

ジェットの物理と 解釈可能な機械学習

Sung Hak Lim, Mihoko Nojiri (JHEP, 2018)

Amit Chakraborty, Sung Hak Lim, Mihoko Nojiri (arXiv 1904.02092)

Machine learning in Jet Physics

今後のコライダー物理

Run III → HL-LHC → FCC?

*シグナルやバックグラウンド

データ量が多い

の理解がより重要

** 効率も大事

hadronization,
PDF, parton shower
modeling, ...

high pT objects (Events in Tail)
soft object, mono something
or something unknown

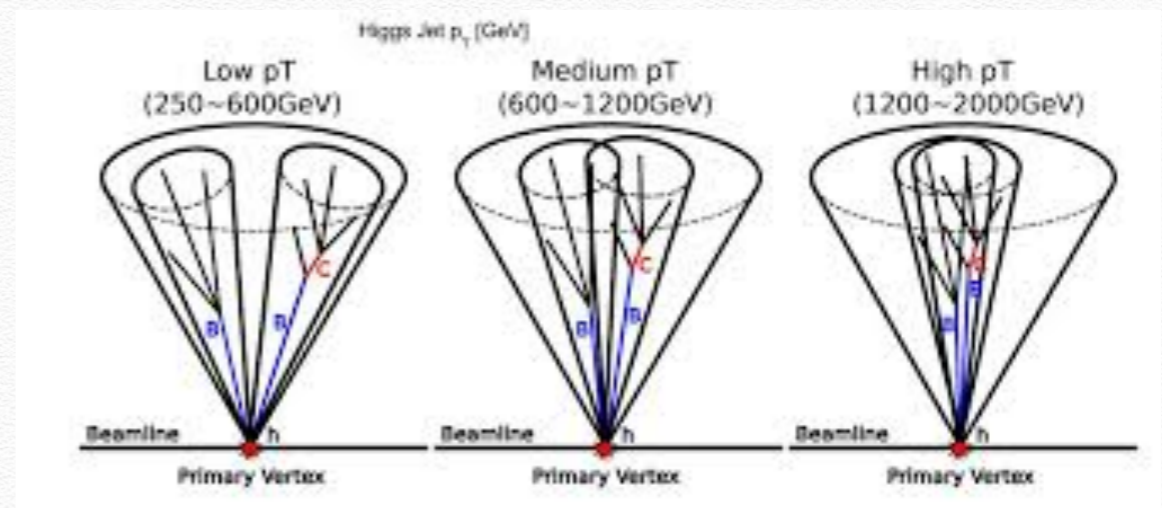


Physics outputs
effective operator,
top partner
dark matter..

ジェットの物理と機械学習

- ❖ Jet physics :QCDがすごく成功した領域 (Theoretical understanding + computation)
 - ❖ QCD に優しいジェット再構成ダイアグラム (kT, CA, antikT)
 - ❖ ブーストした信号への興味 QCD jet (huge background) vs boosted Higgs
→ Jet substructure (mass drop), Top 再構成 (BSM search)
 - ❖ Jet の性質を規定する量(D2とか) and minimal Validation Analysis (leaving optimization to algorithm BDT とか)

理論と計算機の発達が重要

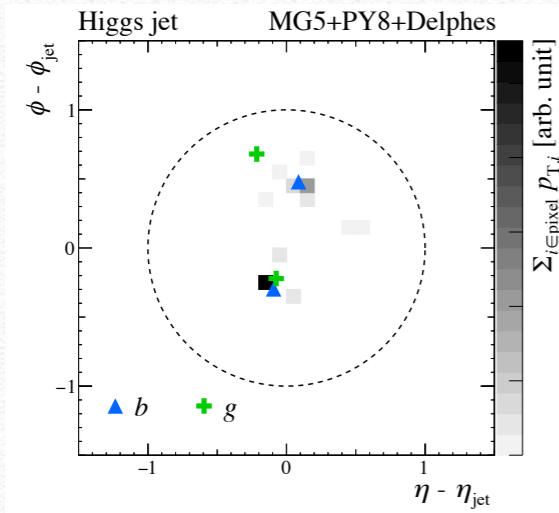


New wave — 機械学習

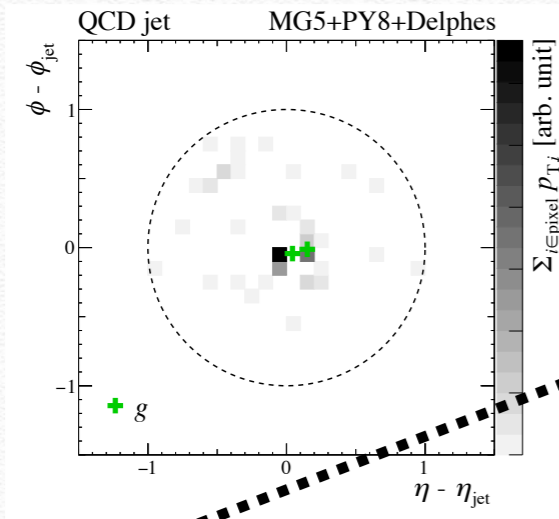
ML:basic units

Input: Jet images

Higgs

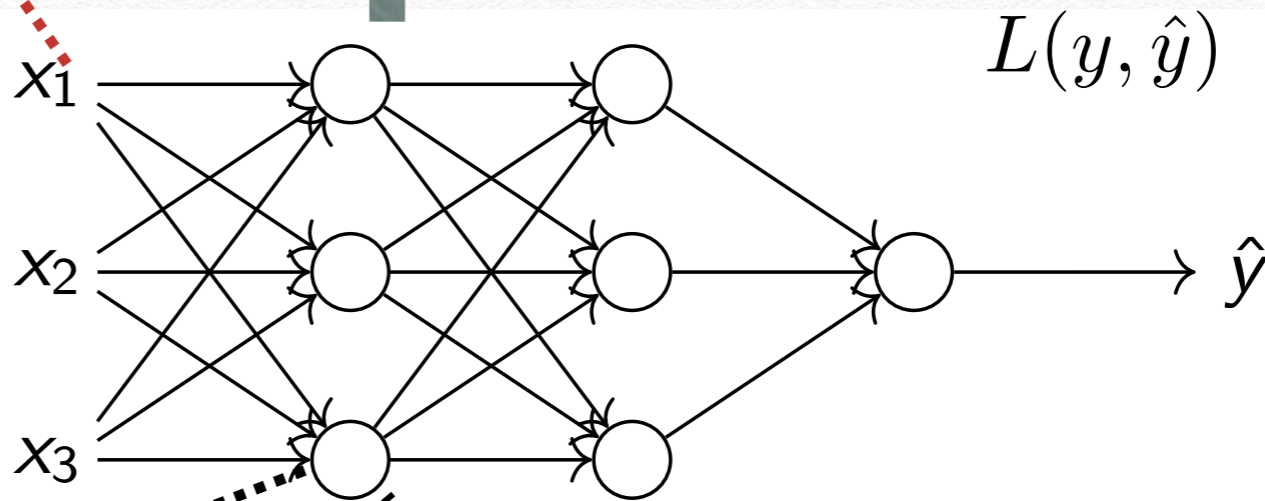


QCD



output: w_{ij}, b_i

損失関数(最小化)



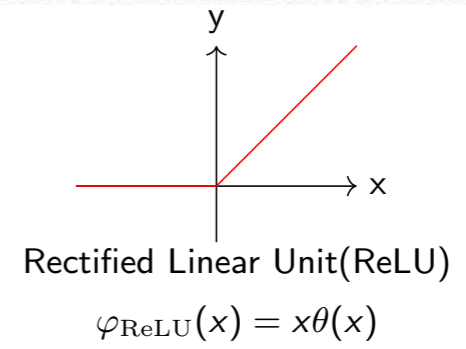
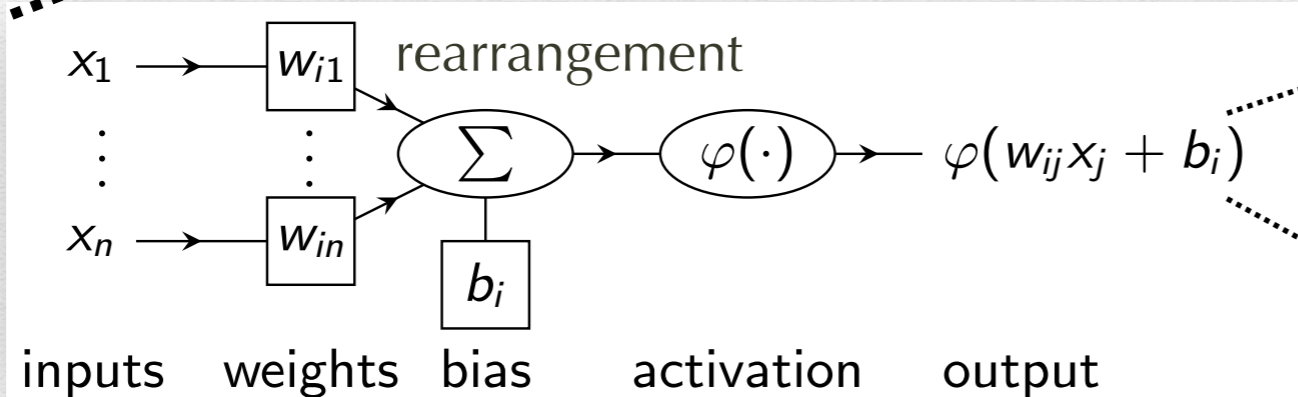
Higgs
 $y = (1, 0)$

