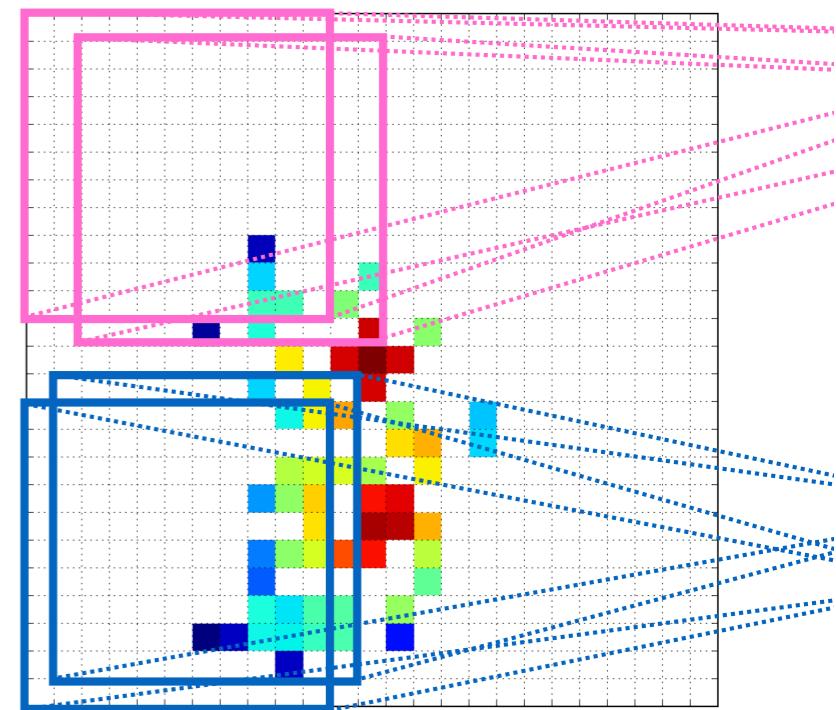


A model agnostic machine learning search for resonant new physics

Benjamin Nachman

Lawrence Berkeley National Laboratory

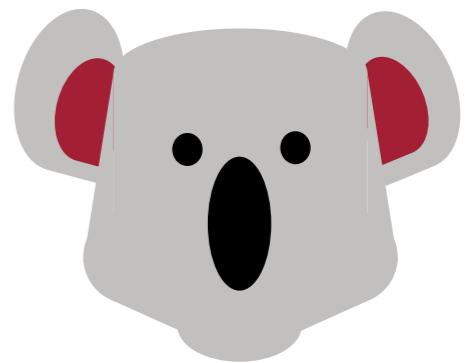


UCI Seminar
January 15, 2019

Outline

- Why model agnostic searches?

Intermezzo: Learning directly from data



- Why it is important to learn directly from data
- How to learn without labels
- BSM Searches background / signal model independent

Most searches at the LHC
(with or without ML)

Data versus SM simulation
("general search")

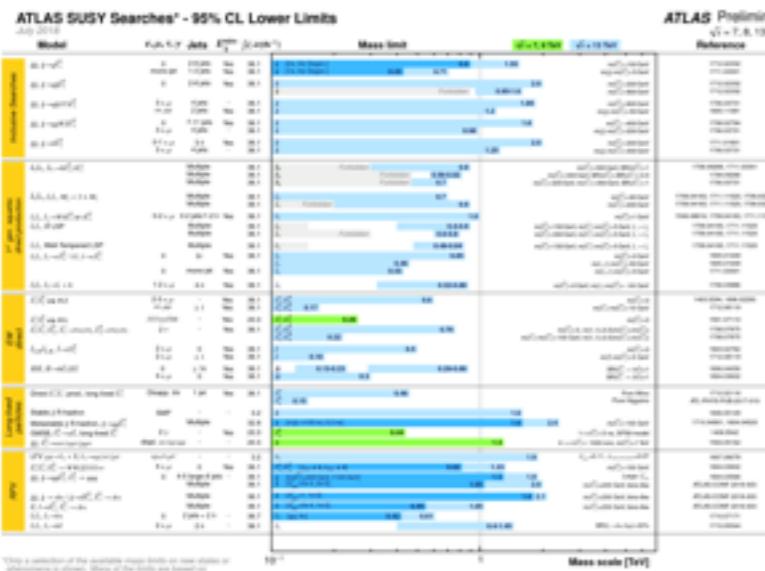
CWoLa Hunting
(focus of this talk)

Autoencoders

- The future

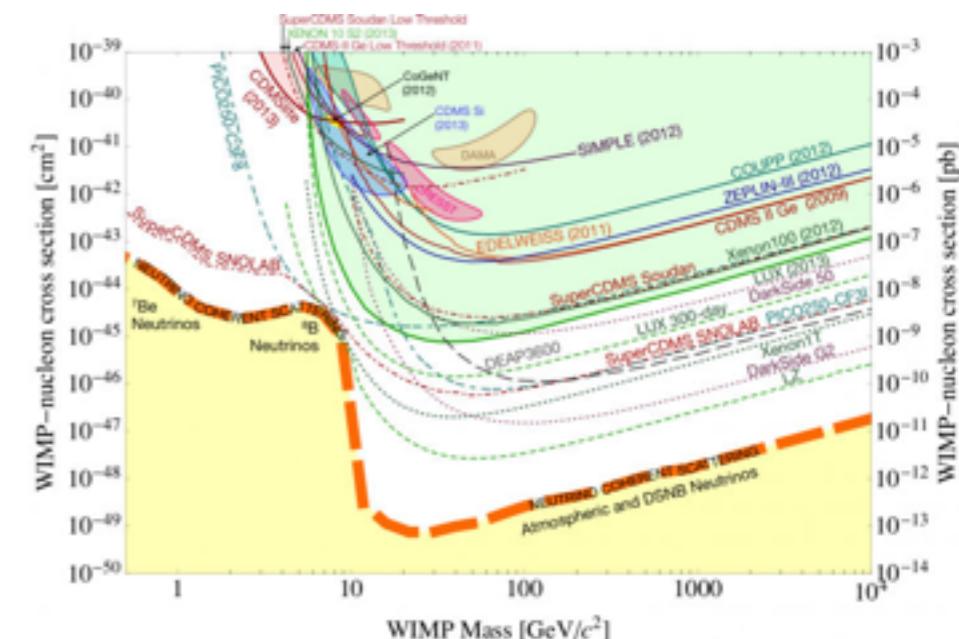
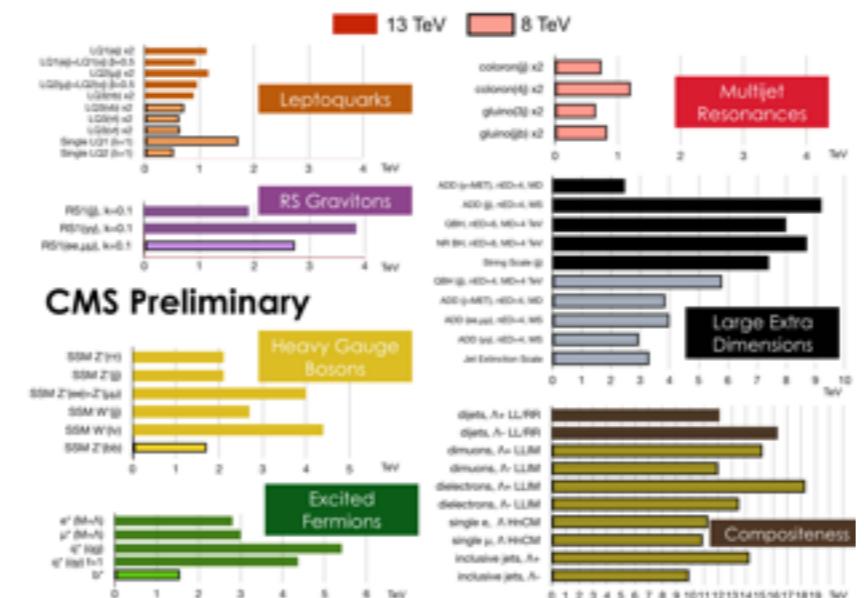
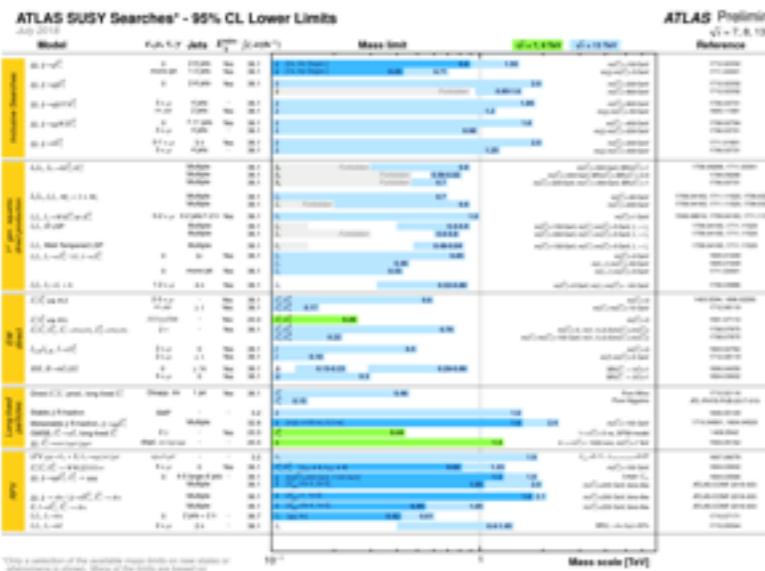
Why model agnostic?

Fact: no conclusive evidence for BSM in HEP*



Why model agnostic?

Fact: no conclusive evidence for BSM in HEP*



There is no BSM (SM all the way to Planck)

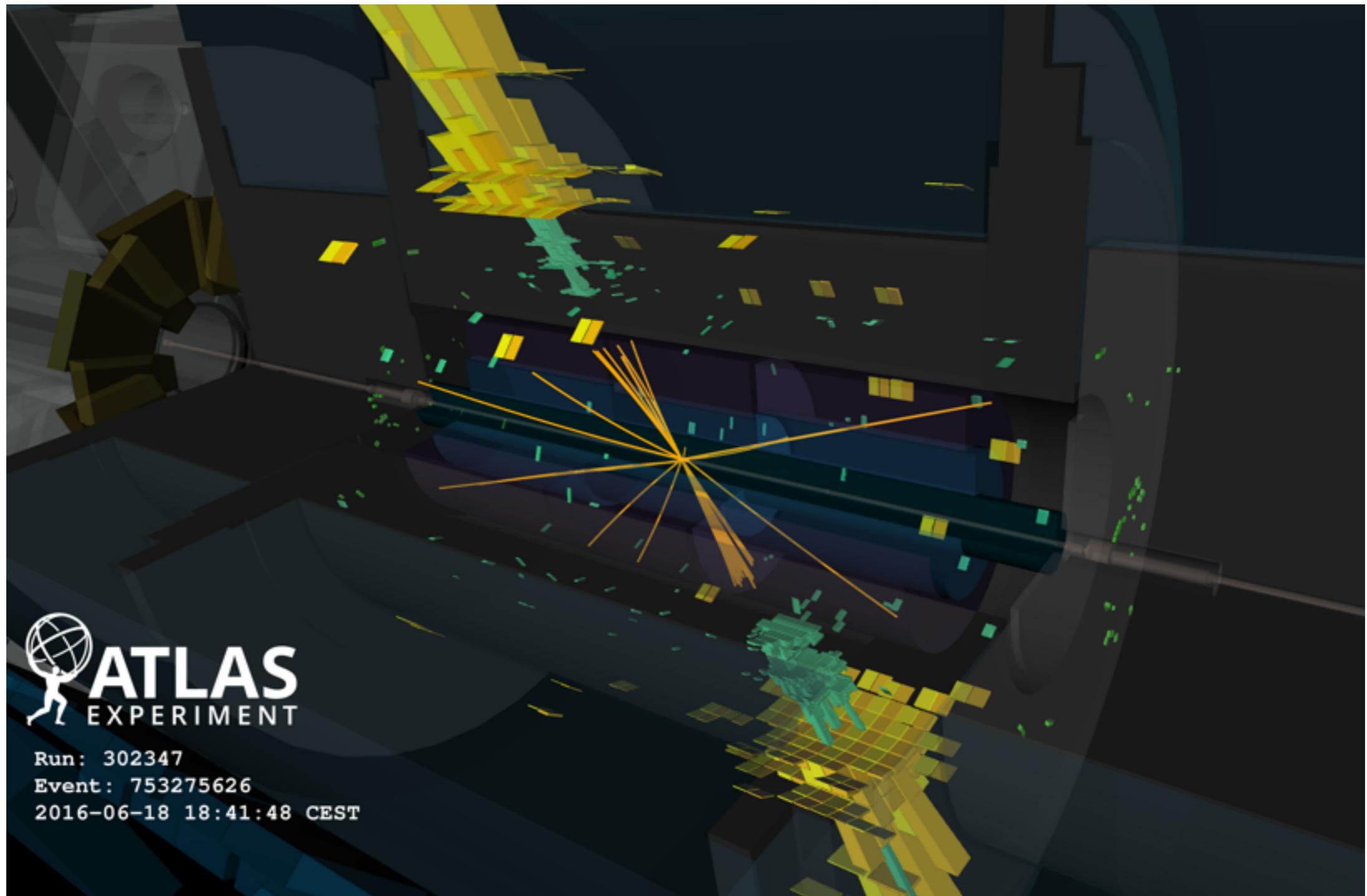
Three possibilities

Patience! (BSM is just rare)

We are just not looking in the right place

*neutrinos are *technically* BSM, but depending on their origin of their mass, may not lead to “new” short-distance physics. We do know there is dark matter, but we don’t even know if it is composed of particles (such as WIMPs). There are also multiple B physics anomalies, but they are not conclusive.

