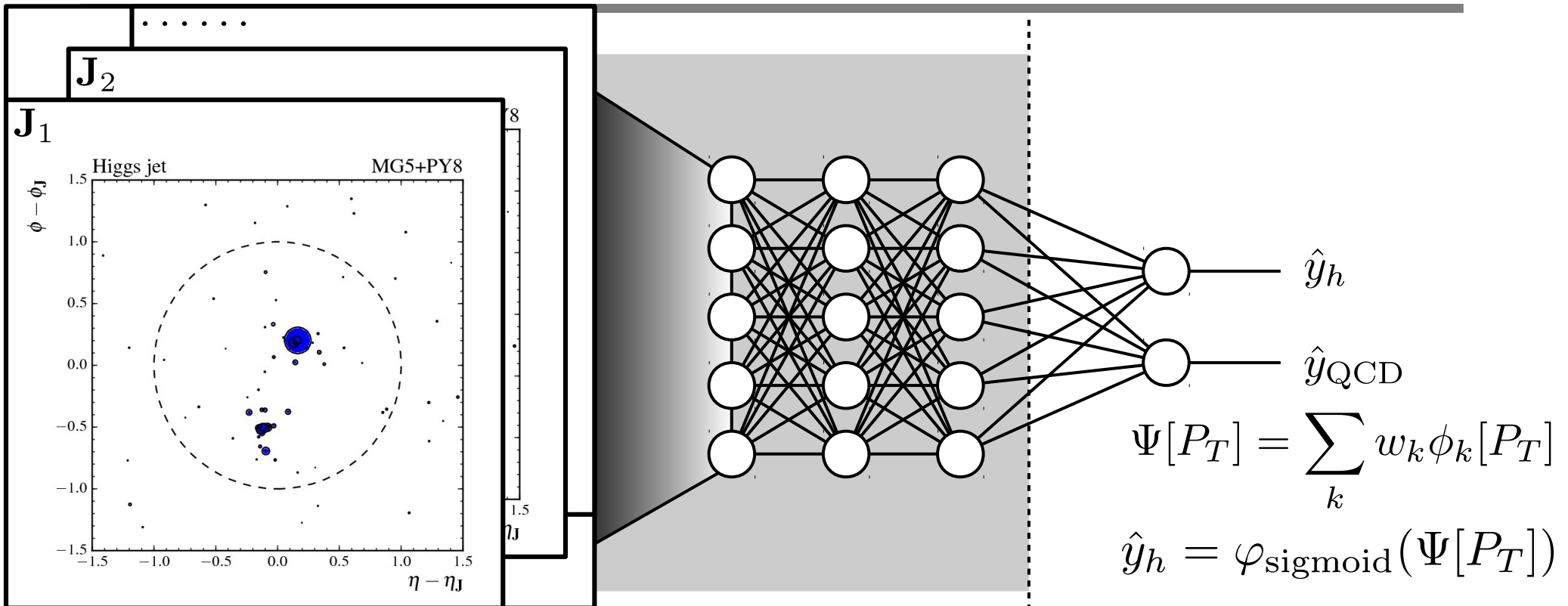
Neural network is often considered as a **black box** because studying its internal information barely gives you an insight about the decision...

Difficulties on understanding the results from neural network



We also want to know **the reasoning** behind the decision!

Basic Structure of a Neural Network Classifier



$$\Psi[P_T] = \sum_k w_k \phi_k[P_T]$$

$$\hat{y}_h = \varphi_{\text{sigmoid}}(\Psi[P_T])$$

Logistic regression

$$P_{T,a}(\vec{R}) = \sum_{i \in J_a} p_{T,i} \delta(\vec{R} - \vec{R}_i)$$

Inputs: energy flow

Lots of inputs

$$\phi_k = \Phi_k[P_T]$$

Feature map

Too many parameters

Talk of Gregor and David for quick review

- N-subjettiness, ECFs...
- MLP, CNN, ParticleNet...
- EFPs, EFN, JUNIPR, ...

The decision boundary is highly nonlinear.
It is hard to **interpret** the NN itself...

Functional Taylor Expansion

Let us consider the “functional Taylor expansion” of the classifier.

$$\Phi[P_{T,a}] = w^{(0)} + \int d\vec{R} P_{T,a}(\vec{R}) w_a^{(1)}(\vec{R}) + \frac{1}{2!} \int d\vec{R}_1 d\vec{R}_2 P_{T,a}(\vec{R}_1) P_{T,b}(\vec{R}_2) w_{ab}^{(2)}(\vec{R}_1, \vec{R}_2) + \dots$$

If we use only relative distance between constituents,

the first nontrivial term is

$$\Phi[P_{T,a}] = \int dR S_{2,ab}(R) w_{ab}^{(2)}(R) + \dots$$

$$S_{2,ab}(R) = \int d\vec{R}_1 d\vec{R}_2 P_{T,a}(\vec{R}_1) P_{T,b}(\vec{R}_2) \delta(R - R_{12})$$

$$w^{(0)} + p_{T,J_a} w_a^{(1)}$$

Reduce the dimension of inputs

[Length/bin width]² → [Length/bin width]

Two-point correlation between constituents at distance R

Two-Point Correlation Spectrum: KEK

Trimmed Spectrum

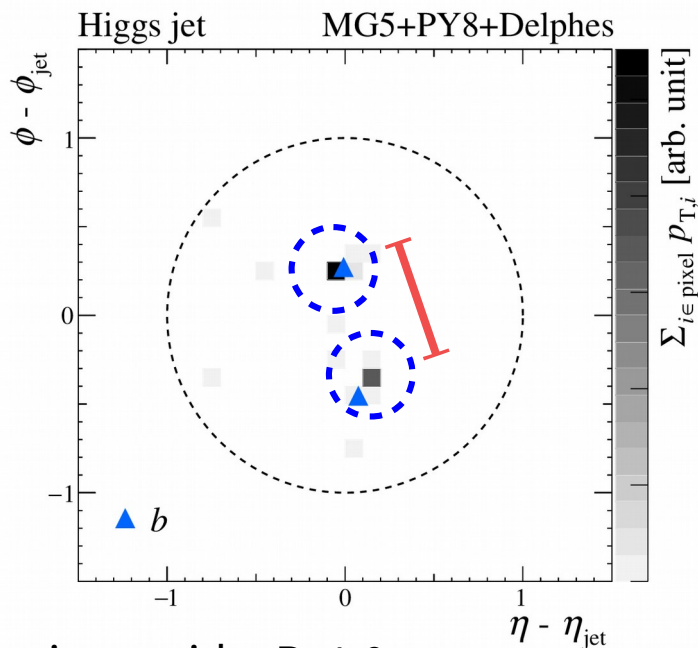
First, let us focus on correlation between hard constituents.

We may consider the two-point correlation spectrum of **trimmed jet**.

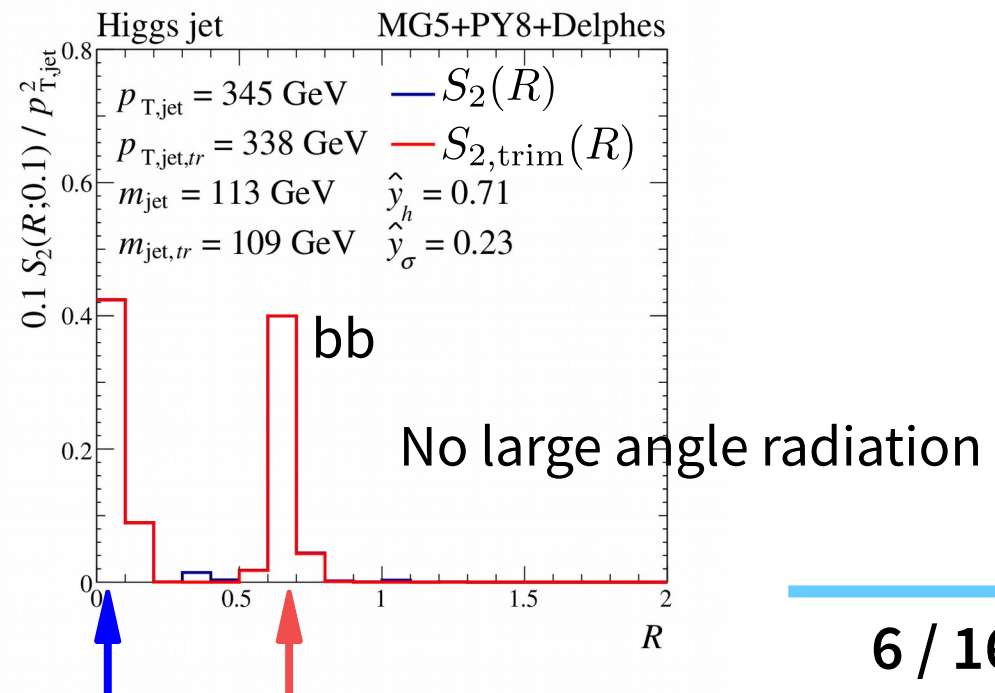
$$S_{2,\text{trim}}(R) = \int d\vec{R}_1 d\vec{R}_2 P_{T,\mathbf{J}_{\text{trim}}}(\vec{R}_1) P_{T,\mathbf{J}_{\text{trim}}}(\vec{R}_2) \delta(R - R_{12})$$

only sensitive to **hard-hard correlations**

For Higgs jet:



Calorimeter jet, anti-kt, R=1.0



Two-Point Correlation Spectrum: KEK

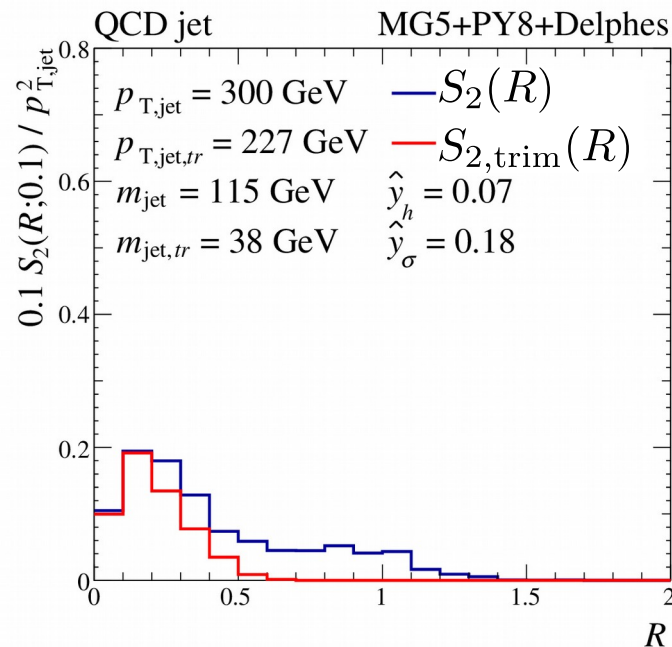
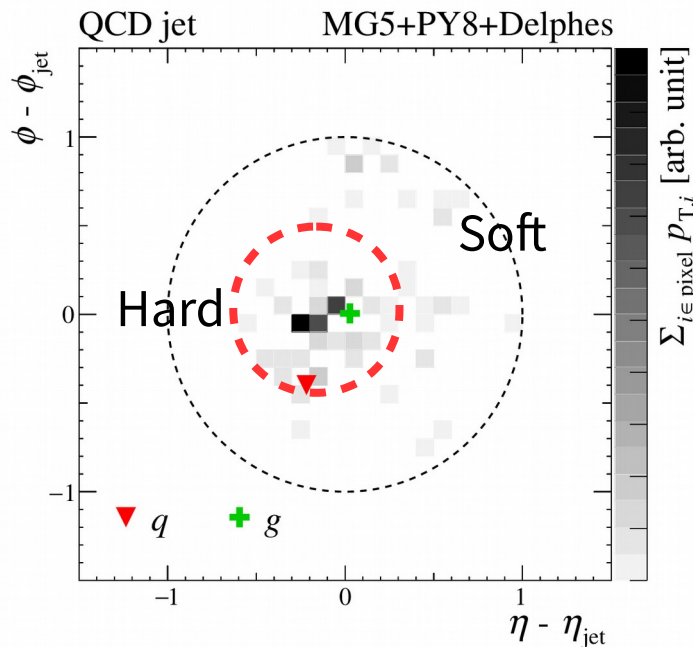
Hard-Soft Correlation

QCD jets have significant soft radiations. We may consider correlation between the soft parts and the hard parts.

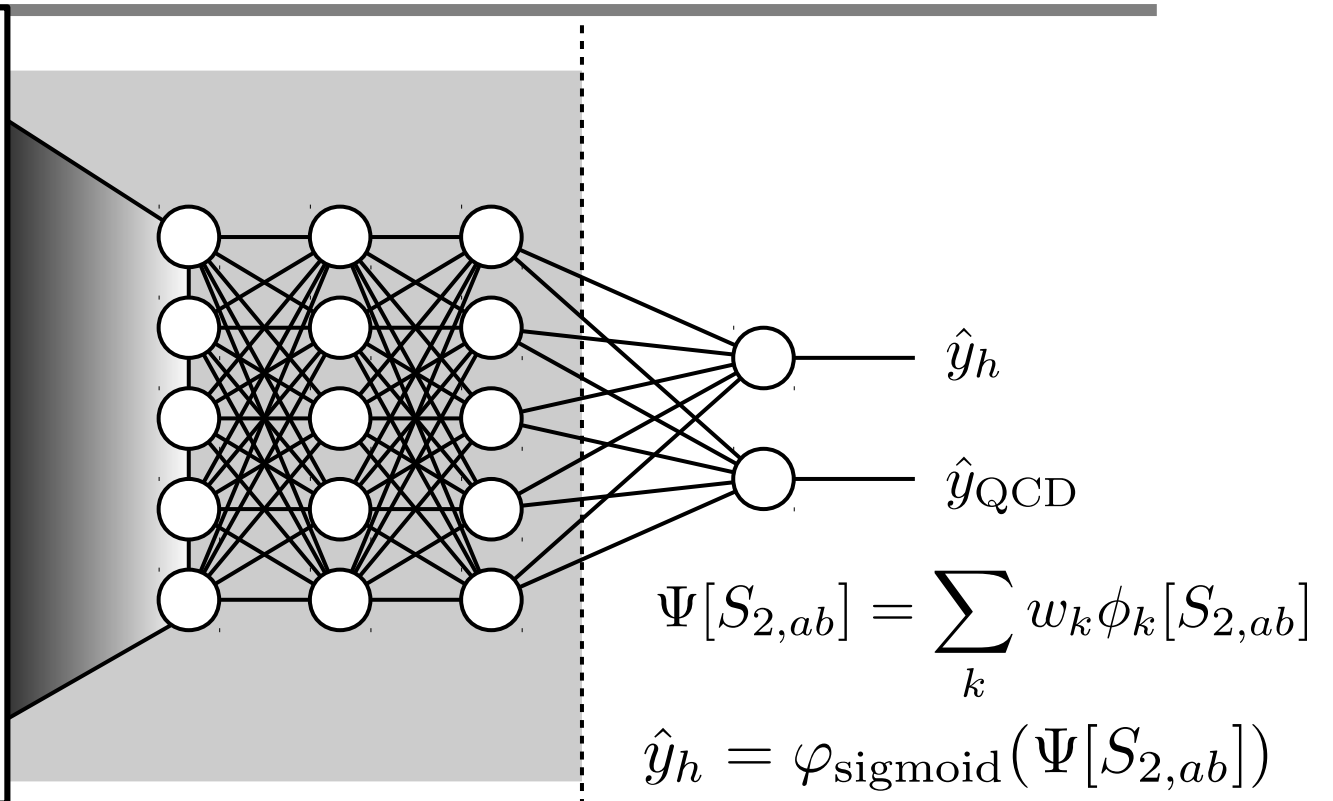
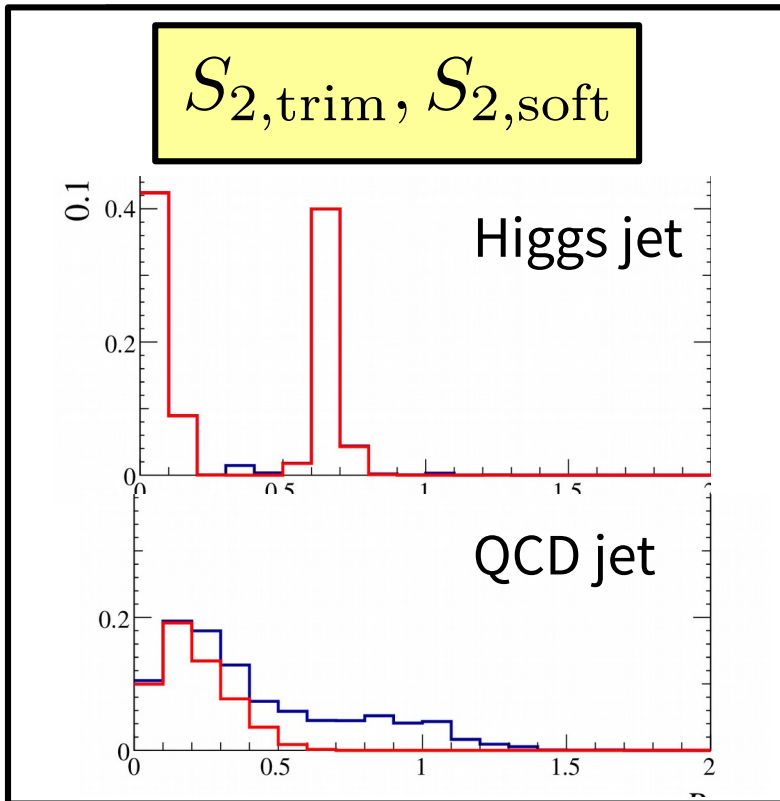
$$S_{2,\text{soft}}(R) = S_2(R) - S_{2,\text{trim}}(R)$$

sensitive to **hard-soft correlations**
subleading **soft-soft correlations**

For QCD jet:



Replacing CNN to MLP+S2



Inputs: two-point
corr. Spectrum
(+ jet mass, PT)

$$\phi_k = \Phi_k[S_{2,ab}]$$

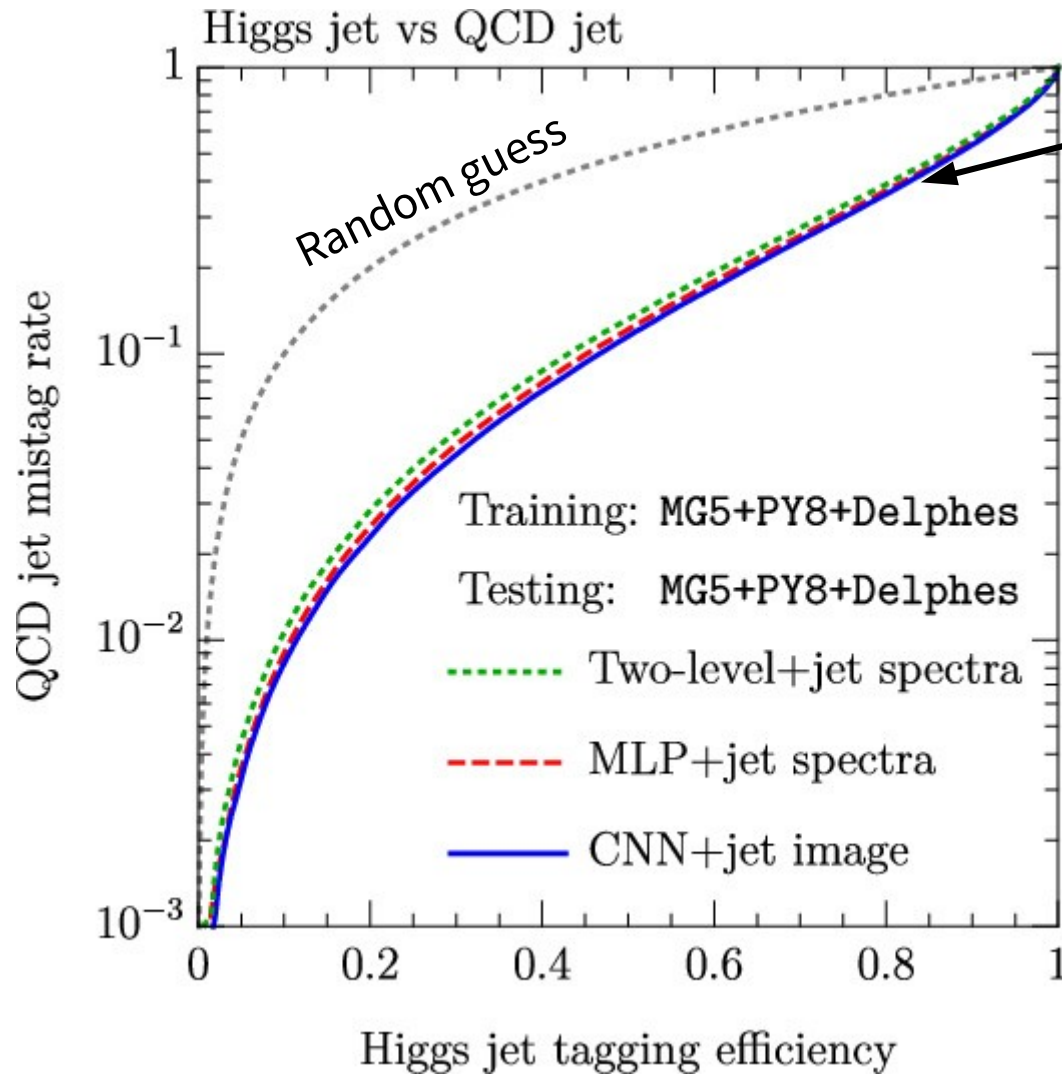
Feature map: MLP

Logistic regression

This architecture represents:

$$\sum_{n=0}^{\infty} \mathcal{O}[P_T^{2n}]$$

Equal Performance between CNN and MLP



Equal performance!!

$$\text{CNN} \sim \sum_{n=0}^{\infty} \mathcal{O} [P_T^n]$$

$$\text{MLP+S2} \sim \sum_{n=0}^{\infty} \mathcal{O} [P_T^{2n}]$$

For Higgs jet vs. QCD jet classification, **MLP+S2** is sufficient.

Interpretable Setup: $\sum_{n=0}^{\infty} \mathcal{O} [P_T^{2n}] \rightarrow \mathcal{O} [P_T^2] + \dots$

We may try the following two-level setup

Level1: substructure analyzer

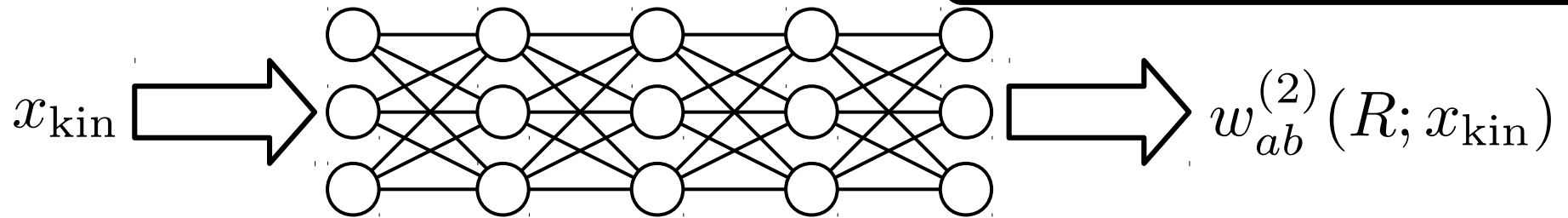
$$\Phi[P_T; x_{\text{kin}}] = \int dR S_{2,ab}(R) w_{ab}^{(2)}(R; x_{\text{kin}}) \quad : \text{IRC safe}$$

$$x_{\text{kin}} = \{p_{T,\text{jet}}, m_{\text{jet}}\}$$

$$\hat{R}_{b\bar{b}} = \frac{2m_h}{p_{T,h}}$$

Use neural network to approximate the weight function.

Level2: kinematics analyzer



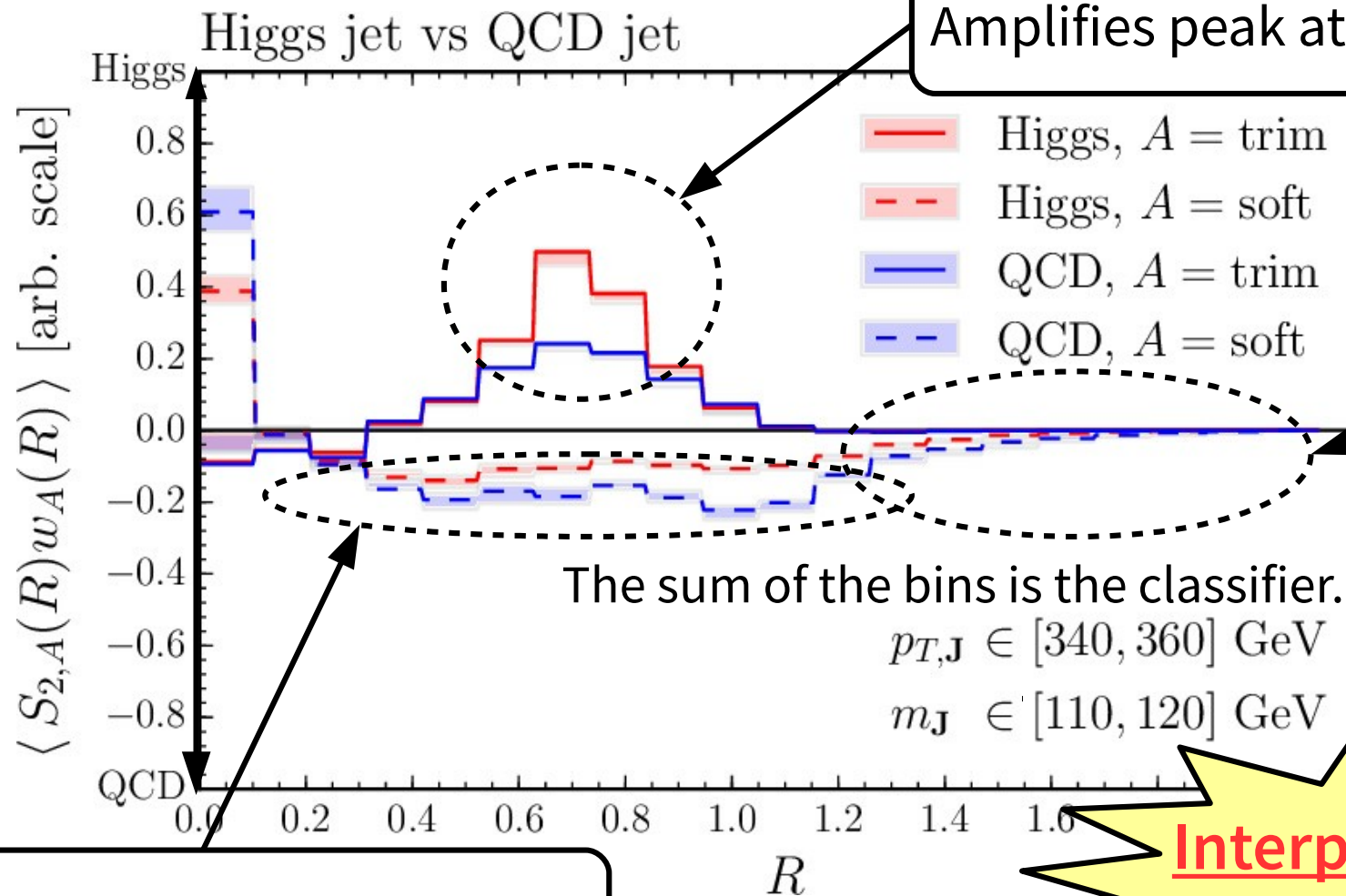
This architecture gives you two interpretable quantities:

$w_{ab}^{(2)}(R; x_{\text{kin}})$ shows the **functional form** of the energy correlator.

$S_{2,ab}(R)w_{ab}^{(2)}(R; x_{\text{kin}})$ shows the **contribution** to the classifier.

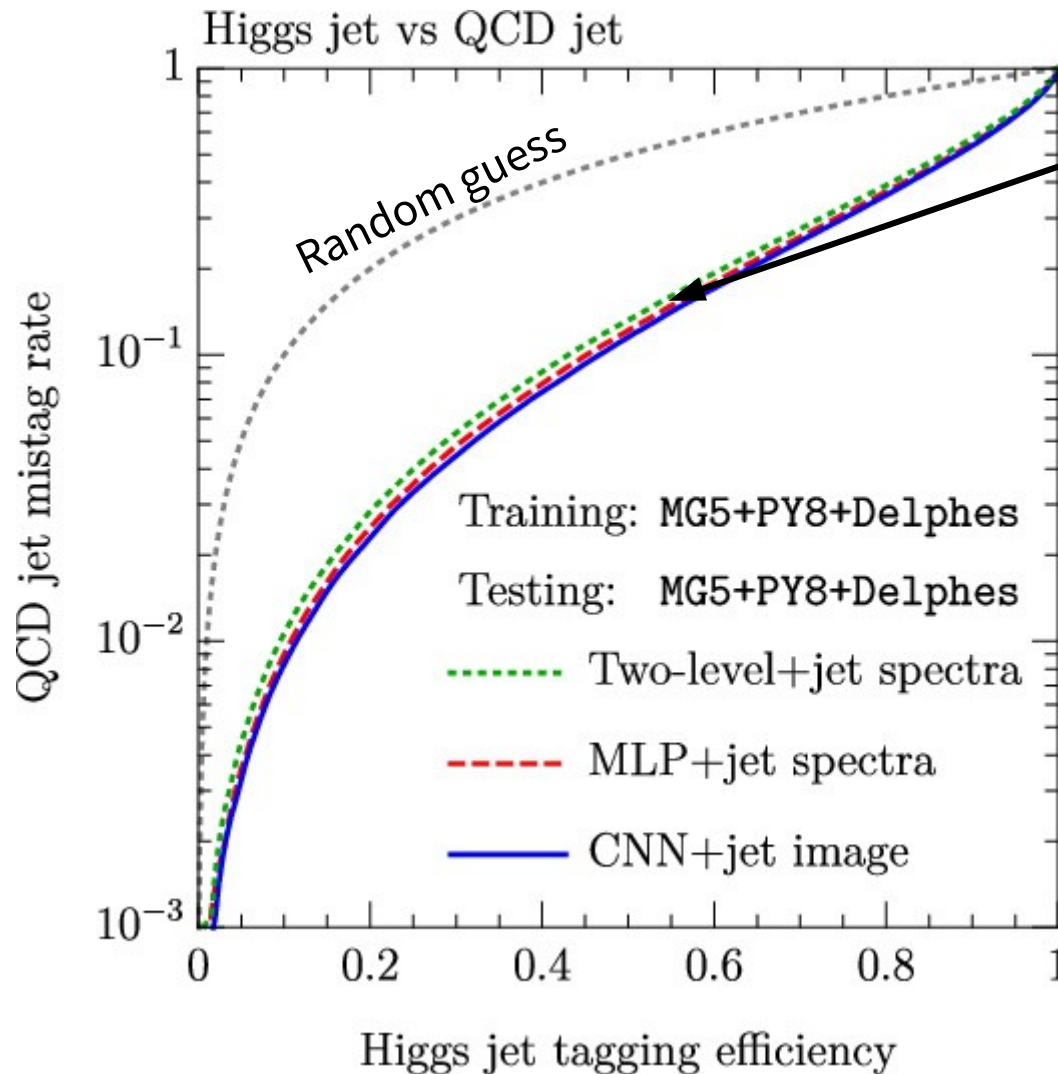
Average of the linear classifier outputs

$$\Phi[S_{2,ab}] = \int dR S_{2,\text{trim}}(R)w_{\text{trim}}^{(2)}(R) + \int dR S_{2,\text{soft}}(R)w_{\text{soft}}^{(2)}(R)$$



Interpretable

Similar Performance between CNN and the interpretable network



Still similar performance!