QCD
 $y = (0, 1)$

hidden layer

ノード

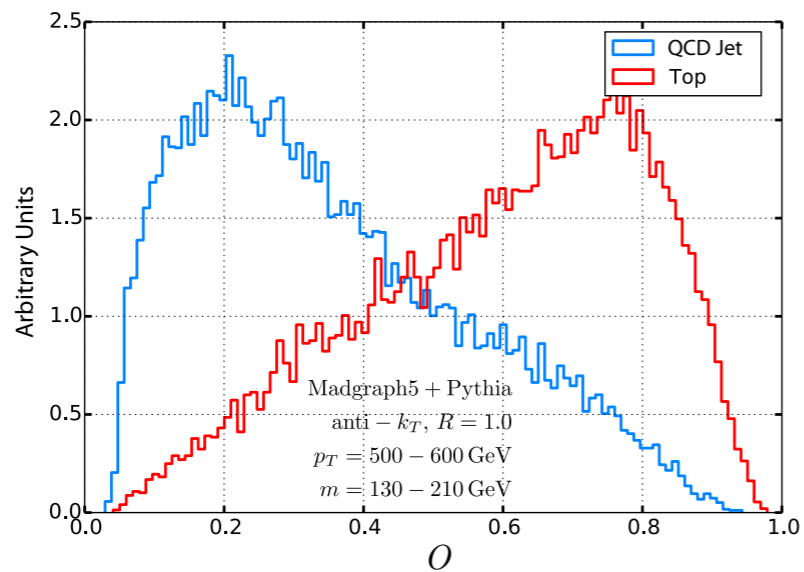
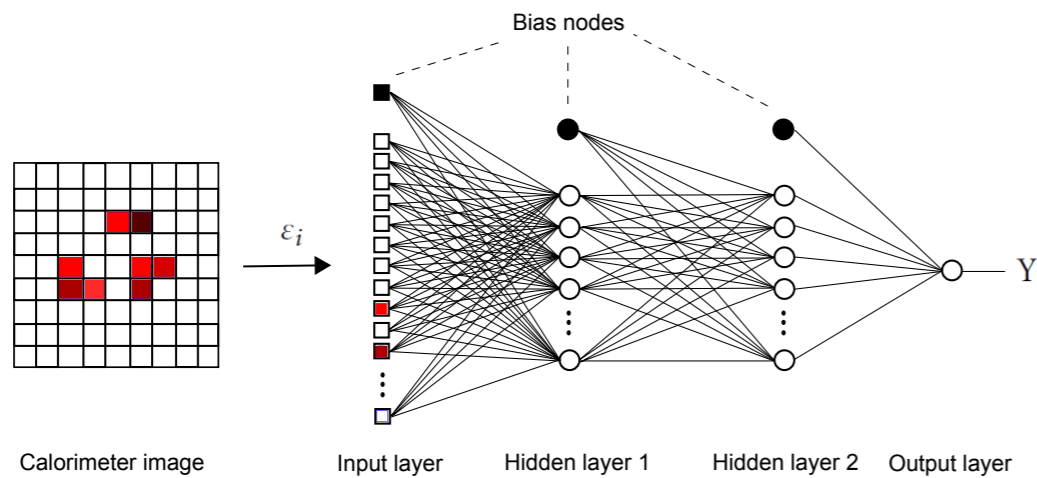
φ : 非線形性



入力の構造化

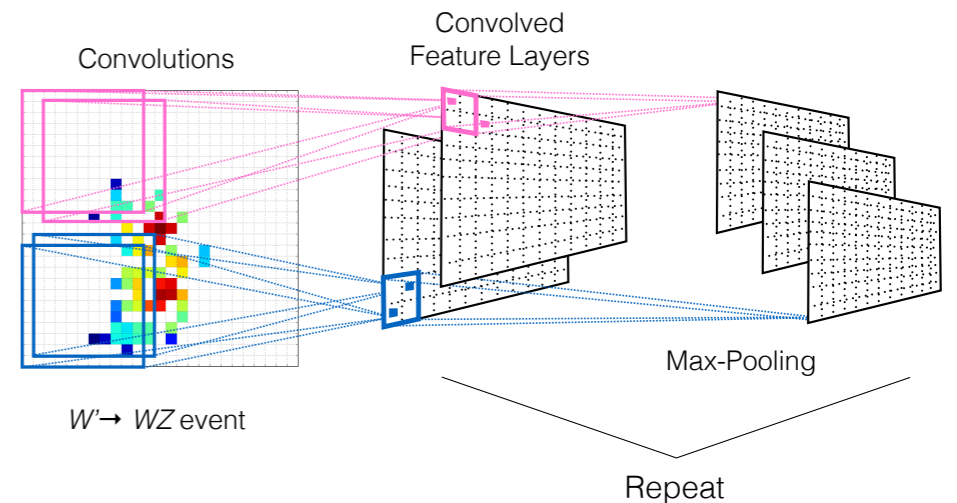
Almeida et al 1501.05968

DNN (all bins in a line)

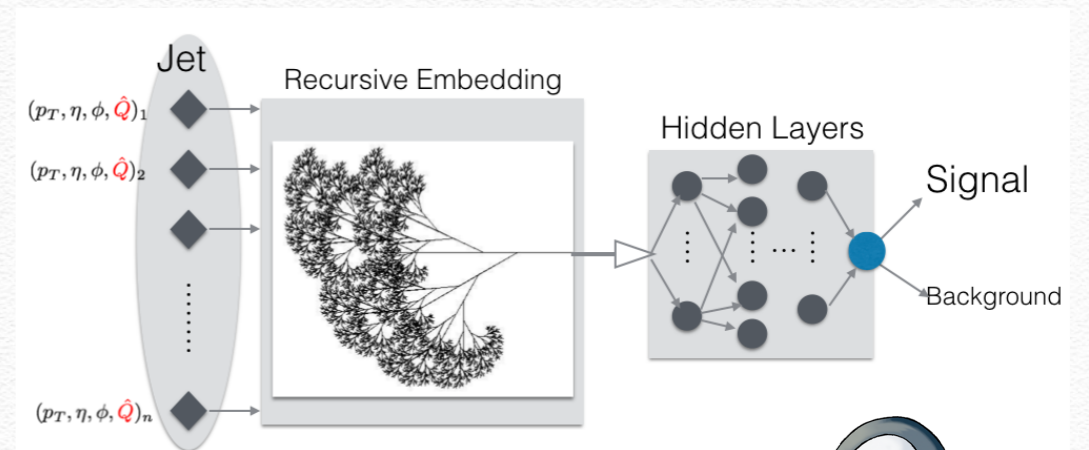


CNN

correct nearby image first 1511.05190



Recursive(Taoli Cheng 1711.02633)



“ Maybe we do not know physics behind, architecture do the job for you”??