Consider the collider di-object resonance search.



Signatures, not models

Consider the collider di-object resonance search.

	e	μ	τ	γ	j	b	t	W	Z	h	
e	$\pm\mp[4], \pm\pm[5]$	$\pm\pm[5, 6]$	$\pm\pm[6, 7]$	[7]	\emptyset	\emptyset	\emptyset	\emptyset	\emptyset	\emptyset	
μ			$\pm\mp[4], \pm\pm[5]$	[7]	\emptyset	\emptyset	\emptyset	\emptyset	\emptyset	\emptyset	
τ				[8]	\emptyset	\emptyset	\emptyset	[9]	\emptyset	\emptyset	
γ					[10]	[11–13]	\emptyset	\emptyset	[14]	[14]	\emptyset
j					[15]	[16]	[17]	[18]	[18]	[18]	\emptyset
b						[16]	[19]	\emptyset	\emptyset	\emptyset	
t							[20]	[21]	\emptyset	\emptyset	
W								[22–25]	[23, 24, 26, 27]	[28–30]	
Z									[23, 25, 31]	[28, 30, 32, 33]	
h										[34–37]	

$p p \rightarrow X \rightarrow \text{SM SM}$

Signatures

1610.09392 (Yvonne Ng, D. Whiteson, et al.)

	e	μ	τ	γ	j	b	t	W	Z	h
e	$Z', H^{\pm\pm}$	$\cancel{R}, H^{\pm\pm}$	$\cancel{R}, H^{\pm\pm}$	L^*	LQ, \cancel{R}	LQ, \cancel{R}	LQ, \cancel{R}	L^*, ν_{KK}	L^*, e_{KK}	L^*
μ		$Z', H^{\pm\pm}$	$\cancel{R}, H^{\pm\pm}$	L^*	LQ, \cancel{R}	LQ, \cancel{R}	LQ, \cancel{R}	L^*, ν_{KK}	L^*, μ_{KK}	L^*
τ			$Z', H, H^{\pm\pm}$	L^*	LQ, \cancel{R}	LQ, \cancel{R}	LQ, \cancel{R}	L^*, ν_{KK}	L^*, τ_{KK}	L^*
γ				H, G_{KK}, Q	Q^*	Q^*	Q^*	W_{KK}, Q	H, Q	Z_{KK}
j					Z', ρ, G_{KK}	W', \cancel{R}	T', \cancel{R}	Q^*, Q_{KK}	Q^*, Q_{KK}	Q'
b						Z', H	W', \cancel{R}, H^\pm	T', Q^*, Q_{KK}	Q^*, Q_{KK}	B'
t							H, G', Z'	T'	T'	T'
W								H, G_{KK}, ρ	W', Q	H^\pm, Q, ρ
Z									H, G_{KK}, ρ	A, ρ
h										H, G_{KK}

Models

There are many uncovered combinations !

(since the writing of this paper, some have been filled in)

Signatures, not models

Consider the collider di-object resonance search.

	e	μ	τ	γ	j	b	t	W	Z	h	BSM
e	$\pm\mp[4], \pm\pm[5]$	$\pm\pm[5, 6]$	$\pm\mp[6, 7]$	[7]	\emptyset	\emptyset	\emptyset	\emptyset	\emptyset	\emptyset	\emptyset
μ				[7]	\emptyset	\emptyset	\emptyset	\emptyset	\emptyset	\emptyset	\emptyset
τ				[8]	\emptyset	\emptyset	\emptyset	[9]	\emptyset	\emptyset	\emptyset
γ					[10]	[11–13]	\emptyset	\emptyset	[14]	[14]	\emptyset
j						[15]	[16]	[17]	[18]	[18]	\emptyset
b							[16]	[19]	\emptyset	\emptyset	\emptyset
t							[20]	[21]	\emptyset	\emptyset	\emptyset
W								[22–25]	[23, 24, 26, 27]	[28–30]	
Z									[23, 25, 31]	[28, 30, 32, 33]	
h										[34–37]	
Signatures											
BSM											

For a generic $p p \rightarrow X \rightarrow B C$, mostly uncovered!

Some of the BSM elements are partially covered with e.g. XH and the photon jets search.

Training Directly on (unlabeled) Data

I hope I've convinced you that it is important to have an open mind when looking at data.

Now, let's take a step back and discuss why it is important to train on data **in general**.

Intermezzo: Learning directly from data

- Why it is important to learn directly from data
- How to learn without labels



The importance of training on data

Let's start with a small thought experiment.

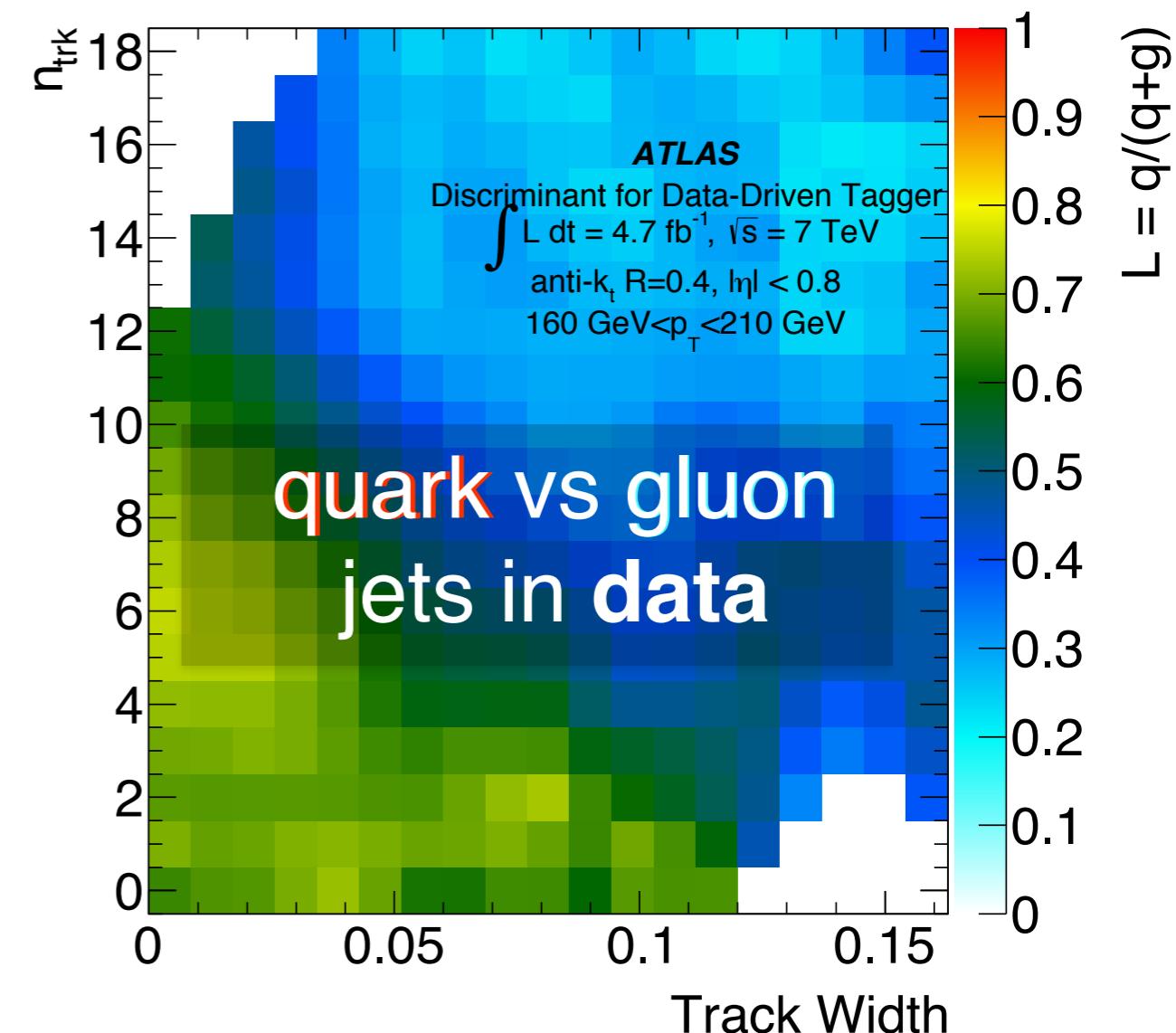
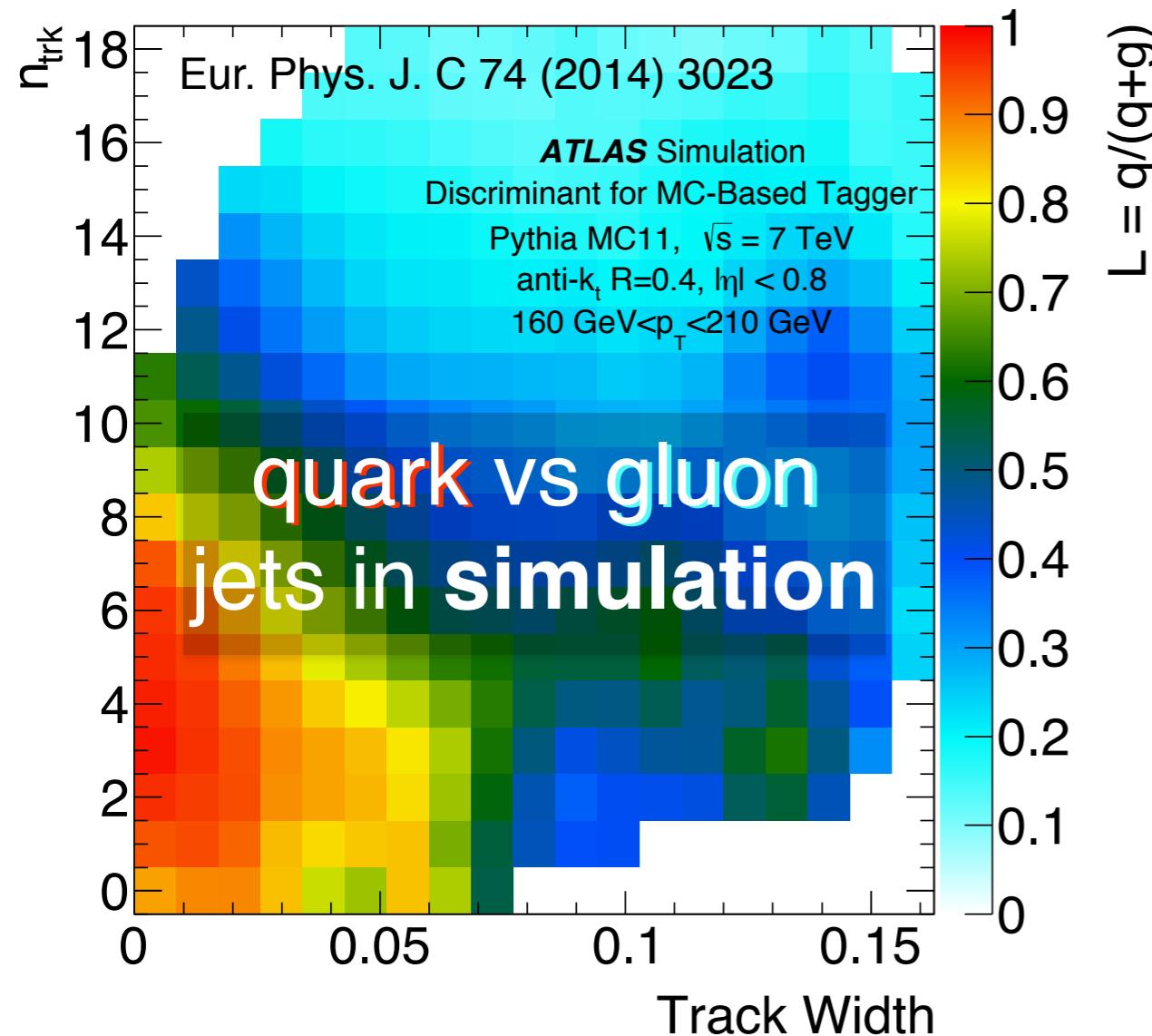
Are these equivalent?