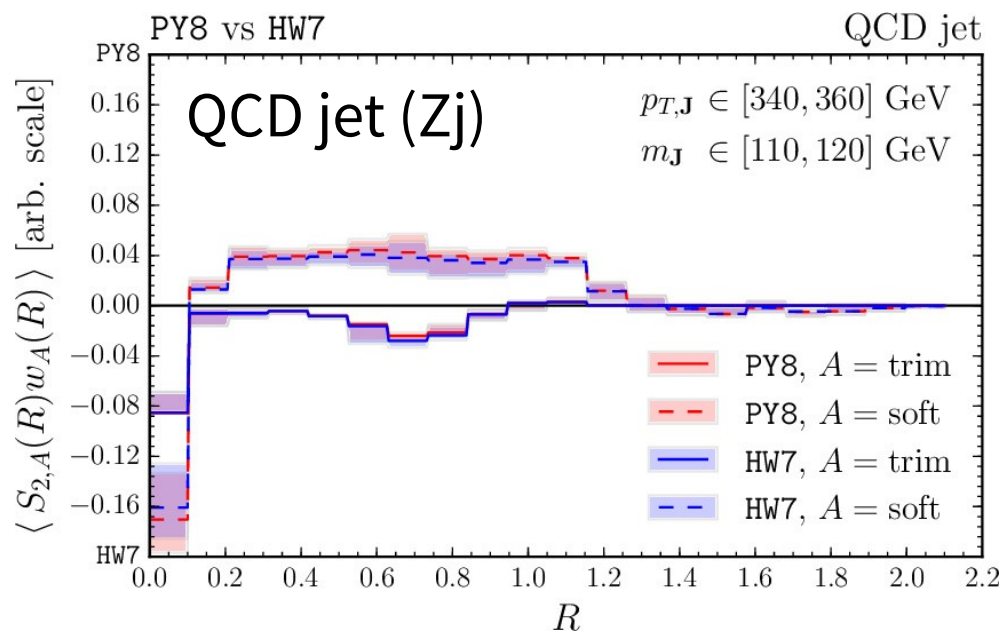
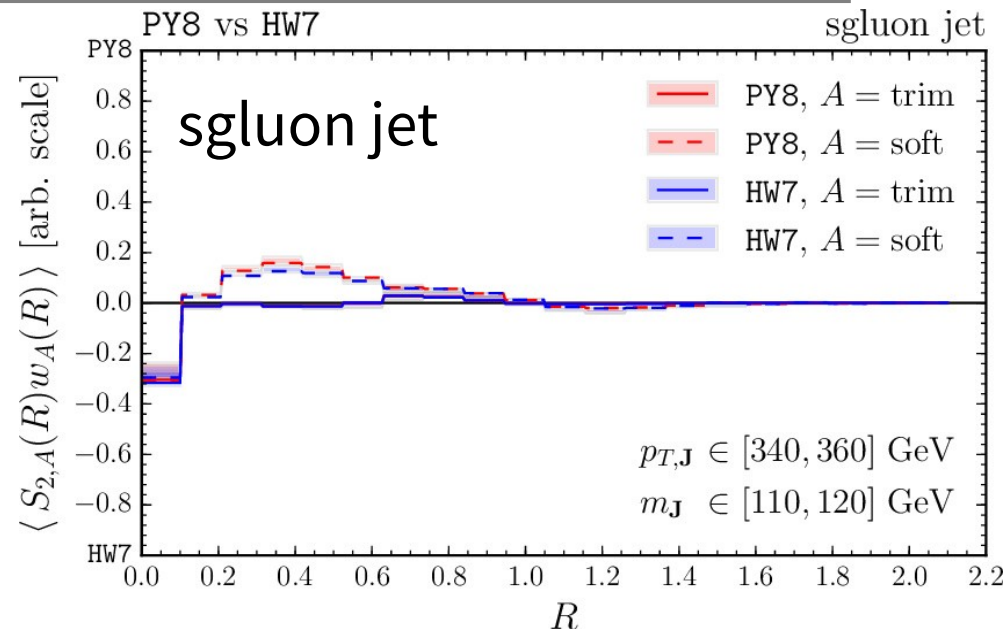
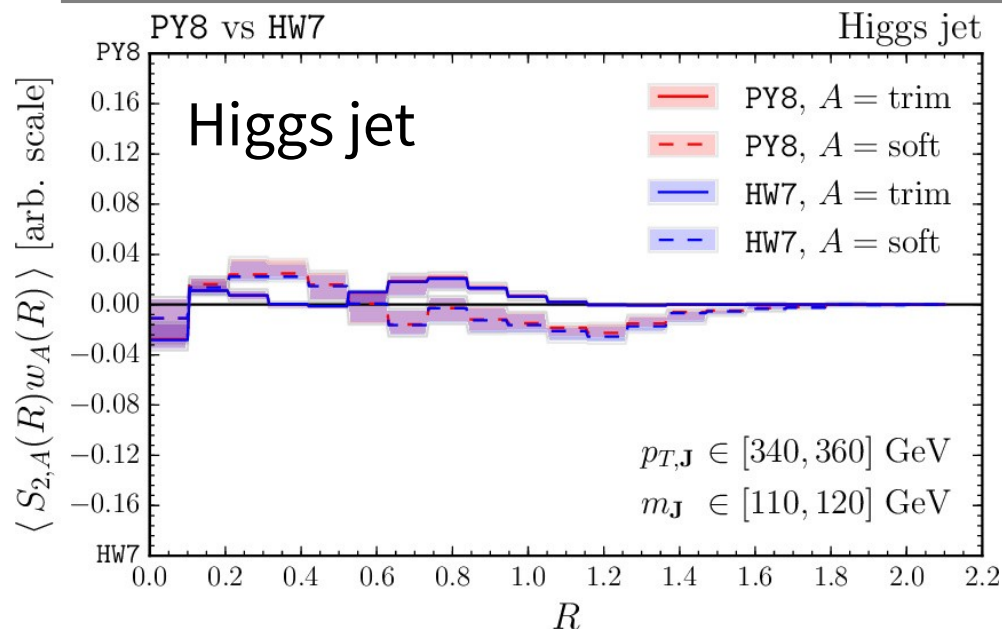
$$\text{MLP+S2} \sim \sum_{n=0}^{\infty} \mathcal{O} [P_T^{2n}]$$

Two-level+S2 \sim

$$\mathcal{O} [P_T^2 \cdot w(p_{T,J}, m_J)]$$

For Higgs jet vs. QCD jet classification, **Two-level+S2** is sufficient.

Comparing Pythia and Herwig



Most of the classifier contribution comes from soft activity

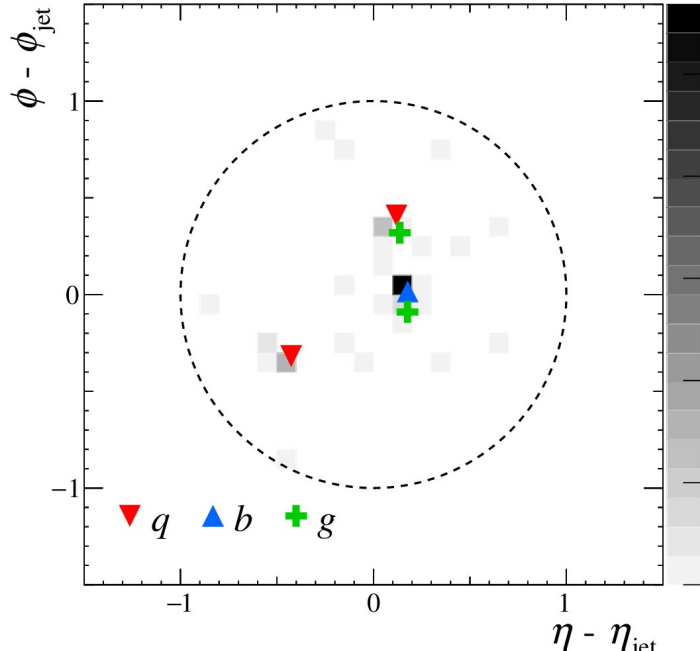
PY8 jets are more compact

Singlet and octet have different behavior

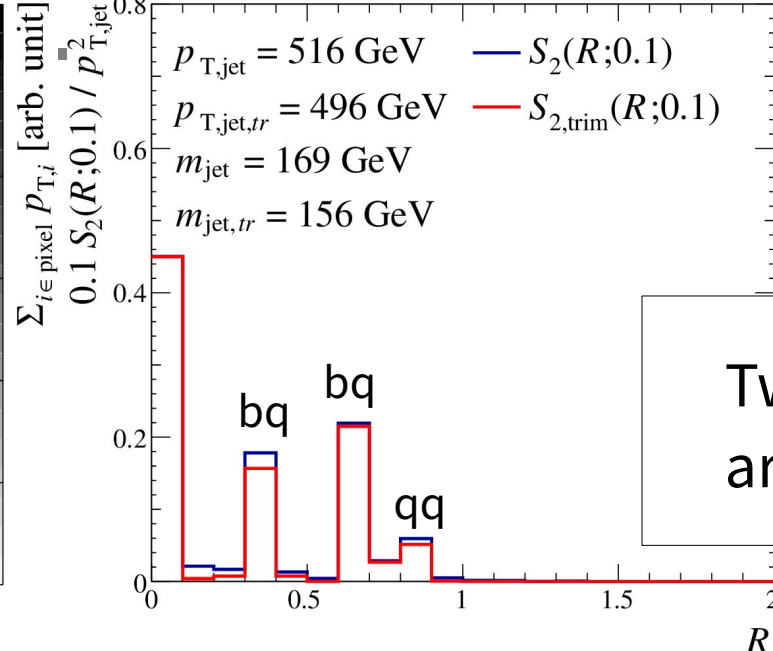
Top Jets (work in progress...)