機械学習は何をやっているか

- ❖ QCD ジェットとHiggs ジェットの分類は得意らしい
- ❖ 何をやってるかわからない（パラメータが多い。冗長性が多いように見える
- ❖ ジェットイメージを使った解析がQCD のこれまでの知識の範囲にあるのか、それともなんかすごいことをやっているのか。

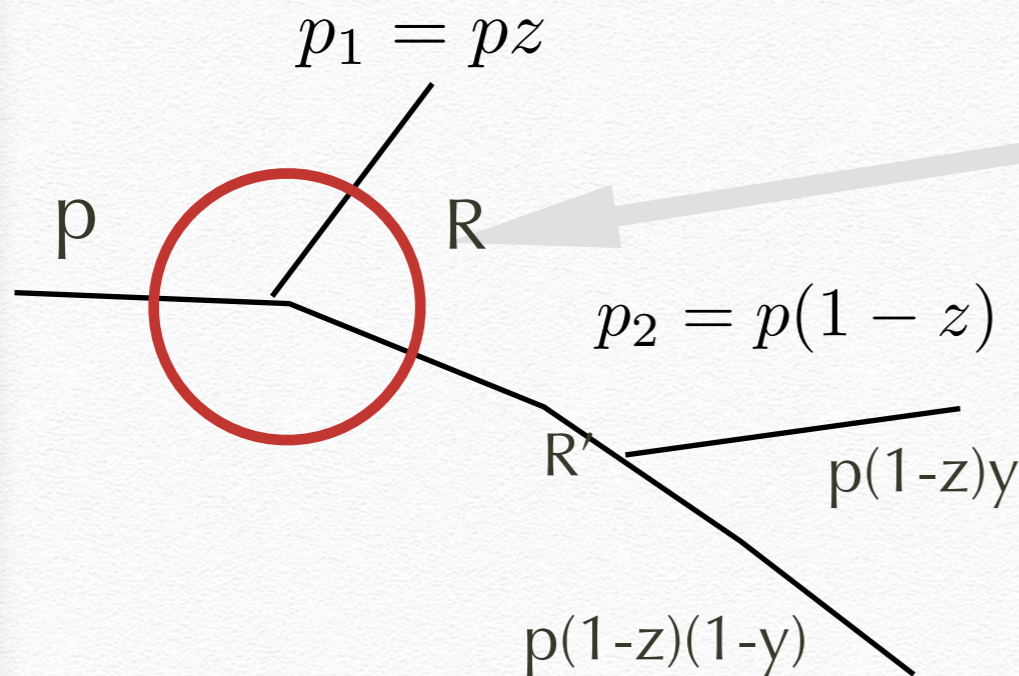
これから話すことのまとめ

jet image(400 程度のデータ) +CNNはやめてもうちょっと動機付けのある量（ジェットのスペクトラム-40くらいの入力）をDNNに入れて見た

1. Higgs jet, color octet のジェット, QCD のjet の分類問題をやらせてみるとCNN と変わらない。**多分CNN-> jet spectrum 構築-> 分類をやっている。**
2. “解釈可能な機械学習モデル” を作ってML が何をやってるか明らかに

最小のインプット (Jet spectrum) とは何か

❖ Monte Carlo → Parton splitting + hadronization



(p, R, z) describe
parton shower splitting

$$p_1 p_2 = p^2 z(1-z)$$

$$p_1^2 + p_2^2 = p^2 [(1-z)^2 + z^2]$$

ジェットはパートンスプリッティングの繰り返しなので、特定のRの範囲に入っている

終状態の全ての粒子のペアについて、以下のものを計算する = ジェットスペクトル

$$S_2(R, \Delta R) = \sum_{ij} p_{Ti} p_{Tj} \text{ for } R < R_{ij} < R + \Delta R \quad \text{for } R = 0, 0.1, 0.2 \dots$$

IRC safeなコンビネーション (C-correlator)

ジェットイメージとの関係について

$$\text{jet} = \text{energy flow (+ ...)} \quad P_T(\vec{R}) = \sum_{i \in \mathbf{J}} p_{T,i} \delta(\vec{R} - \vec{R}_i),$$

*classifier using energy flow

$$h_i = \hat{\Psi}_i[P_T]$$

calorimeter hit position

$$h_i = w_i^{(0)} + \int d\vec{R} P_{T,a}(\vec{R}) w_{i,a}^{(1)}(\vec{R}) + \frac{1}{2!} \int d\vec{R}_1 d\vec{R}_2 P_{T,a}(\vec{R}_1) P_{T,b}(\vec{R}_2) w_{i,ab}^{(2)}(\vec{R}_1, \vec{R}_2) + \dots$$

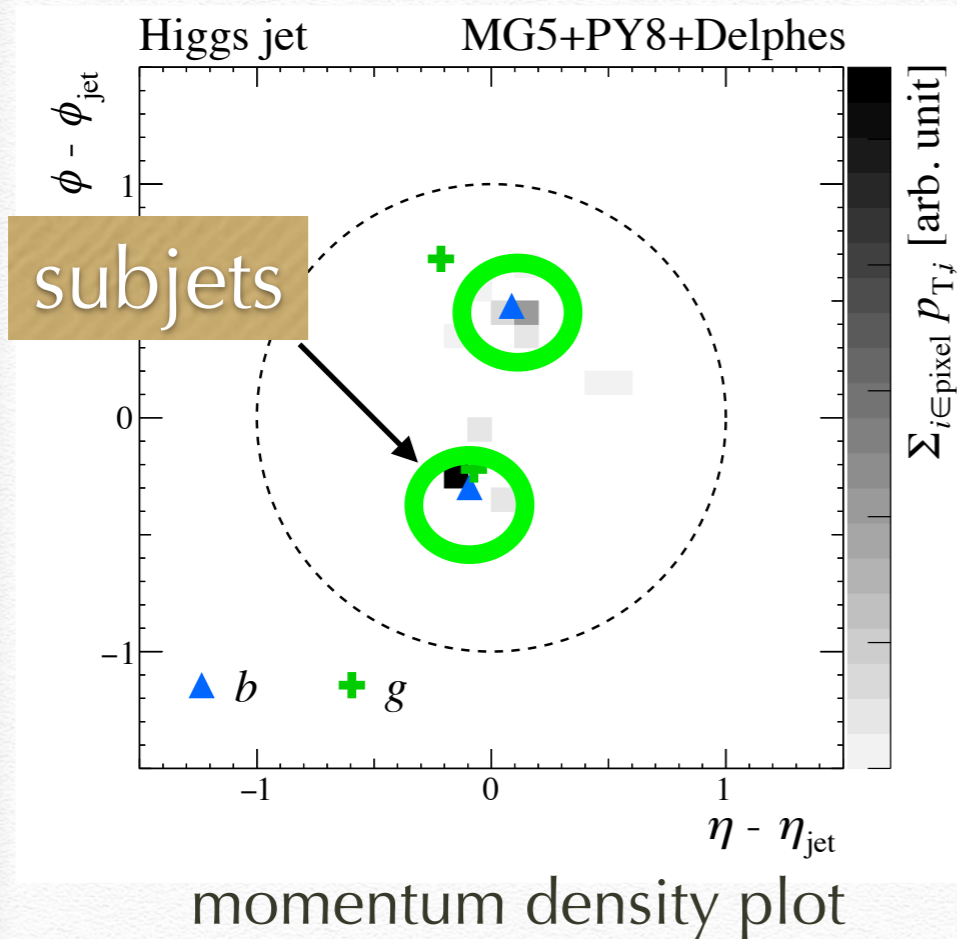
If w only depends on R_{12}

$$\frac{1}{2!} \int dR S_{2,ab}(R) w_{i,ab}^{(2)}(R) + \dots$$

*classifier using Jet spectrum

$$h_i = w_i^{(0)}(\vec{x}_{\text{kin}}) + \int dR S_{2,A}(R) \frac{w_{i,A}^{(2)}(R; \vec{x}_{\text{kin}})}{2} + \frac{1}{2} \int dR_1 dR_2 S_{2,A}(R_1) S_{2,B}(R_2) \frac{w_{i,AB}^{(4)}(R_1, R_2; \vec{x}_{\text{kin}})}{12} + \dots$$

ジェットの中のソフトとハードなアクティビティ



jet内部にハードなsubjects(IRC safe)がある
subjects に対する pt cut

=> trimmed jet

$$\mathbf{J}_{\text{trim}} = \bigcup_{\substack{a \\ \frac{p_{T,J_a}}{p_{T,J}} \geq f_{\text{trim}}}} \mathbf{J}_a .$$

$S_{2\text{trim}}$

$$S_{2,\text{trim}}(R; \Delta R) = \frac{1}{\Delta R} \sum_{i,j \in \mathbf{J}_{\text{trim}}} p_{T,i} p_{T,j} \cdot I_{[R,R+\Delta R)}(R_{ij}),$$