- (1) train in simulation & calibrate in data
- (2) train in data

← what we mostly do right now

Background and Motivation

Usual paradigm: train in simulation, validate on data, test on data.

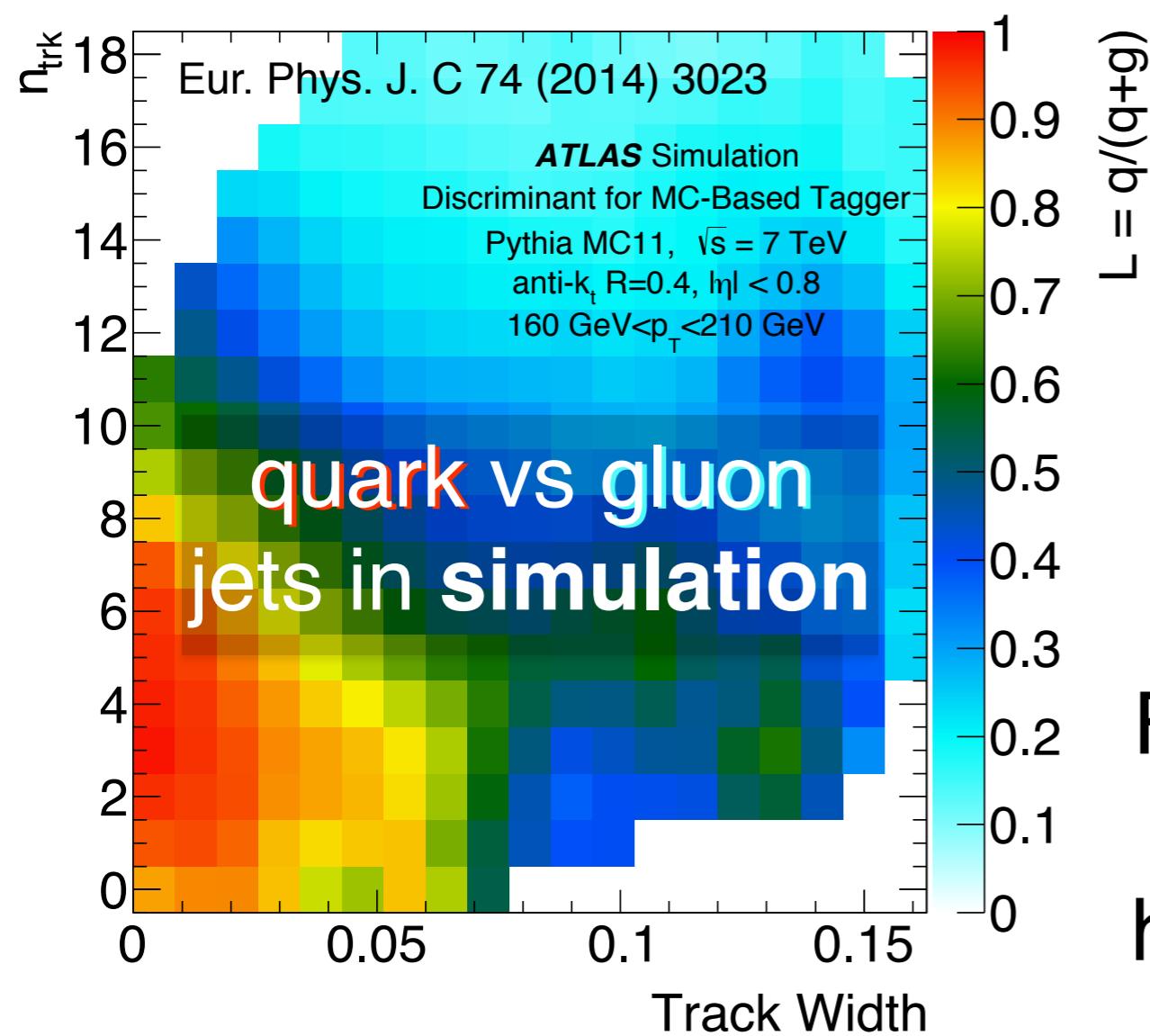


If data and simulation differ, this is sub-optimal!

(answer from previous slide is no!)

Background and Motivation

Usual paradigm: **train in simulation**, validate on data, test on data.



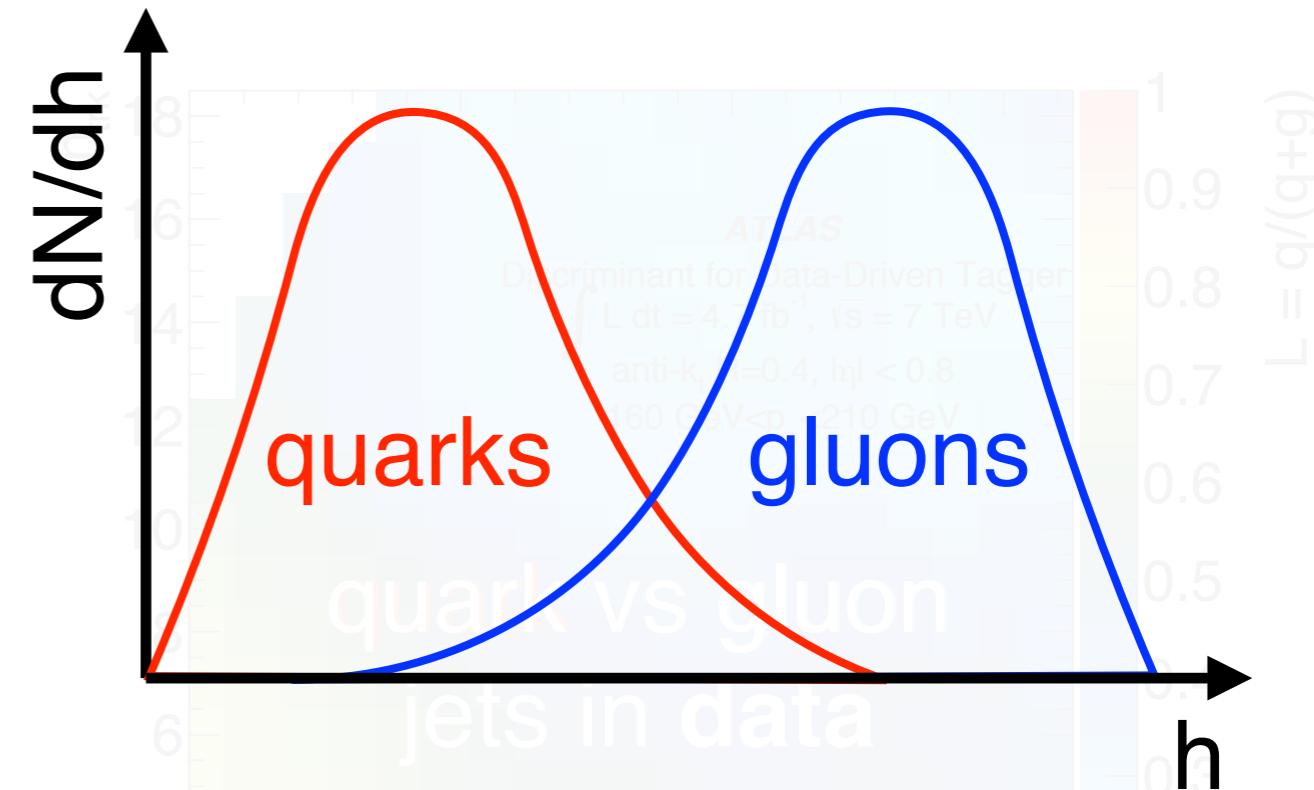
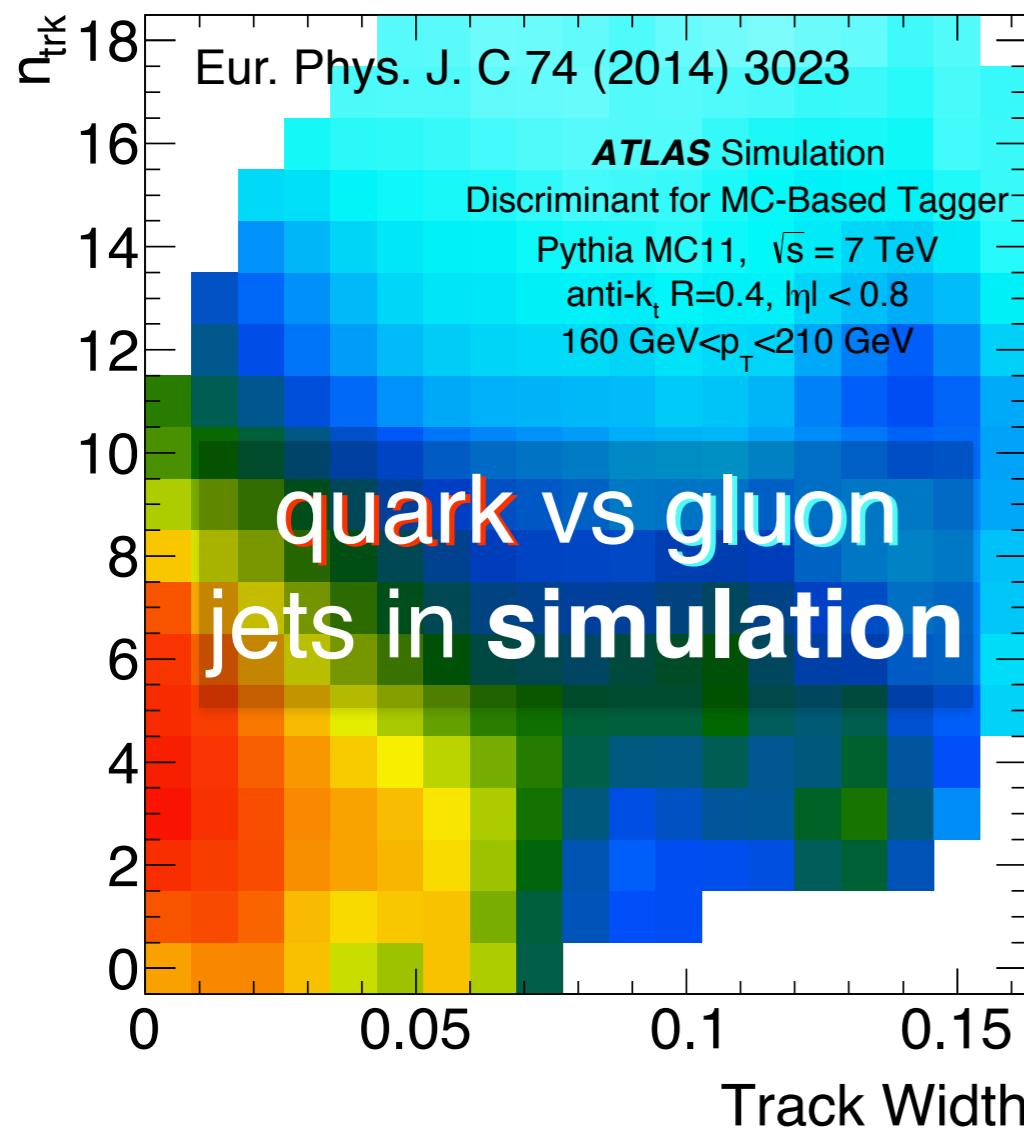
Recall: optimal classifier (by Neyman-Pearson NP) is a threshold cut on the likelihood ratio.