Preliminary

Top jet MG5+PY8+Delphes

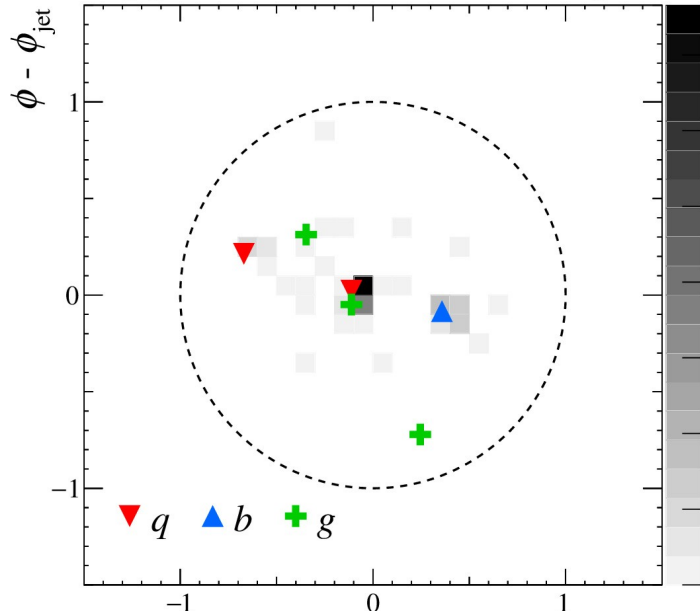


Top jet MG5+PY8+Delphes

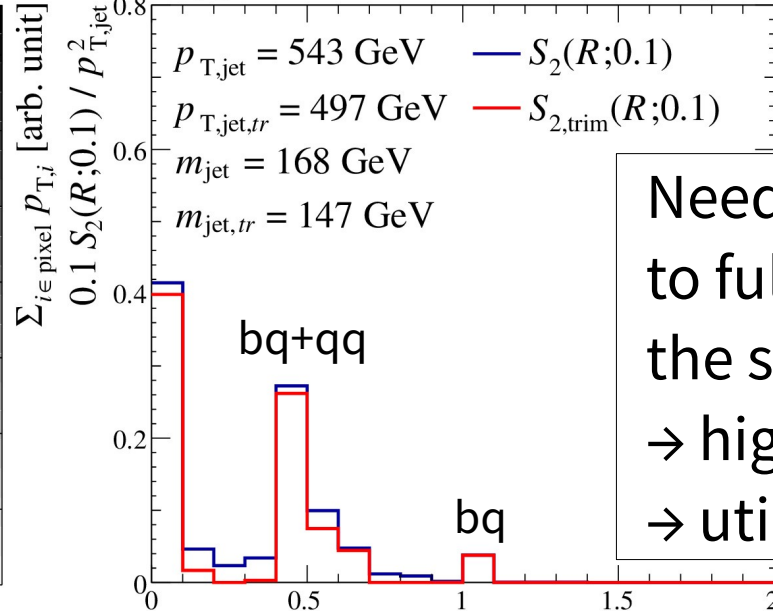


Two-point correlations are enough.

Top jet MG5+PY8+Delphes



Top jet MG5+PY8+Delphes



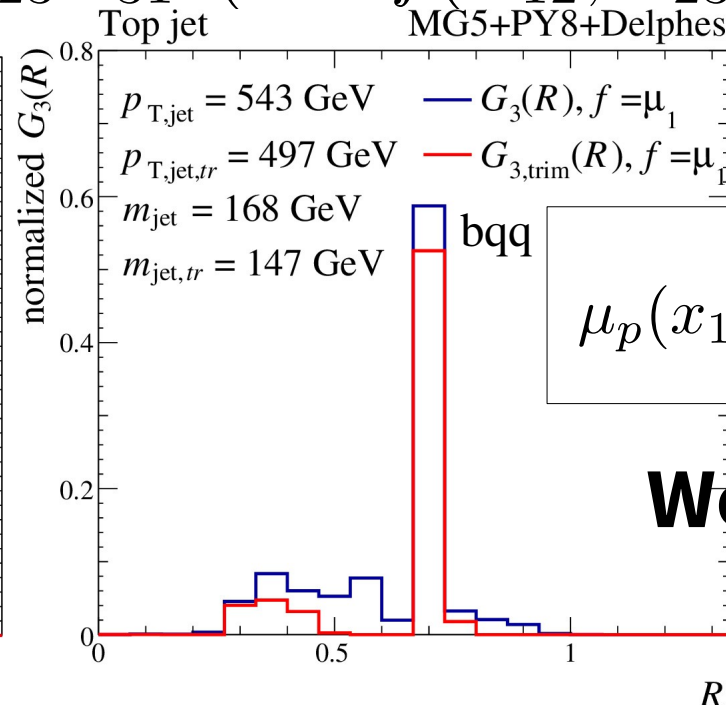
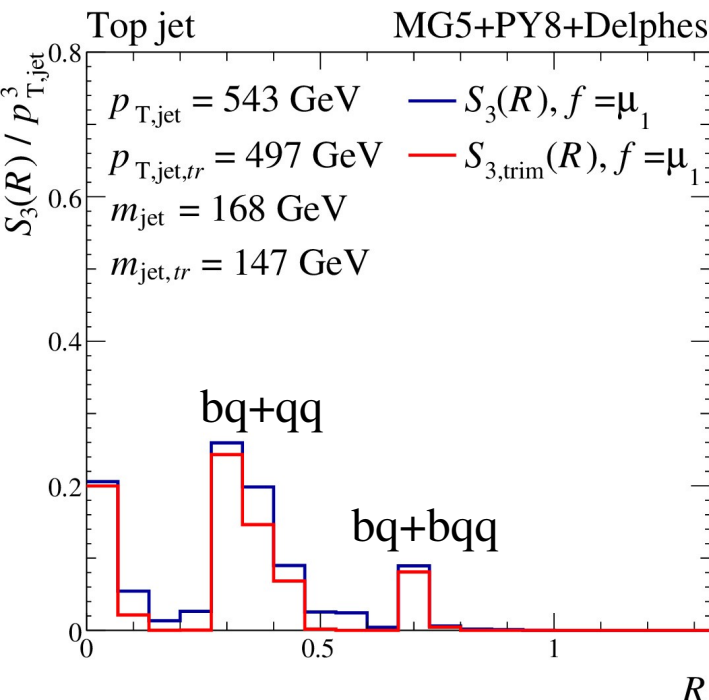
Need more information to fully encode the substructures.
 → higher order correlation?
 → utilize subjet information?

Three-point correlation spectrum

Preliminary

$$S_{3,abc}(R) = \int d\vec{R}_1 d\vec{R}_2 d\vec{R}_3 P_{T,a}(\vec{R}_1) P_{T,b}(\vec{R}_2) P_{T,c}(\vec{R}_3) \times \delta(R - f(R_{12}, R_{23}, R_{31}))$$

$$G_{3,abc}(R) = \int d\vec{R}_1 d\vec{R}_2 d\vec{R}_3 P_{T,a}(\vec{R}_1) P_{T,b}(\vec{R}_2) P_{T,c}(\vec{R}_3) \times R_{12} R_{23} R_{31} \delta(R - f(R_{12}, R_{23}, R_{31}))$$



$$\mu_p(x_1, \dots, x_N) = \left(\frac{\sum_i x_i^p}{N} \right)^{\frac{1}{p}}$$

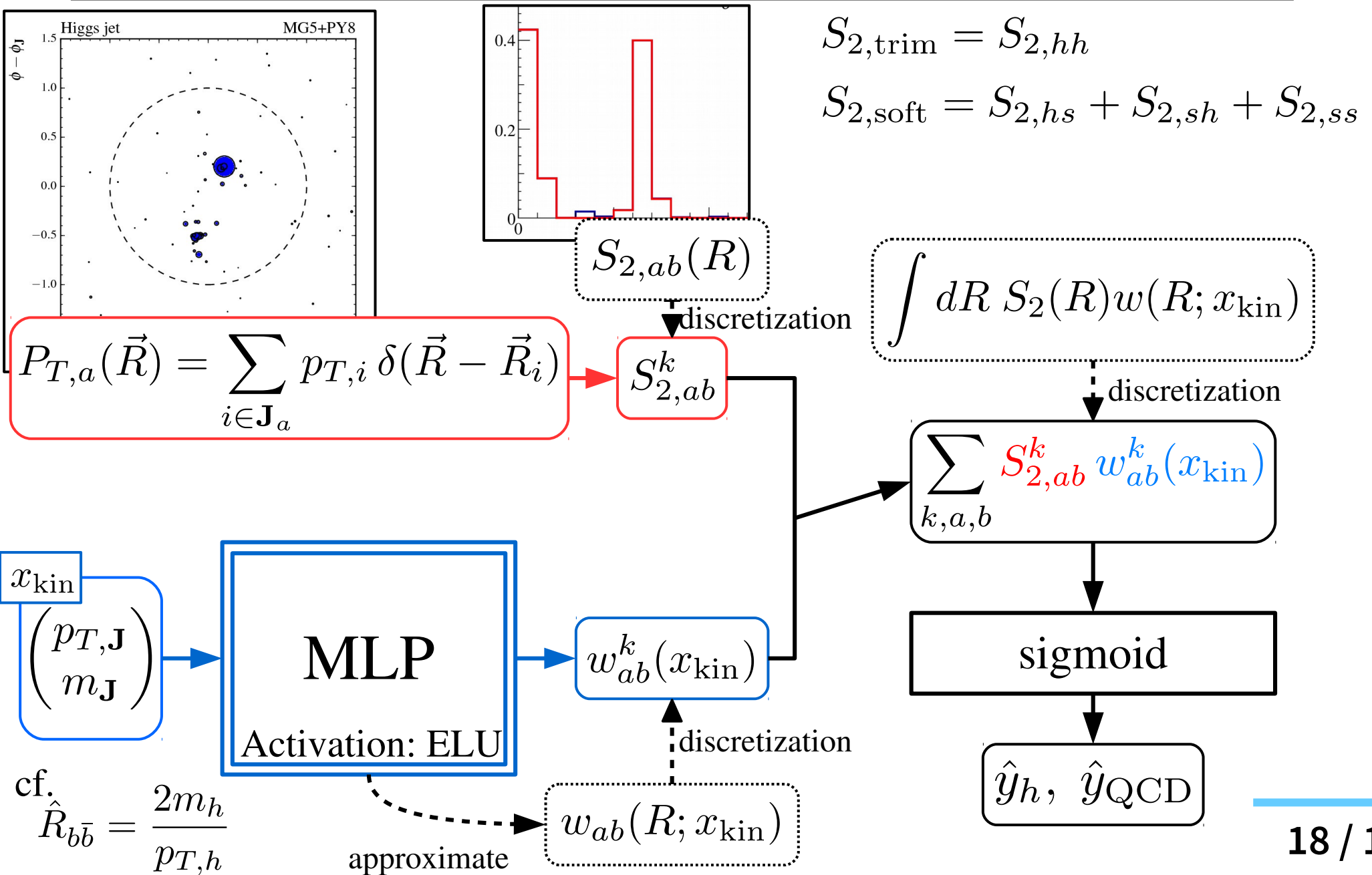
Work in progress...