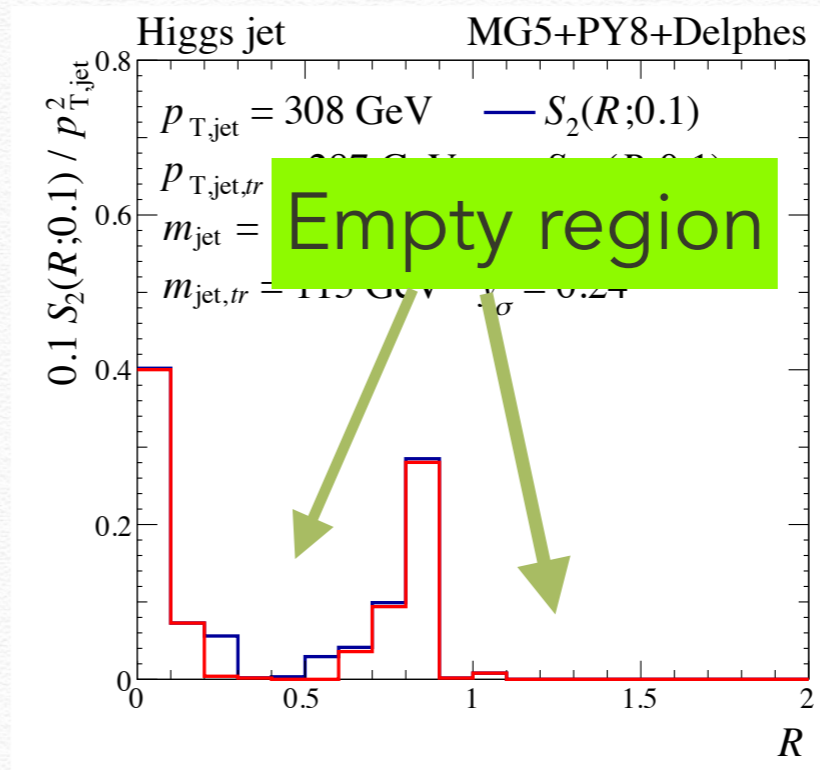
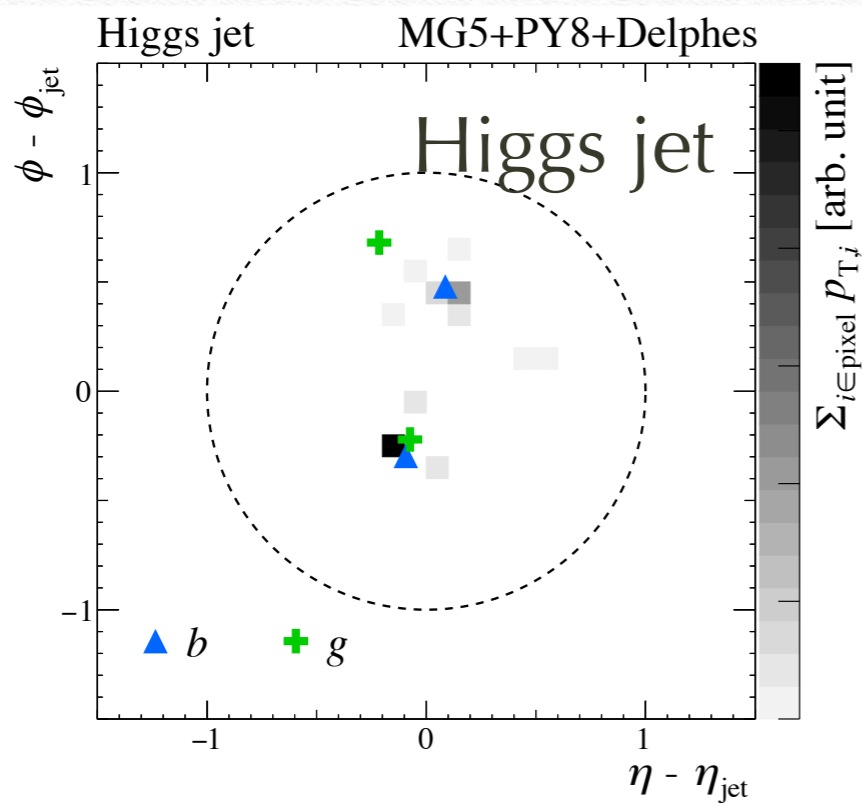
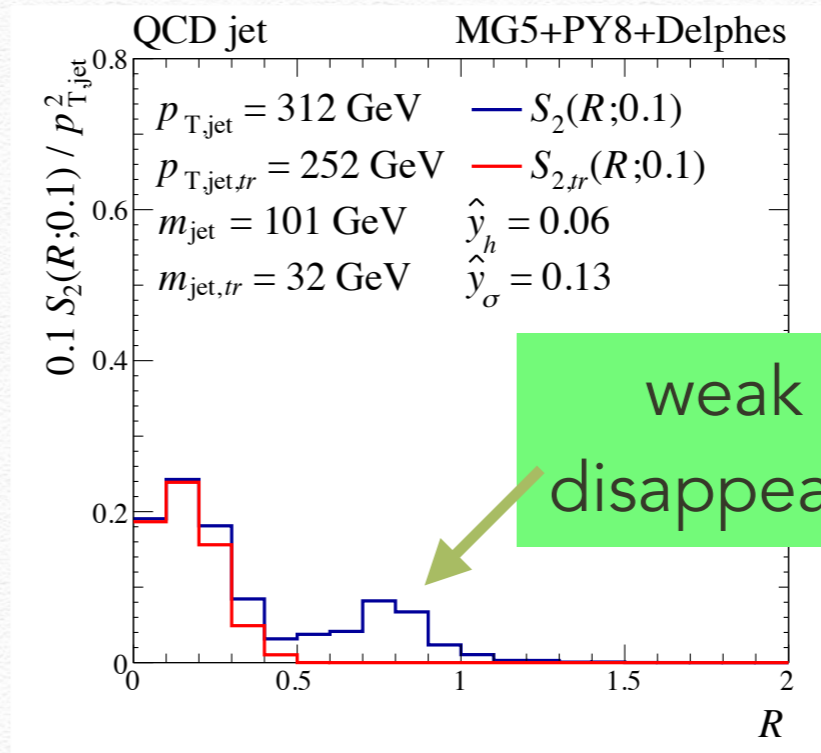
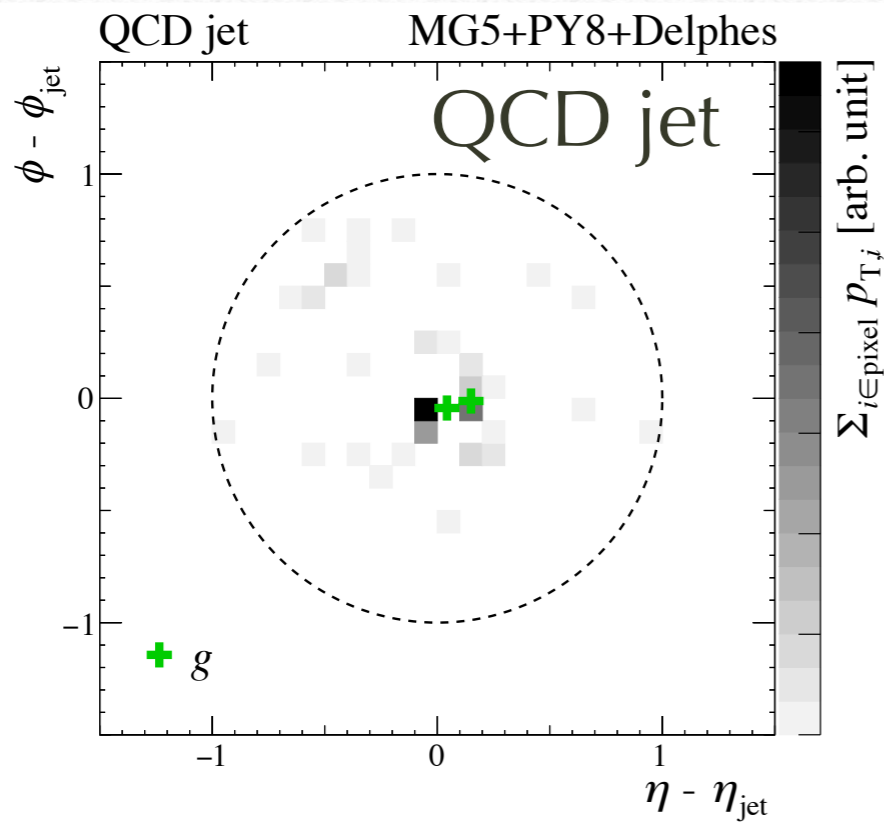
= "hard-hard" correlation