For a 2D feature space, no need for a NN or BDT - can use a histogram to “train” the **classifier**.

$$h(n_{\text{trk}}, \text{Track Width}) \rightarrow [0, 1]$$

Background and Motivation

Usual paradigm: **train in simulation**, validate on data, test on data.

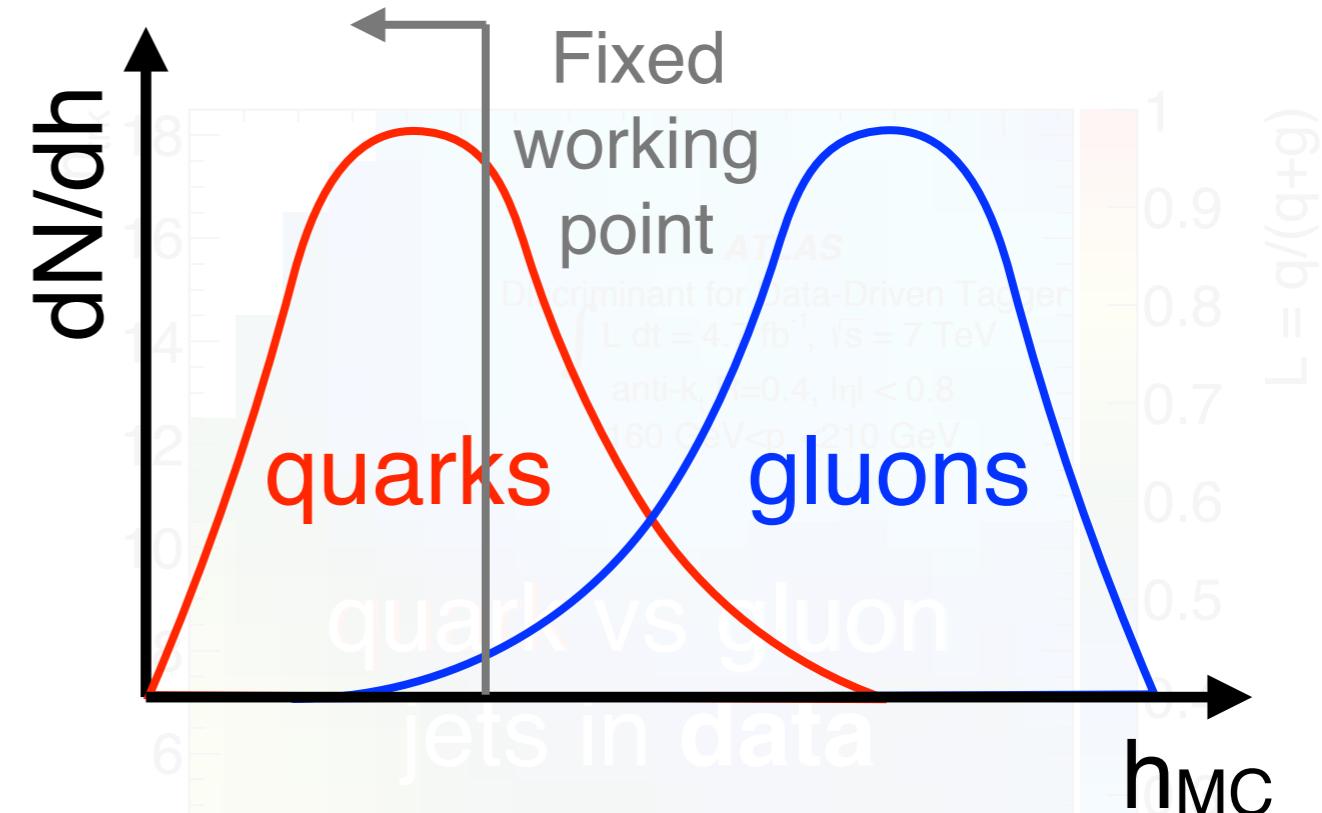
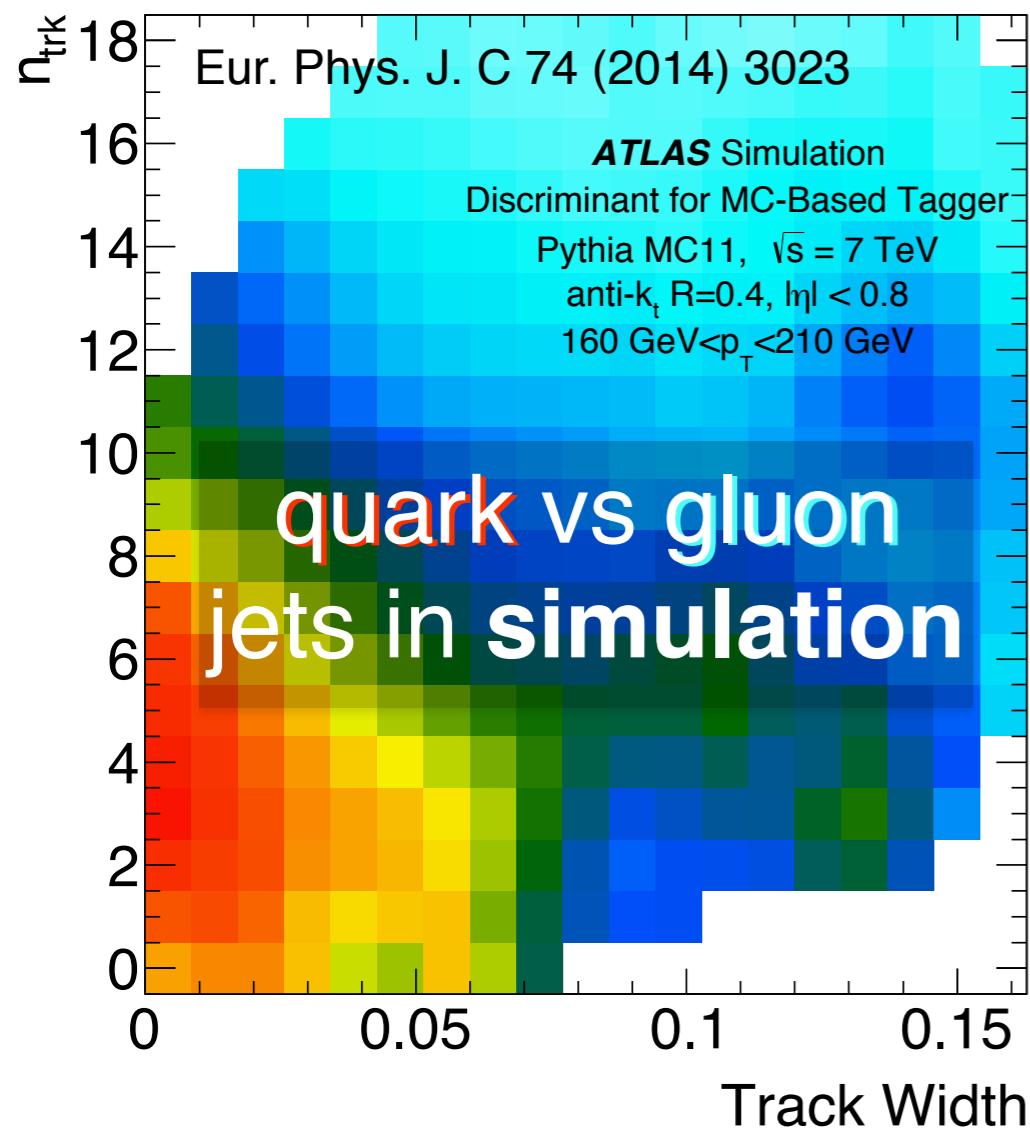


For a 2D feature space, no need for a NN or BDT - can use a histogram to “train” the **classifier**.

$$h(n_{\text{trk}}, \text{Track Width}) \rightarrow [0, 1]$$

Background and Motivation

Usual paradigm: **train in simulation**, validate on data, test on data.

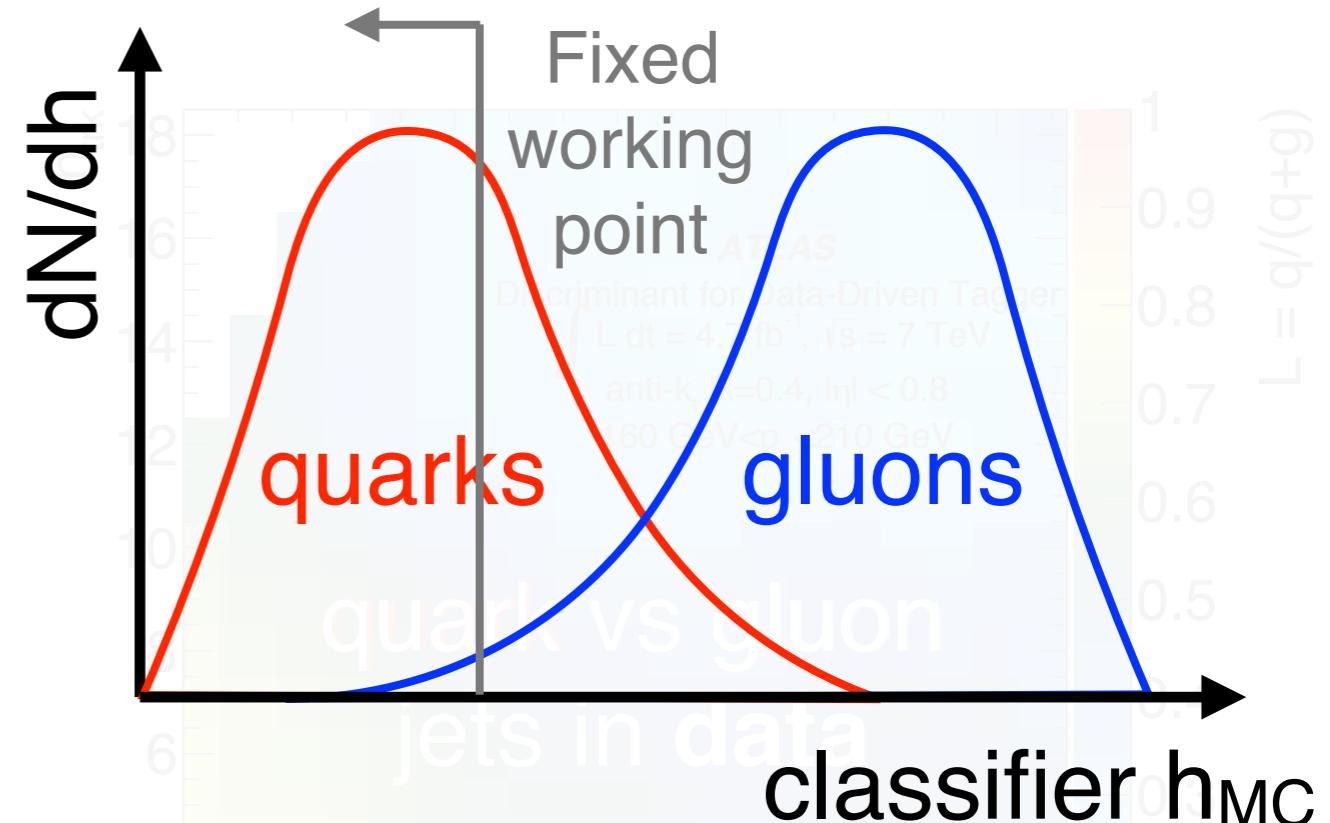
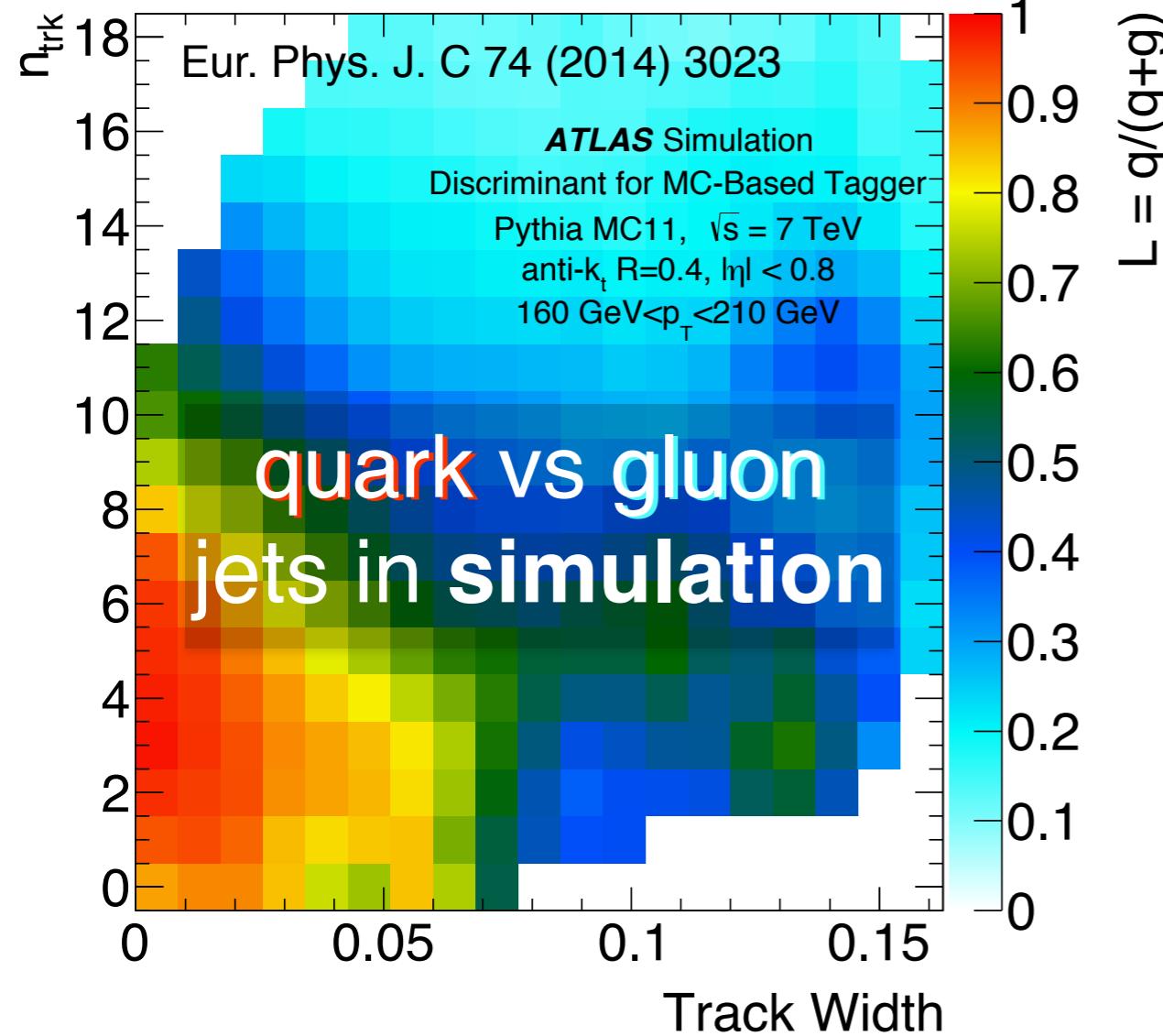