Summary

- For the next run of LHC and future colliders, we need a quick and reliable jet substructure analysis framework.
- We developed a machine learning framework using **two-point correlation spectrum** for analyzing jet substructures.
- The spectrum is derived from the jet image analysis and the corresponding two-level model is **interpretable**.
- This analysis strategy is not limited to Higgs jet vs. QCD jet classification, but we may use it for comparing Monte-Carlo simulations (and real data).
- Analyses with more complex objects (top jets...) are ongoing.

Please stay tuned!

BACKUP



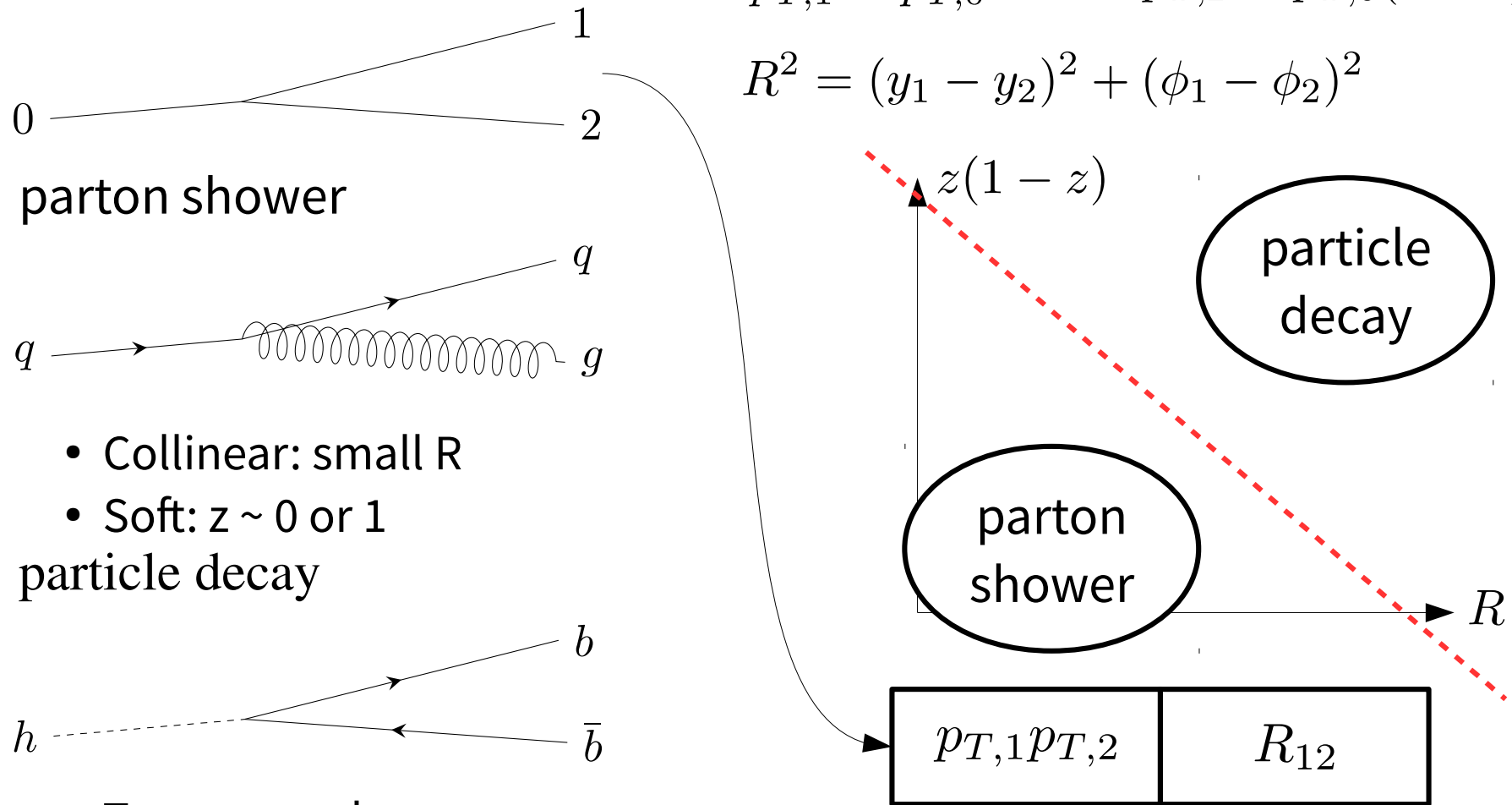
Why the linear classifier works?

The parameter set $(p_{T,0}, z, R)$ determines the kinematics of parton evolution.

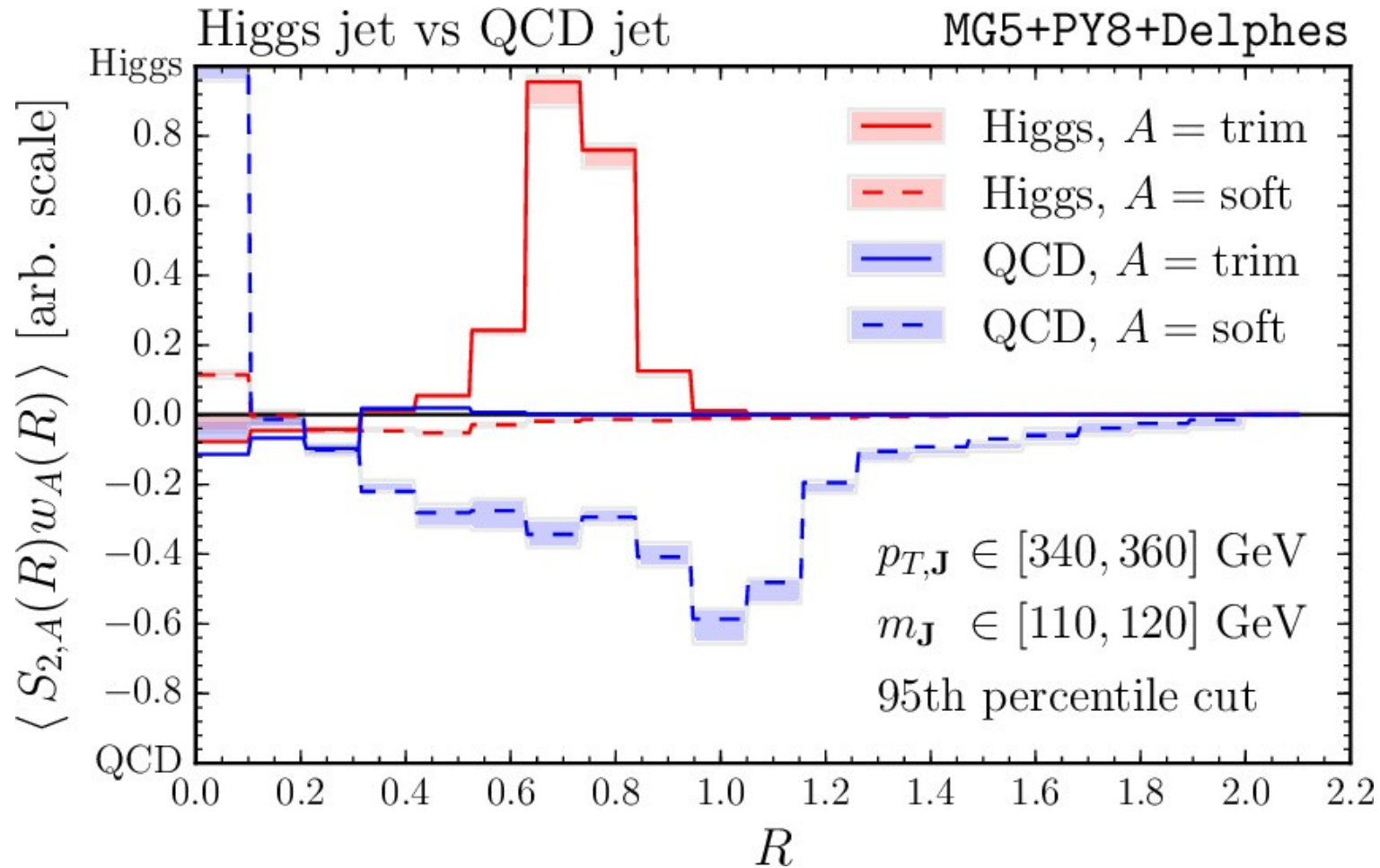
$$p_{T,1} = p_{T,0}z \quad p_{T,2} = p_{T,0}(1 - z)$$

$$R^2 = (y_1 - y_2)^2 + (\phi_1 - \phi_2)^2$$

- parton shower
 - Collinear: small R
 - Soft: $z \sim 0$ or 1
- particle decay
 - Transverse decay: $R \sim 2m/pt$ and $z \sim 0.5$