$S_{2\text{soft}}$

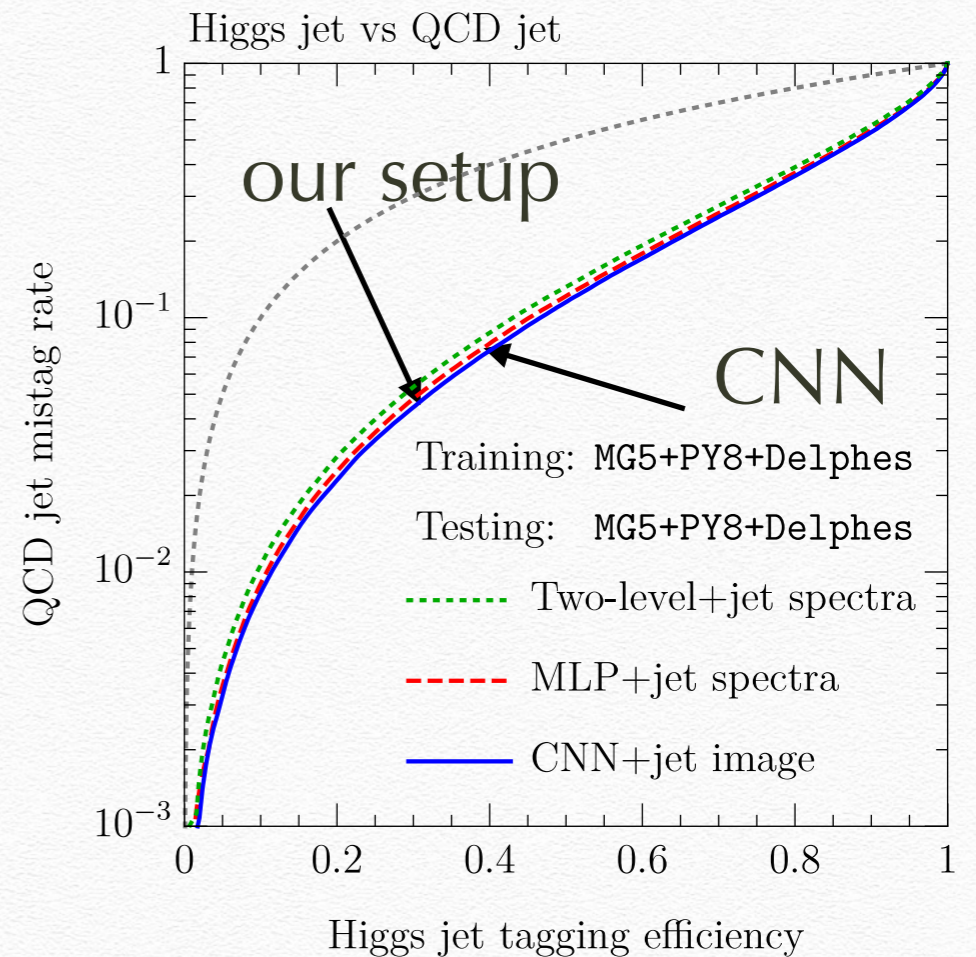
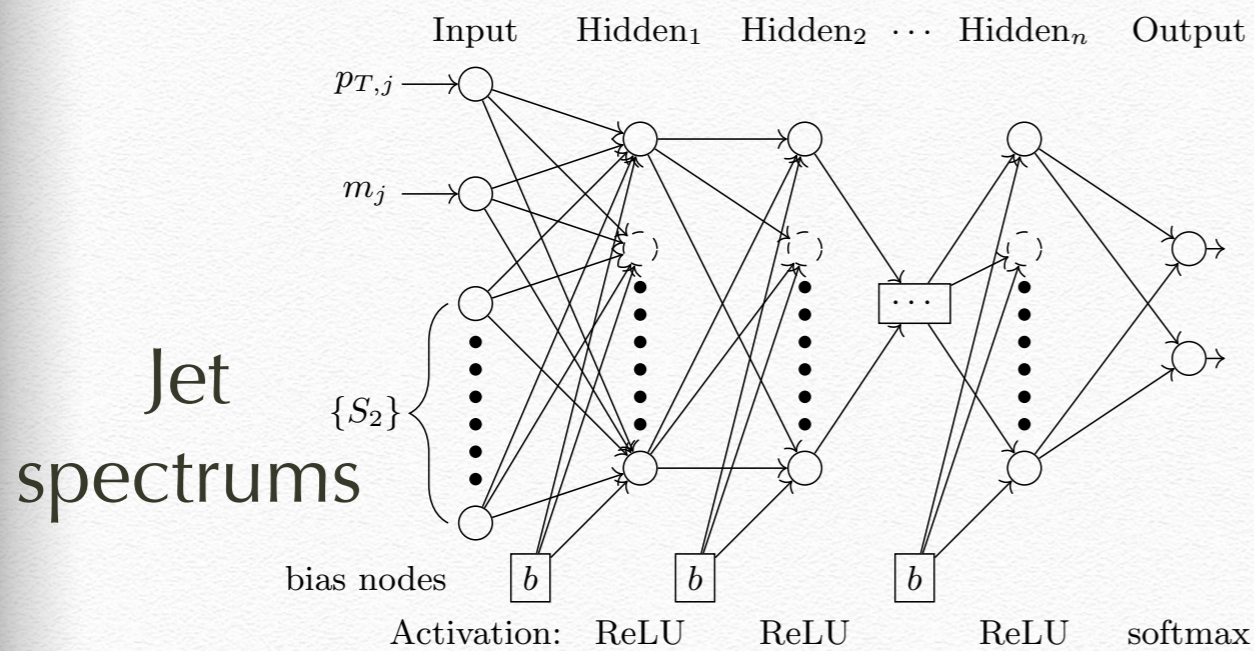
$$S_{2,\text{soft}}(R; \Delta R) = S_2(R; \Delta R) - S_{2,\text{trim}}(R; \Delta R).$$

= "soft-hard correlation" + "soft-soft" correlation

ジェットスペクトル $S_2(R)$ の例



<CNN> vs <DNN with $S_{2\text{trim}}(R)$ and $S_{2\text{soft}}(R)$ >



for $300\text{GeV} < p_T < 400\text{GeV}$ and $100\text{GeV} < m_j < 150\text{GeV}$

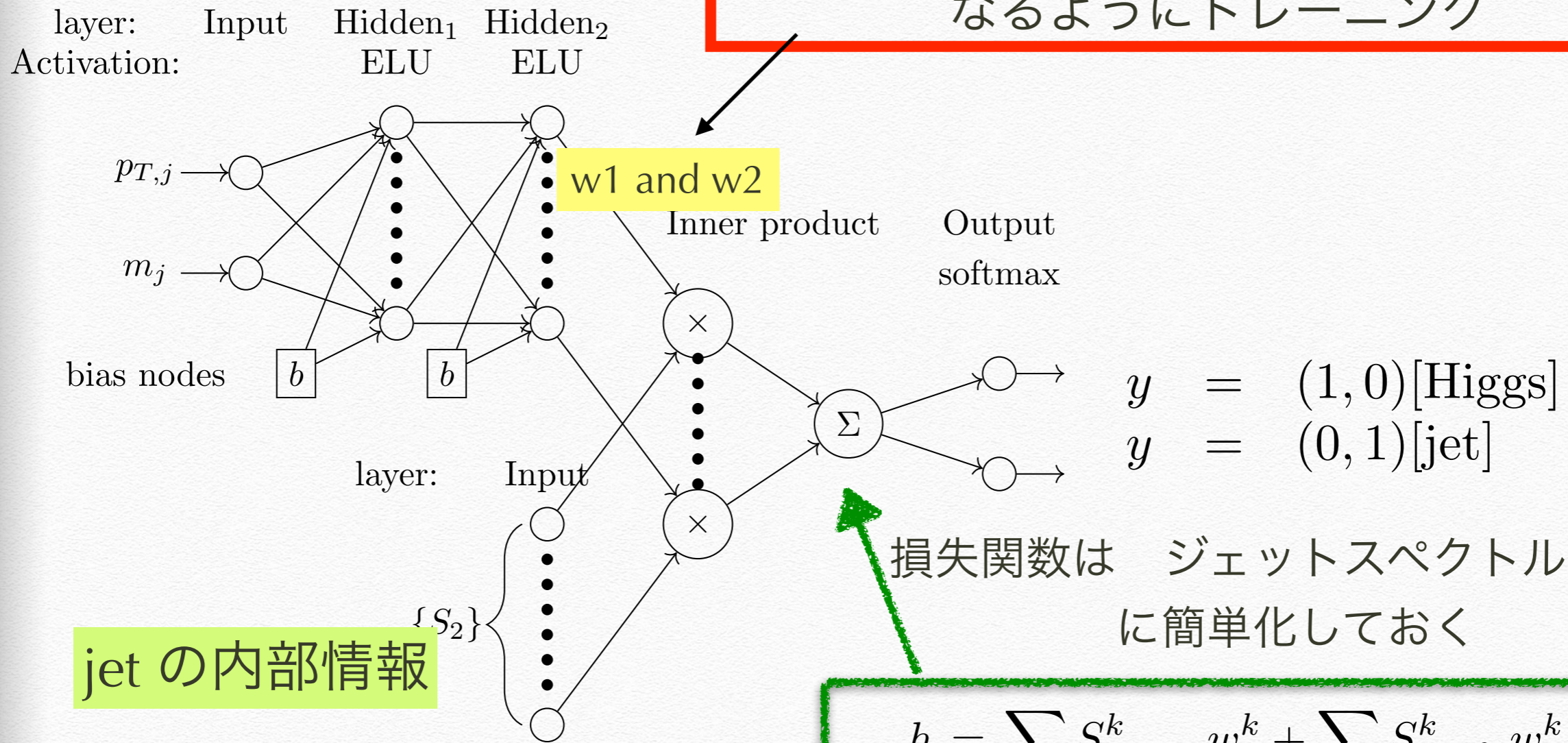
ほぼ同じ効率

- ❖ ほとんど2点相関で分類できるのではないかと (Or NN with jet image "find" two-point correlation by itself (not proven))
- ❖ 入力も少なくなったので、それぞれの寄与を調べることもできるのではないかと