For a 2D feature space, no need for a NN or BDT - can use a histogram to “train” the **classifier**.

$$h_{\text{MC}}(n_{\text{trk}}, \text{Track Width}) \rightarrow [0, 1]$$

Background and Motivation

Usual paradigm: **train in simulation**, validate on data, test on data.



WP in simulation:
 $\epsilon_{\text{signal,MC}}, \epsilon_{\text{back,MC}}$

Background and Motivation

Usual paradigm: train in simulation, **validate on data**, test on data.

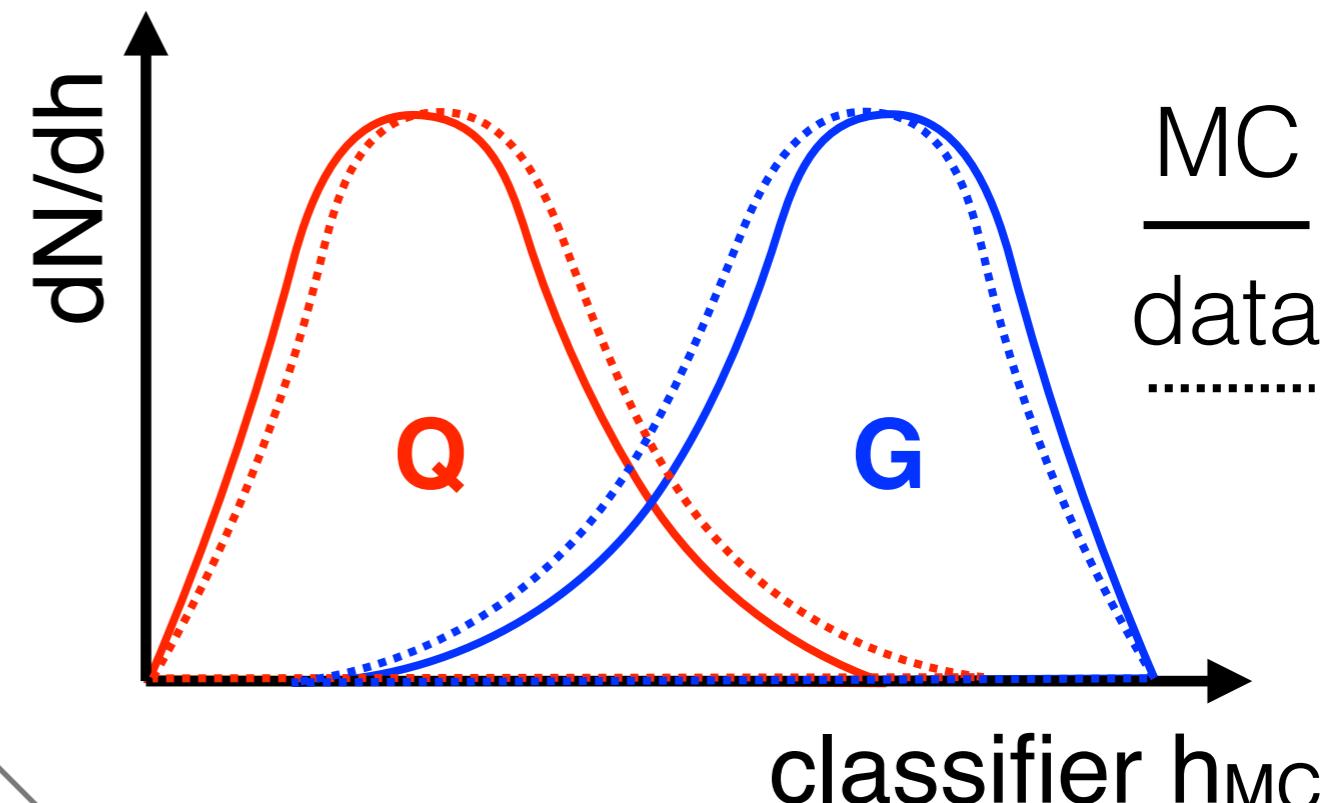
Determine the performance
of the WP in data.
Eur. Phys. J. C 74 (2014) 3023

How did we get this?

$$\text{dijets} = f_q \times \mathbf{Q} + (1-f_q) \times \mathbf{G}$$

$$Z+\text{jets} = g_q \times \mathbf{Q} + (1-g_q) \times \mathbf{G}$$

2 equations, 2 unknowns (\mathbf{Q} , \mathbf{G})



two event samples with
different q/g fractions

(N.B. f & g from simulation and selection can't bias Q and G)

Background and Motivation

Usual paradigm: train in simulation, **validate on data**, test on data.

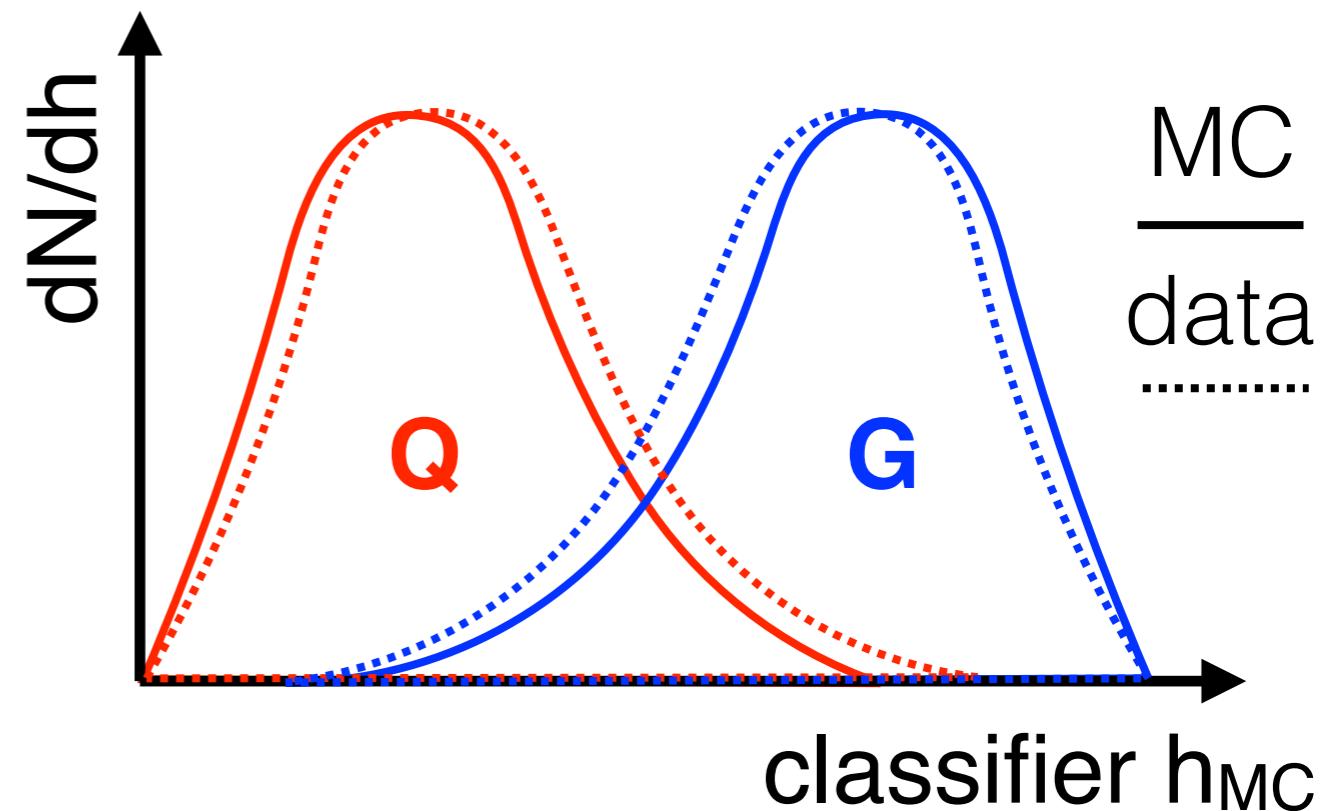
Determine the performance of the WP in data.
Eur. Phys. J. C 74 (2014) 3023

How did we get this?

$$\text{dijets} = f_q \times \mathbf{Q} + (1-f_q) \times \mathbf{G}$$

$$Z+\text{jets} = g_q \times \mathbf{Q} + (1-g_q) \times \mathbf{G}$$

2 equations, 2 unknowns (**Q**, **G**)



WP in data:

$\epsilon_{\text{signal,data}}, \epsilon_{\text{back,data}}$

Can correct the MC to have the same performance as data.

Background and Motivation

Usual paradigm: train in simulation, validate on data, **test on data**.

Once we have scale factors (& their uncertainty), we can ensure that our analysis will be **accurate**.

...so what is the problem?

remember my claim from earlier:

If data and simulation differ, this is sub-optimal!

This is an **accuracy** versus **precision** problem. It is “easy” to achieve accuracy through calibration, but the results may not be the best one possible.

Background and Motivation

In this 2D feature space, we can actually derive h_{data} .