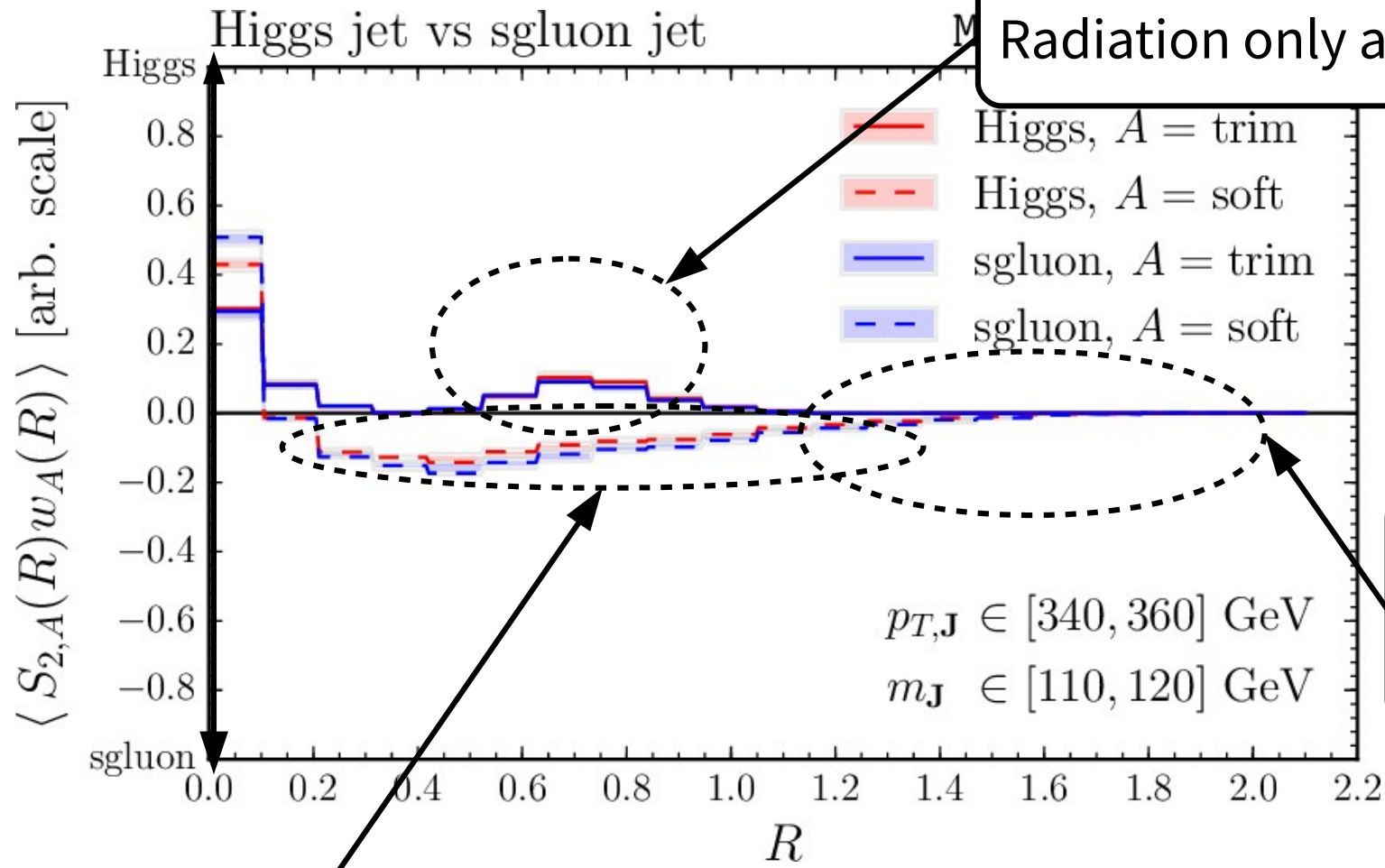


Average of the linear classifier outputs



Two-Prong jet: 1 vs 8

Identification of color of originating parton: **1 vs 8**

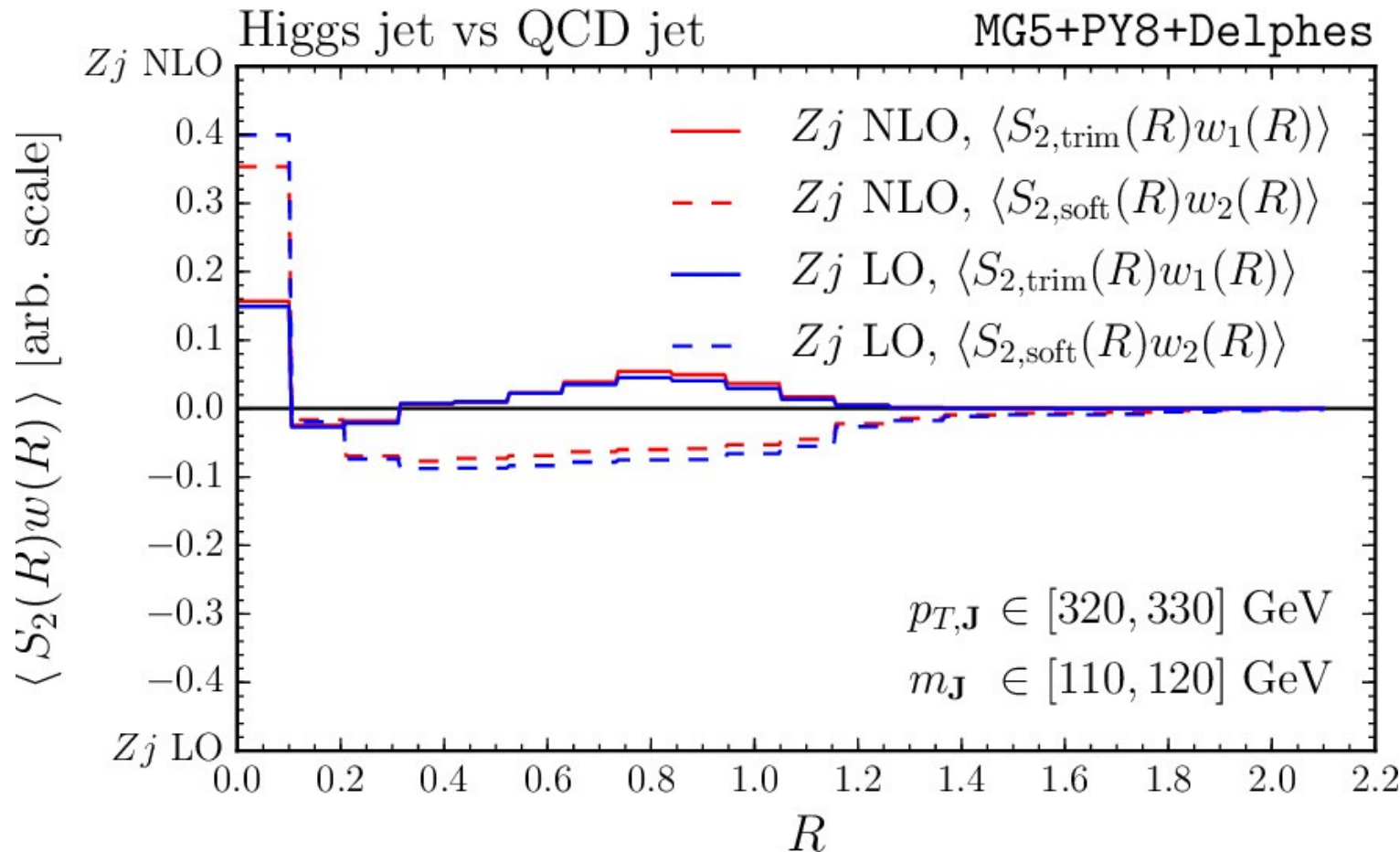


Radiation only at Rbb: **Higgs jet**

Large angle radiation:
sgluon jet

More soft activity: **sgluon jet**

Training with MC@NLO events



We can check the sanity of training with negative weights.