[運動量処理系]・[radiation 処理系]

= [解釈可能モデル]

w1 and w2(スペクトル係数) を分類が正しくなるようにトレーニング

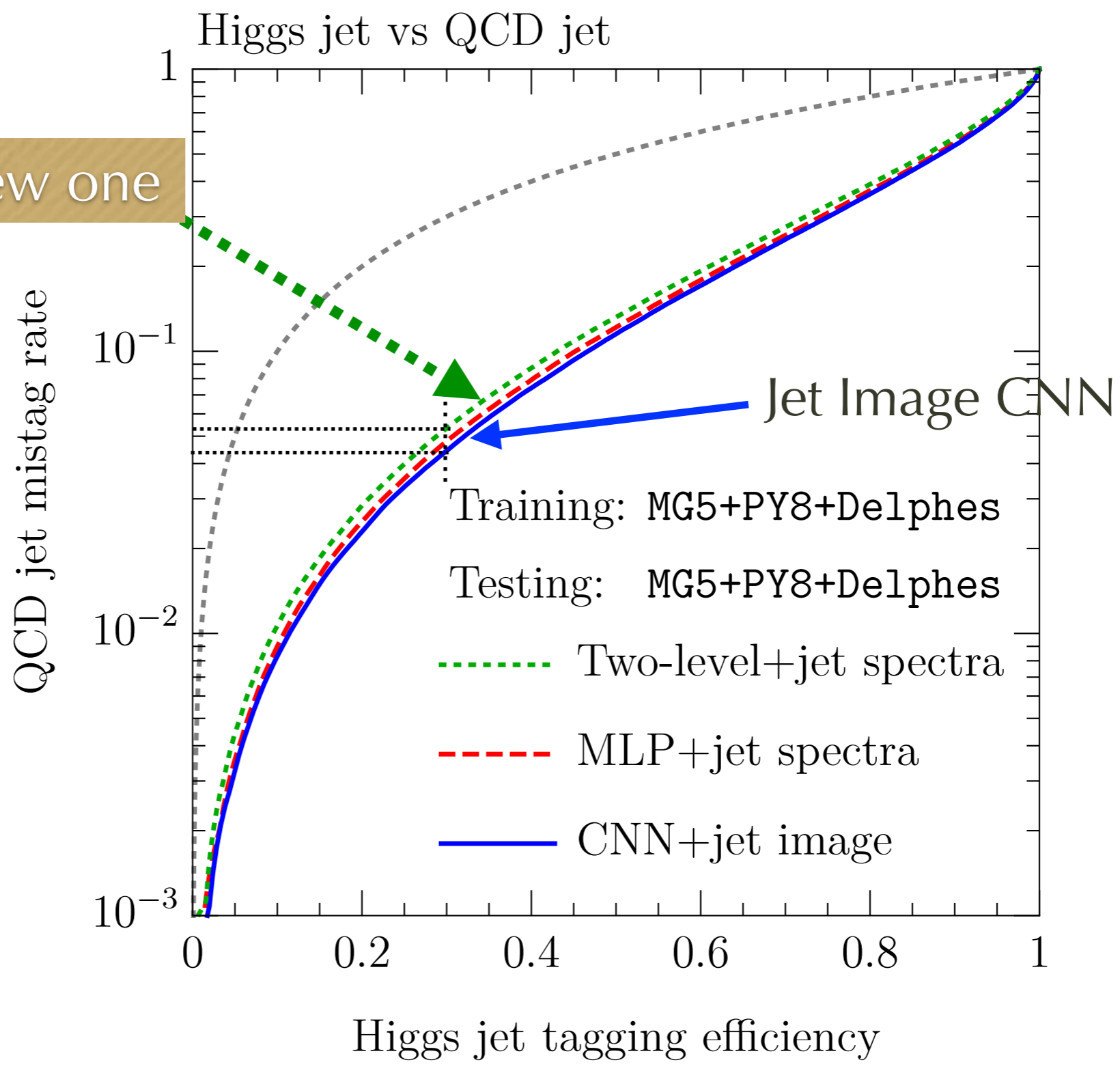


損失関数は ジェットスペクトルの一次に簡単化しておく

$$h = \sum_k S_{2,trim}^k w_1^k + \sum_k S_{2,soft}^k w_2^k,$$

この効率も元とほとんど変わらない

New one



損失関数の「係数」 w の意味

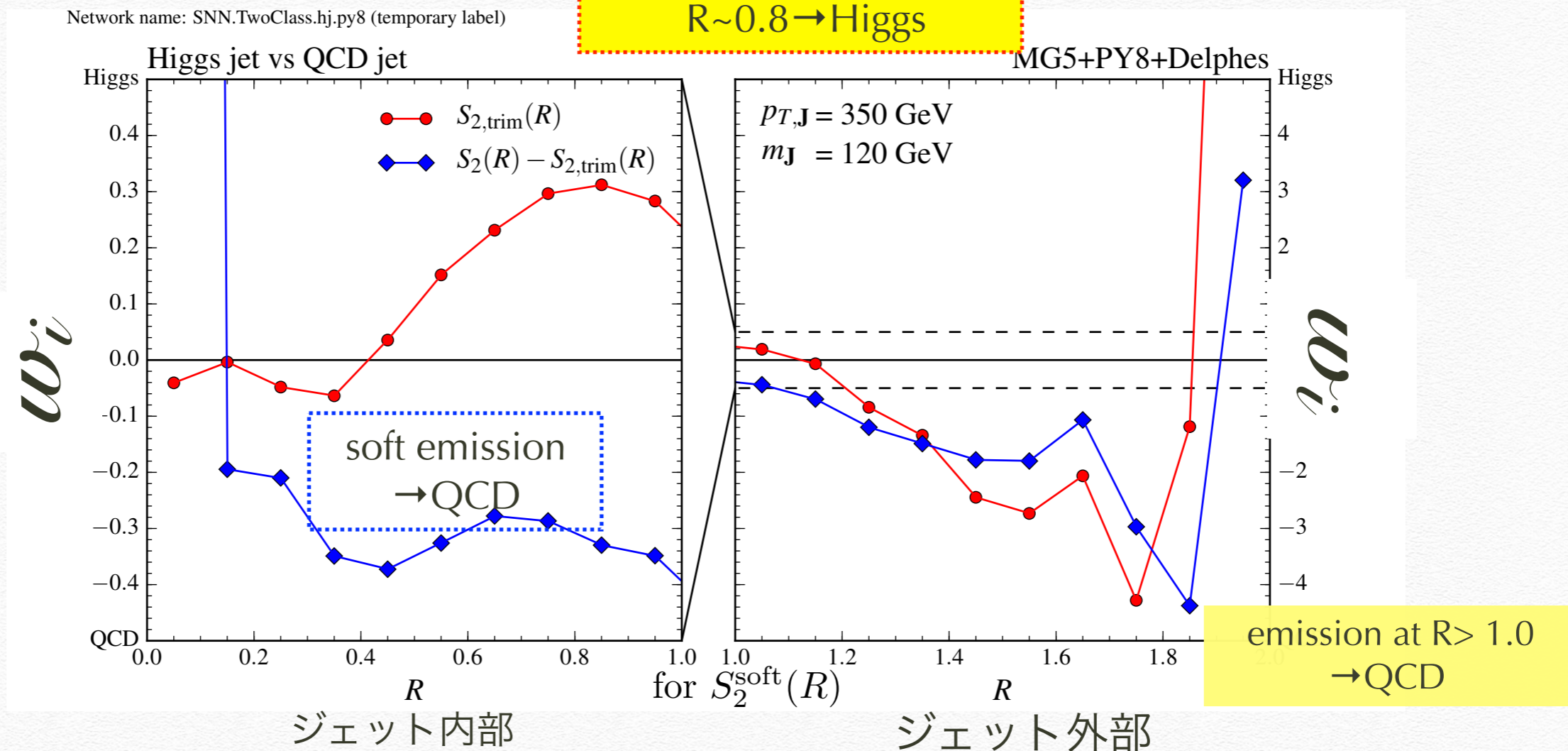
❖ 特定の R における deposit の分類に対する重要性 = w の絶対値

$$h = \sum_k S_{2,\text{trim}}^k w_1^k + \sum_k S_{2,\text{soft}}^k w_2^k,$$

coefficient w_1 and w_2 is instructed to depend on mass and momentum but not jet spectrum