Using the same trick as earlier:

$$\text{dijets} = f_q \times Q + (1-f_q) \times G$$

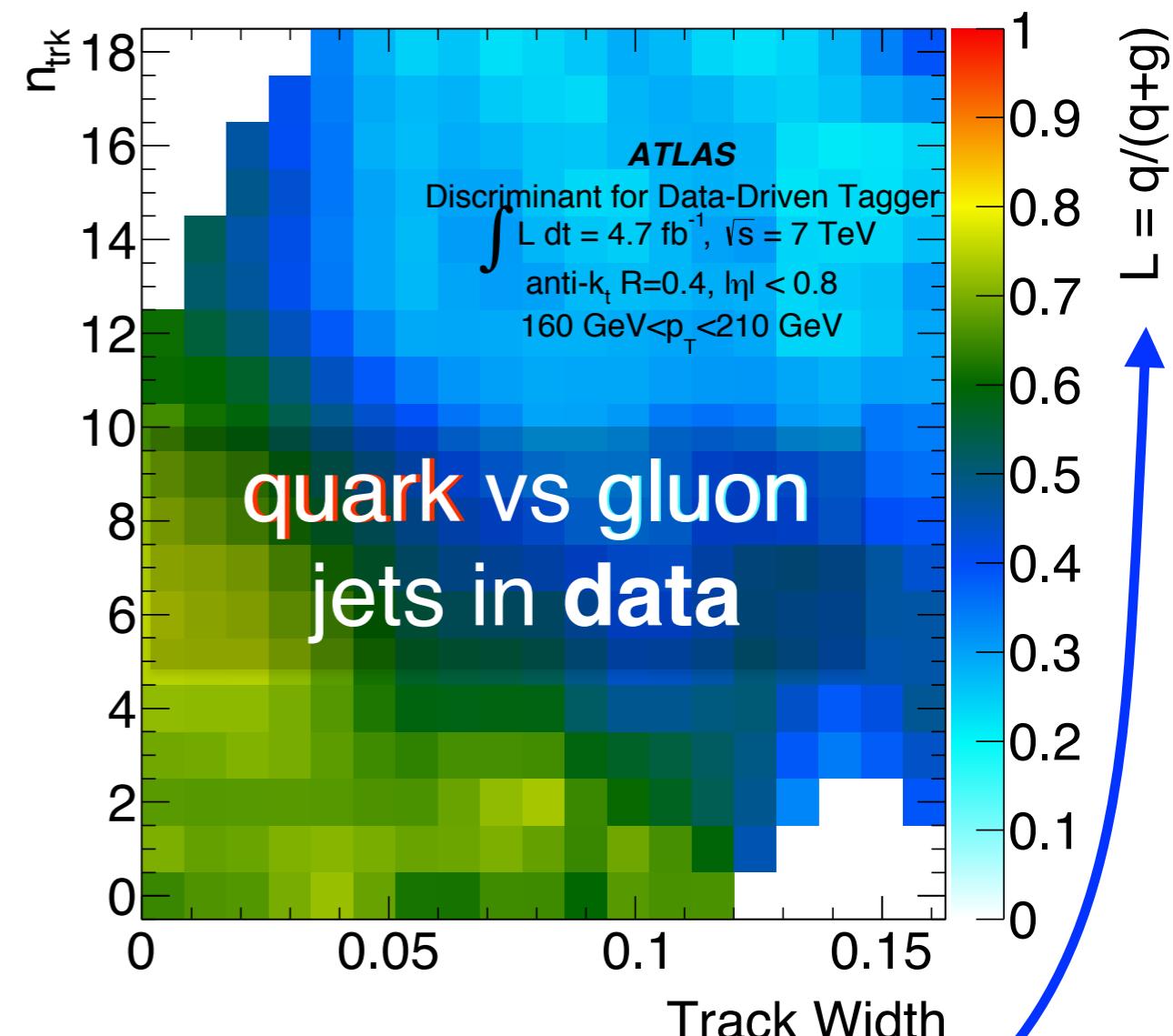
$$Z+\text{jets} = g_q \times Q + (1-g_q) \times G$$

2 equations, 2 unknowns (Q, G)

(now Q and G are 2D histograms)

in general:

$$h_{\text{MC}}(n_{\text{trk}}, \text{Track Width}) \neq h_{\text{data}}(n_{\text{trk}}, \text{Track Width})$$

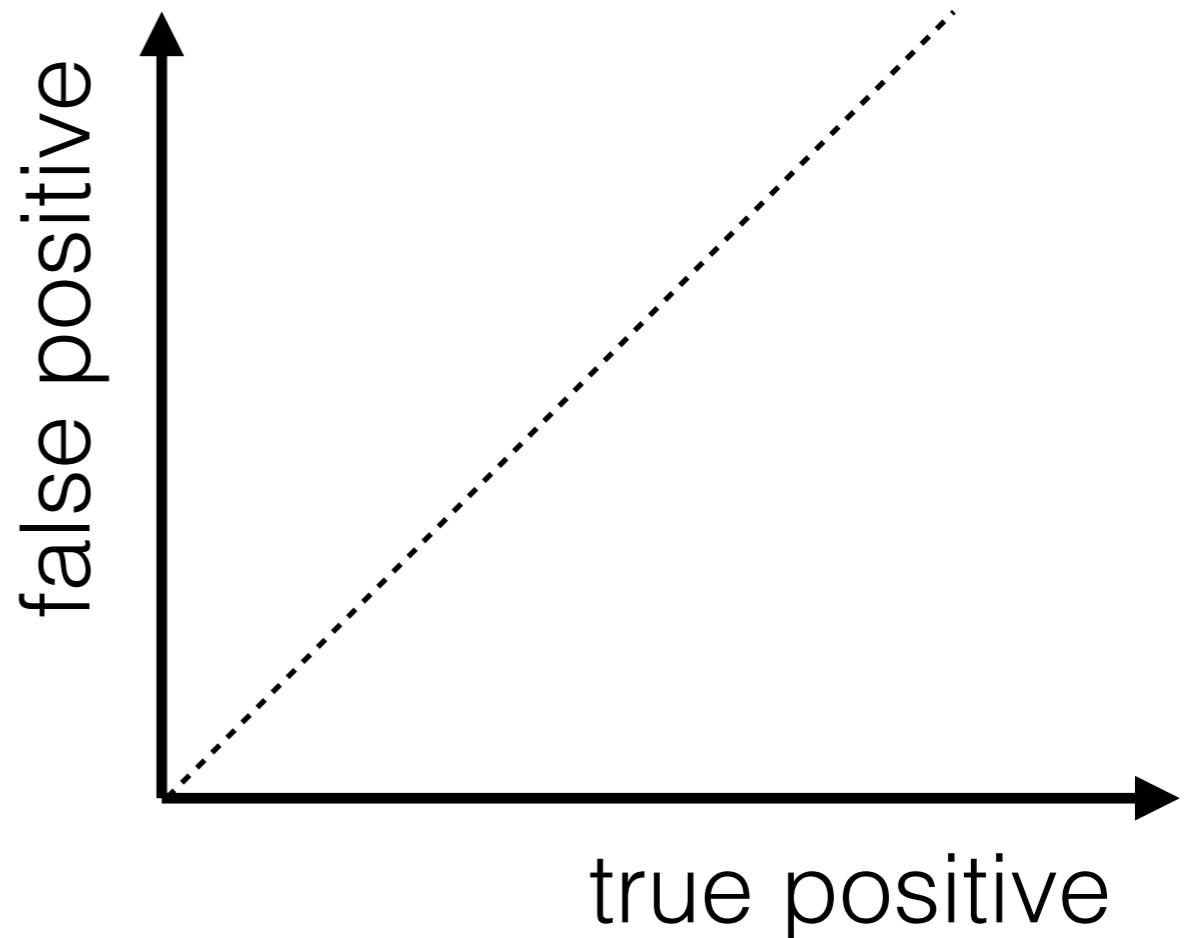


Take it to the extreme

To stress this point, suppose that h_{MC} is the random classifier:

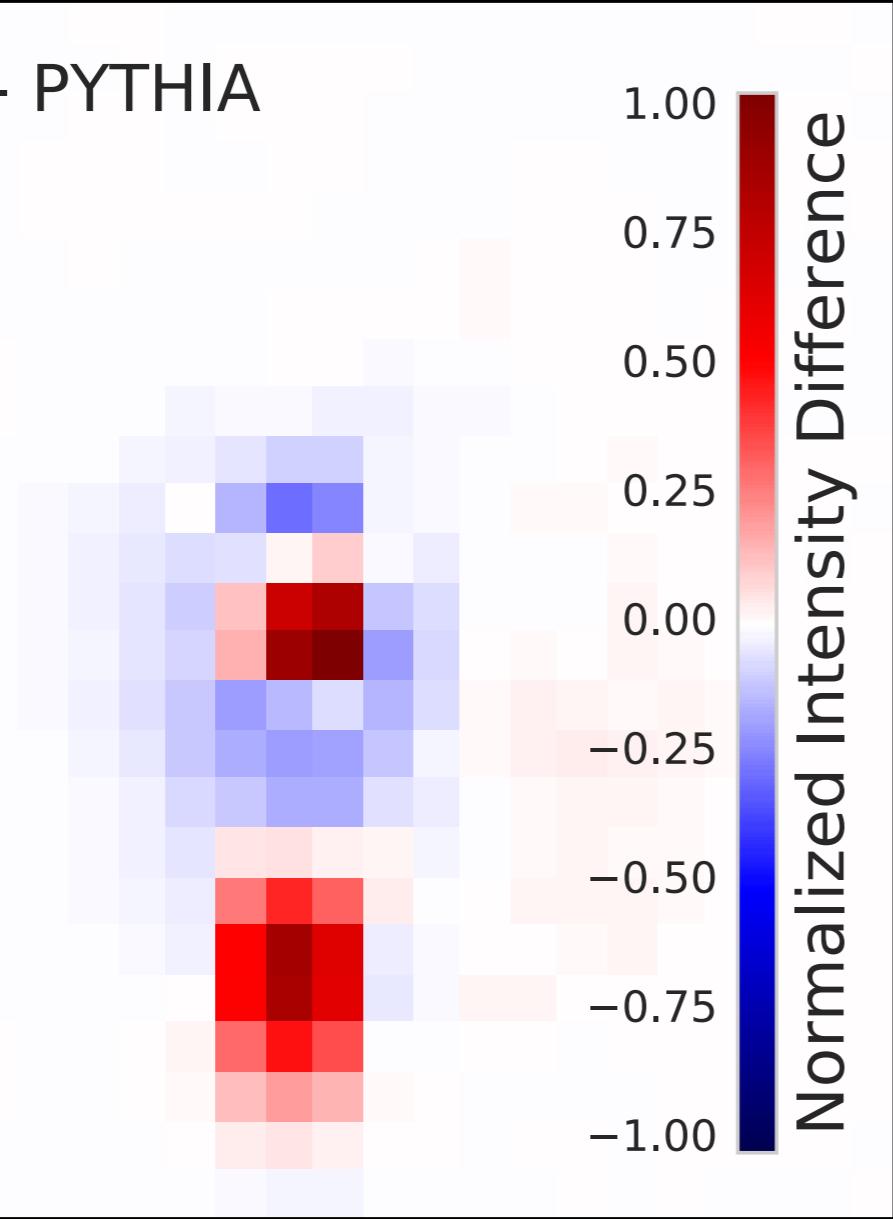
$h_{MC} = 0$ if you pick a random number x in $[0, 1]$ and $x < \varepsilon$
1 otherwise

We can calibrate this classifier in data, but clearly, it is sub-optimal !!



One more slide about why it matters

Sherpa - PYTHIA



Especially important for
deep learning using subtle
features → hard to model!

*W boson radiation
pattern - same physics,
different simulators!*

Achieving the Optimal Classifier

Two ways around the problems mentioned earlier:

- (1) Derive the classifier in MC, but don't let it use information that is not present in data.

“Learning to pivot”

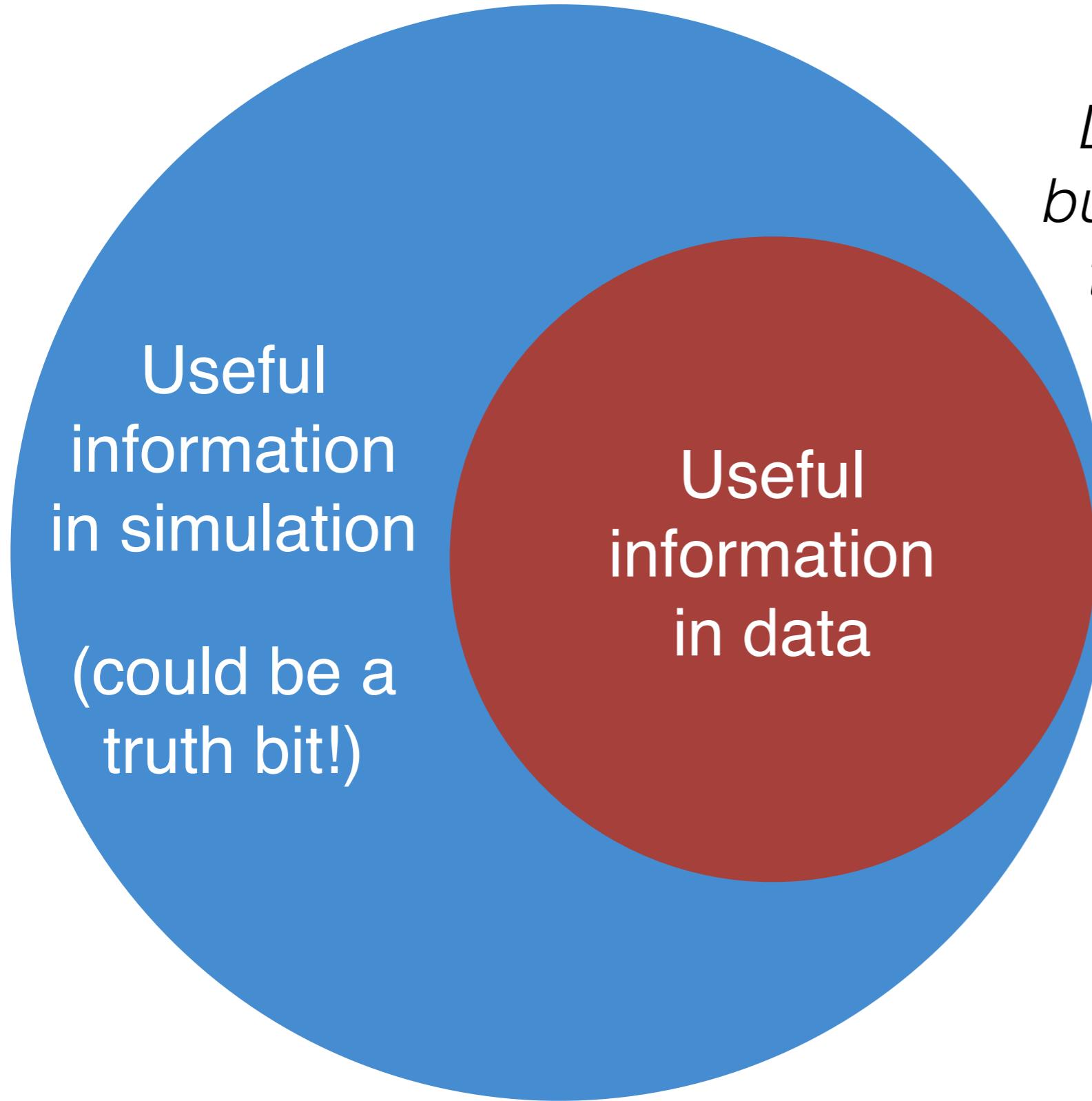
G. Louppe, M. Kagan, K. Cranmer, 1611.01046

- (2) Train on unlabeled data.

“Weak supervision”

*L. Dery, BPN, F. Rubbo, A. Schwartzman, JHEP 05 (2017) 145
E. Metodiev, BPN, J. Thaler, JHEP 10 (2017) 174*

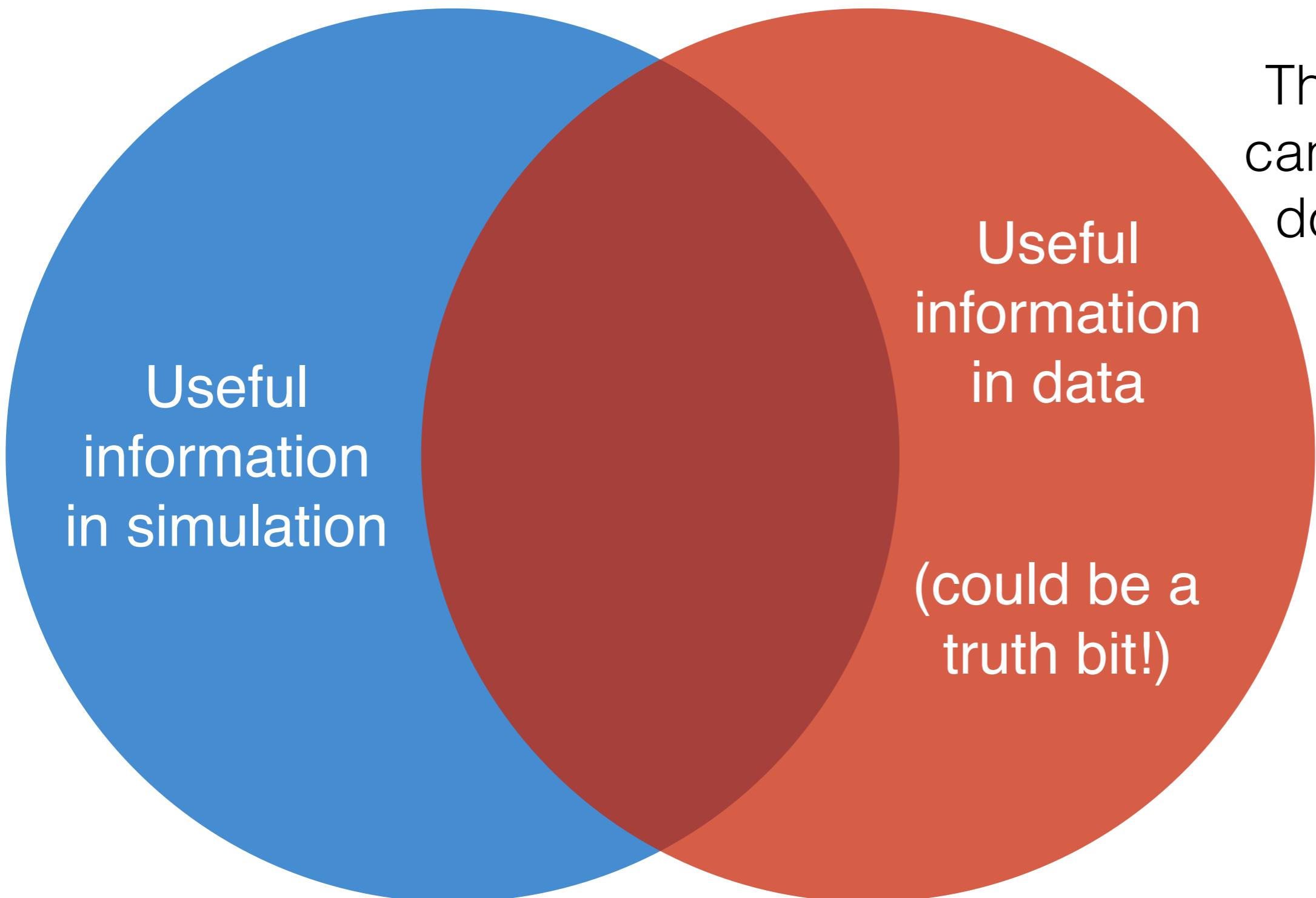
When option 1 works



Option 1:
*Derive the classifier in MC,
but don't let it use information
that is not present in data.*

Classifier can't use
information from
simulation that is not
useful in data and
can learn the (data)
optimal classifier.

When option 1 is sub-optimal

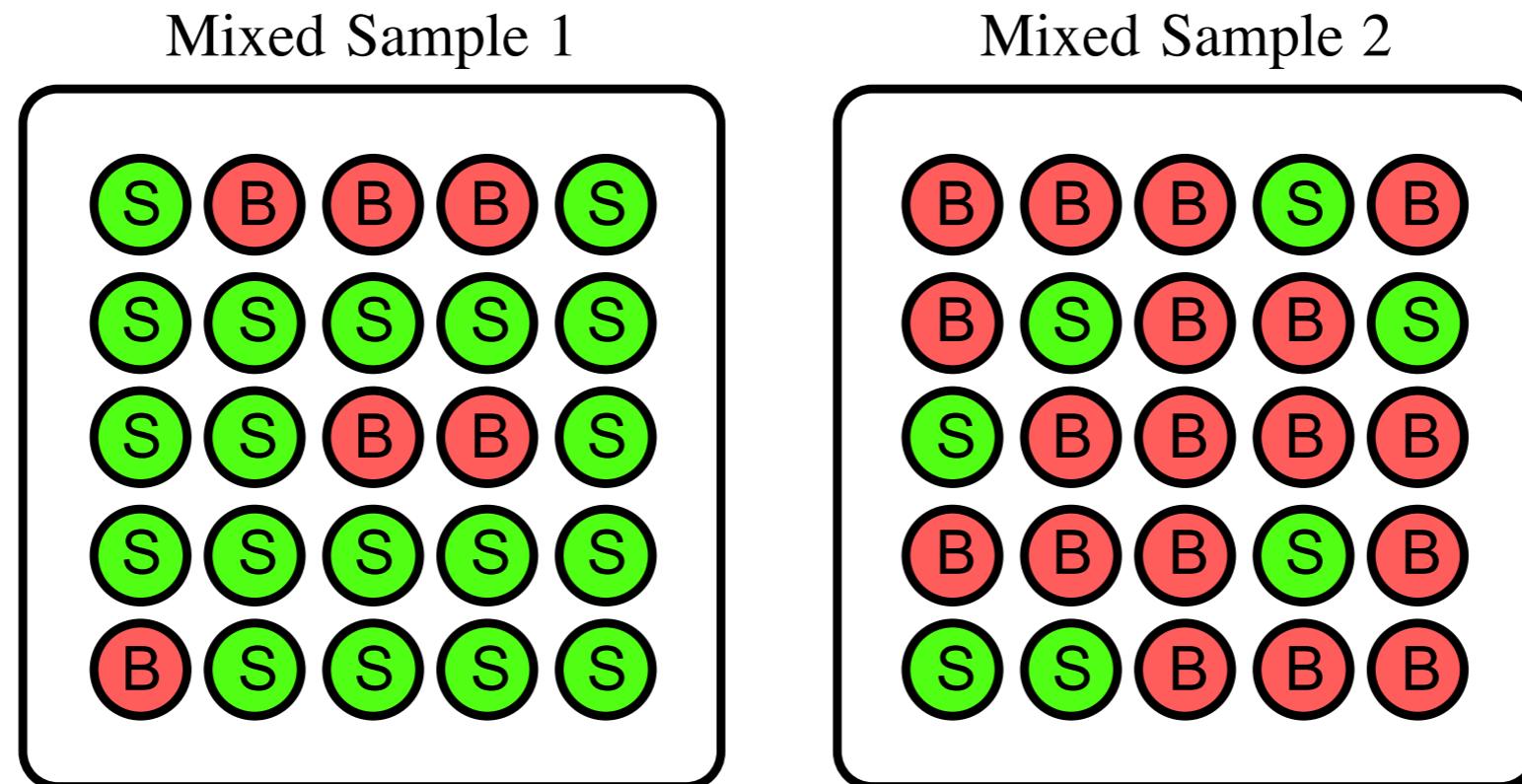


...many other applications of this approach, such as reducing sensitivity to systematic uncertainties, unwanted correlations between features, etc.

Option 2: Learn from data!

Challenge is that data are **unlabeled**.

The setup: suppose you have (at least) two mixed samples, each composed of two classes (say **q** and **g**).