QCD jet vs Higgs jet

trimmed jet emission at $R \sim 0.8 \rightarrow$ Higgs



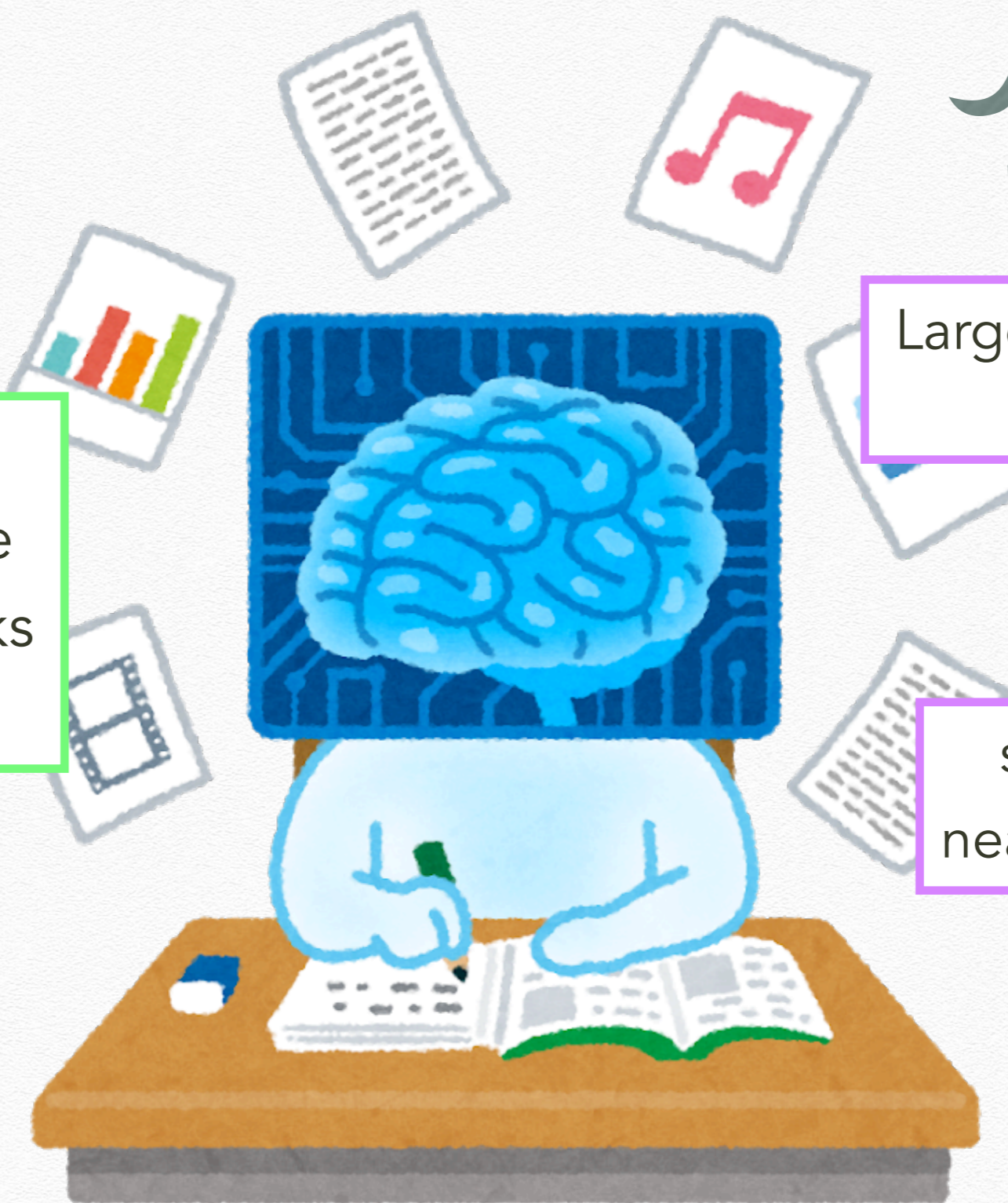
"Reasons" behind jet classification



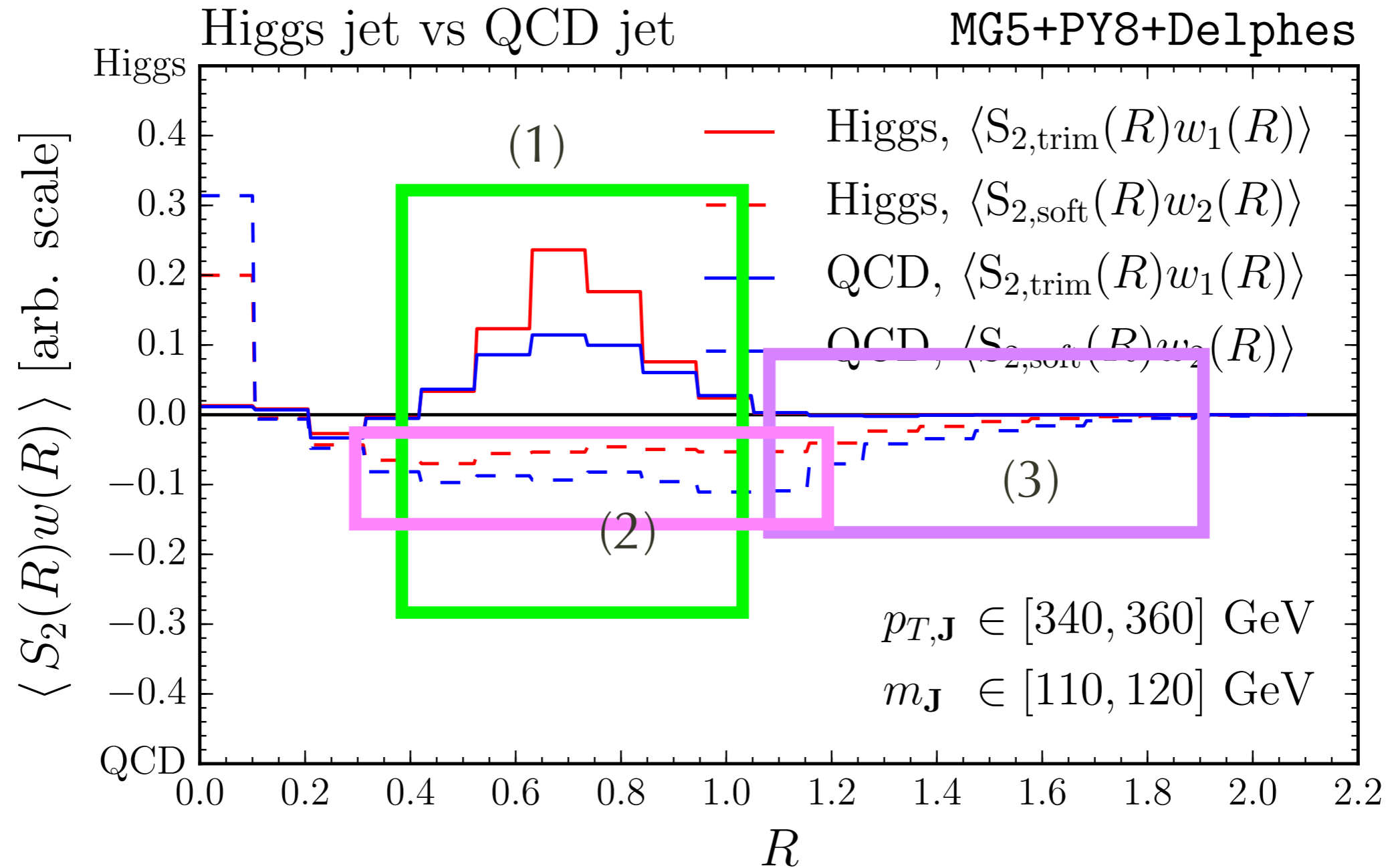
two hard
substructure
making peaks
in $S_{2\text{trim}}$

Large Hard-and-soft
correlation

soft emission
near jet boundary



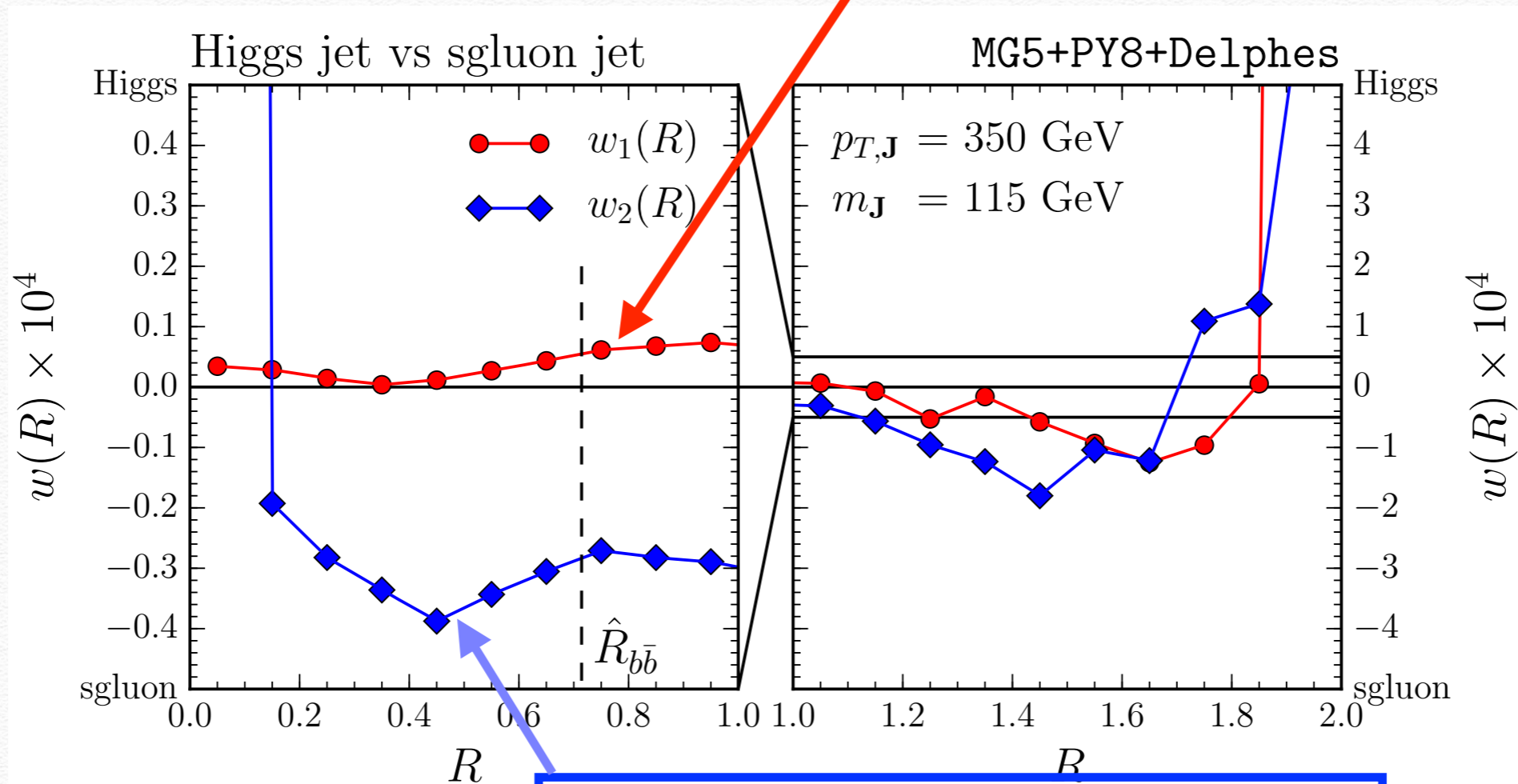
平均的な損失関数 $\langle h \rangle = \langle S_2 w \rangle$



親粒子のカラー依存性

- ❖ Higgs (singlet scalar) vs color octet particle. 質量は同じで 2b に崩壊するようにしておく

質量が同じなので $S_{2\text{trim}}$ にピークは出ない

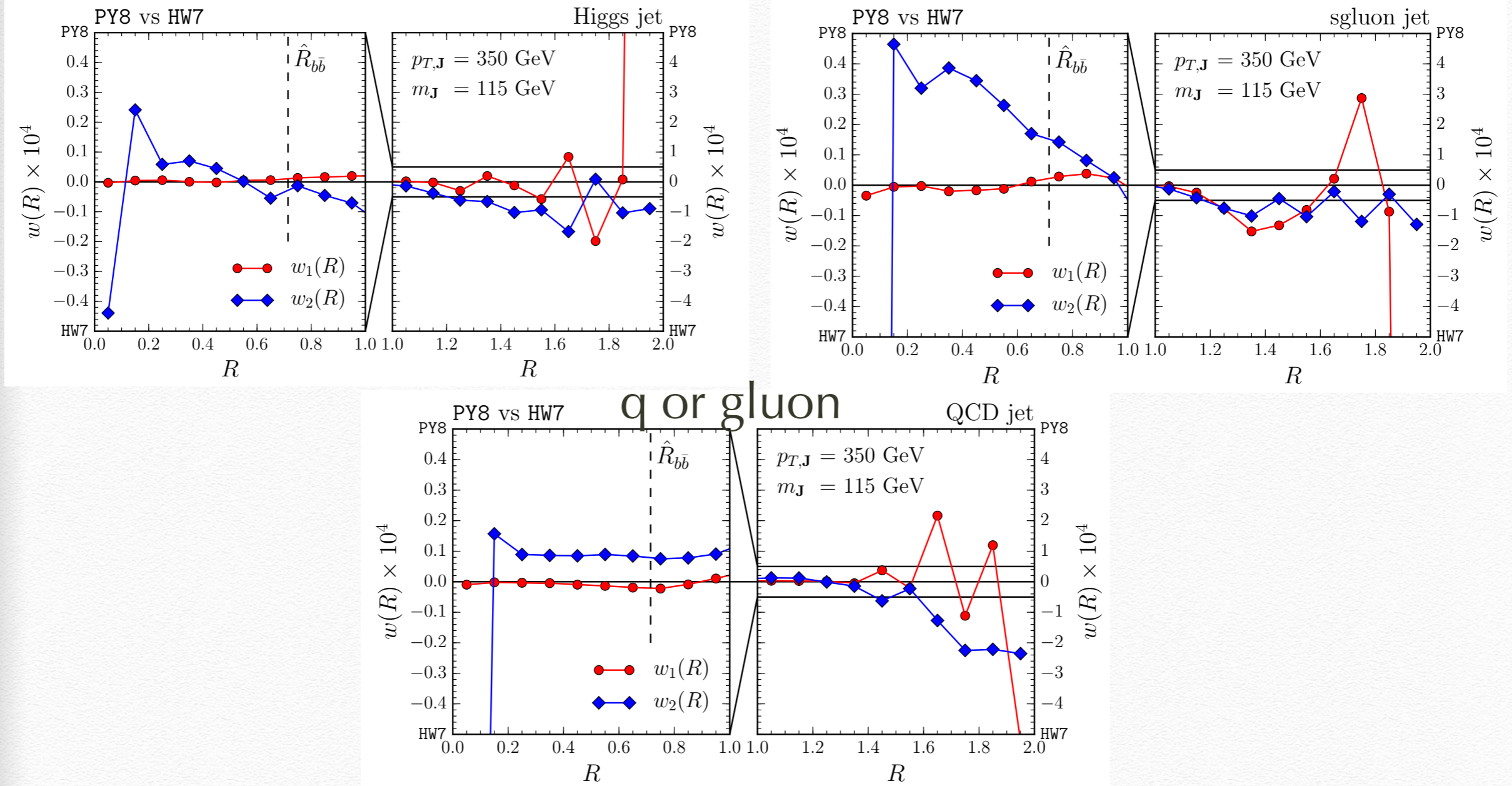


trimmed jet の外のアクティビティは カラー 8 重項っぽい

MC を比べてみる (PY8 vs HW7)

Higgs boson

Scalar gluon



* Herwig のほうが大きな R まで出している。(angular ordering vs pT ordering?)

* Trimmed distribution は同じで信頼できる (Hard(er) Physics is common)

いままでこういう比較をしたことがなかったので少しびっくりした。

まとめとか今後とか

まずは普通の物理

- ❖ **ジェットの分類におけるコアな量を見つけた。Input: は $N \times N$ から N ($N \sim 20$ for our case.) くらいにはなった。もちろん計算も早い。**
- ❖ **ソフトなアクティビティを分類に使っているようだ。**（色々心配）
- ❖ **実際のイベントはMC と違うので、もっと面白いことがあるかもしれない。**
- ❖ **トップの場合は？**

機械学習として

- ❖ 犬と猫の分類の場合何を原理として判断しているかというのは、定量化しにくい。
- ❖ ジェットの物理の場合は効率のよいインプットがなにか物理的な推測で決めることができる。
- ❖ CNN は本当の答えに収束していくプロセスを、「正解」を利用して追跡できないか
- ❖ CNN はジェットスペクトルを見つけるために計算しているのか？ → **直接的には示していない**ので別のモデルを作って調べようとしている。