E. Metodiev, BPN, J. Thaler, JHEP 10 (2017) 51

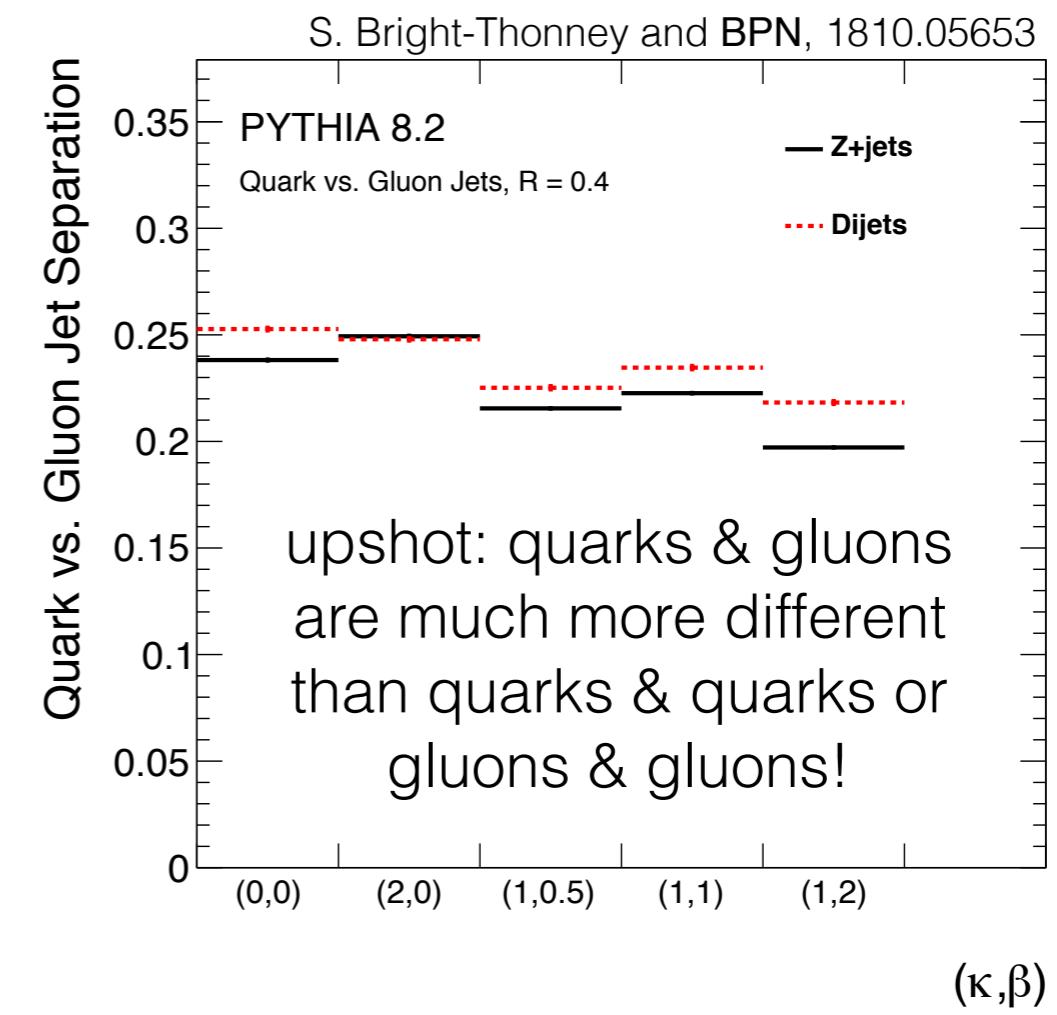
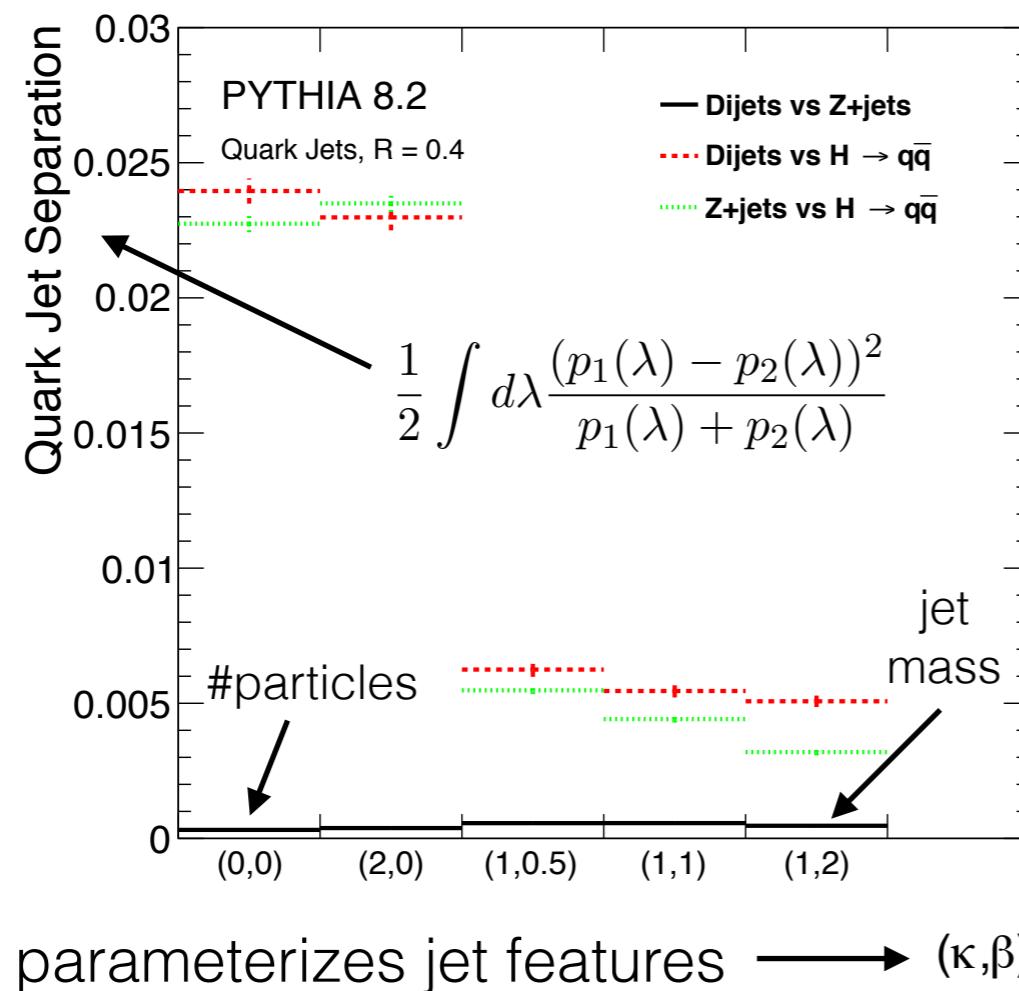
Assumption: The two classes are well-defined i.e. **q** in sample 1 is statistically identical to **q** in sample 2.

There is an interesting connection between what I'm calling “weak supervision” and the topic of “label noise”.

Weak supervision, caveats up front

The two classes are well-defined (i.e. q in sample 1 is statistically identical to q in sample 2).

This is often not exactly true, but is often nearly true.



Weak sup. option 1: Use class proportions

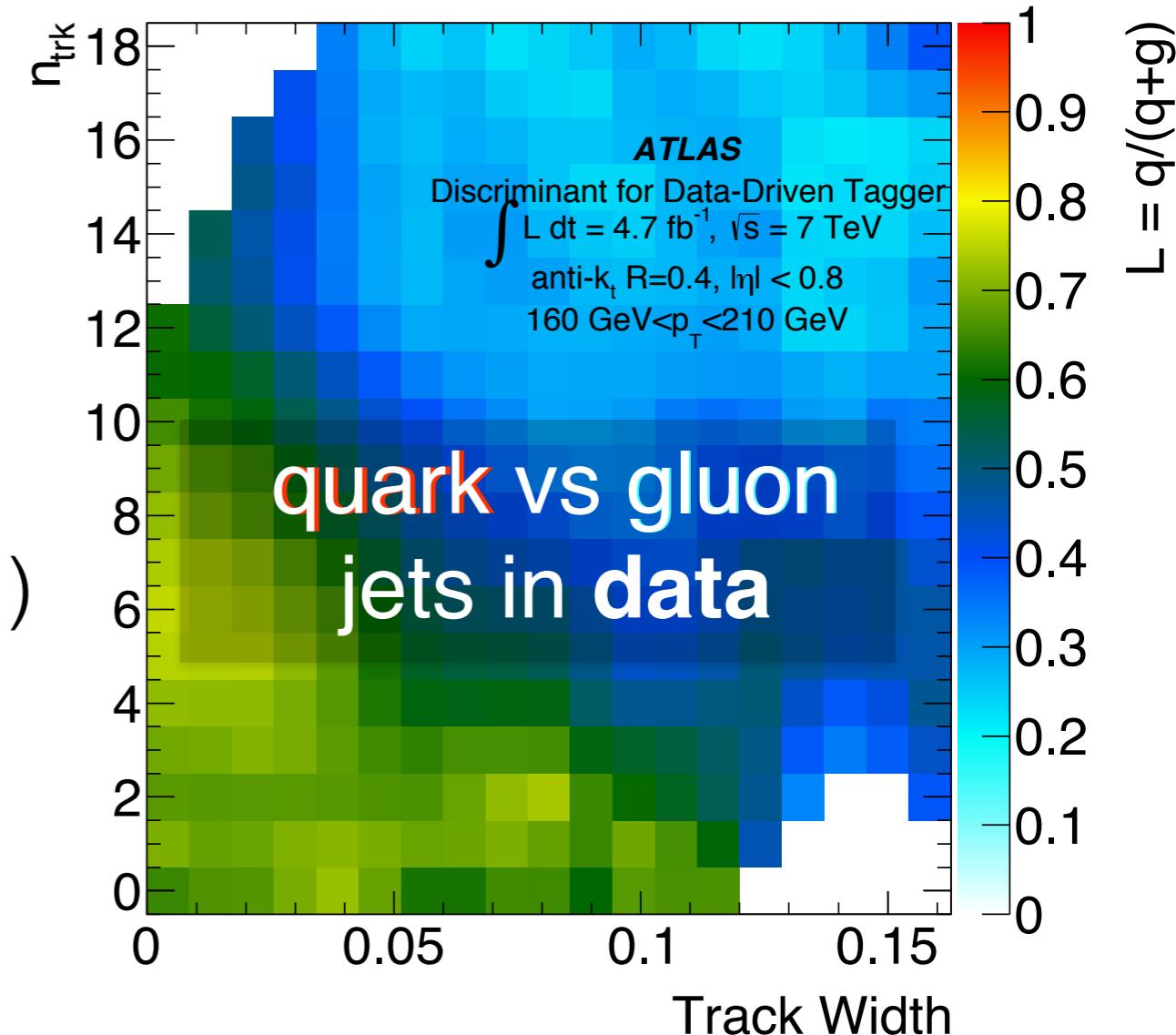
Remember this plot?

$$\text{dijets} = f_q \times Q + (1-f_q) \times G$$

$$Z+\text{jets} = g_q \times Q + (1-g_q) \times G$$

two equations, two unknowns (Q, G)

We often know f, g
 (from ME + PDF) much better than
 full radiation pattern inside jets.



This doesn't work well when you have more than 2 observables because the templates become sparse.

Method 1: Learn from Proportions



LoLiProp

*Learning from
Label Proportions*

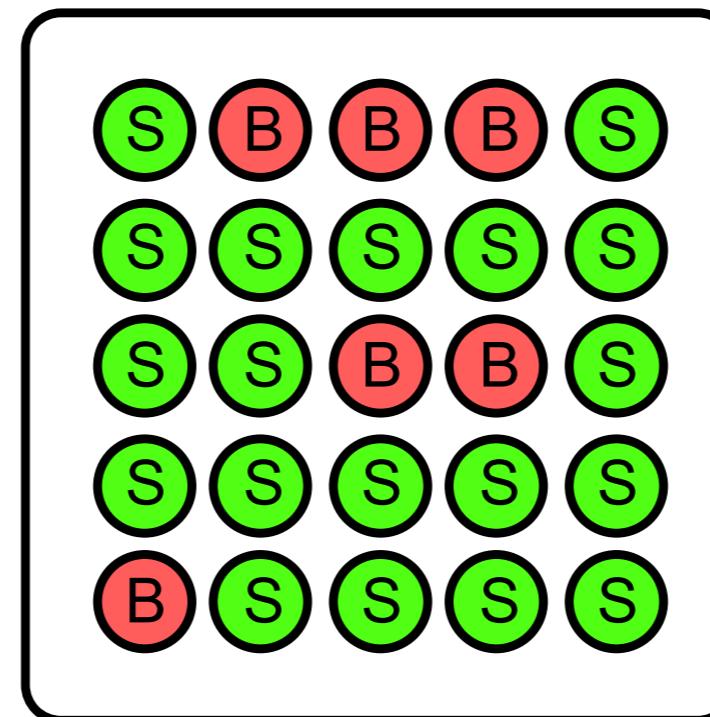
Solution: Train using
class proportions.

Work “on average”

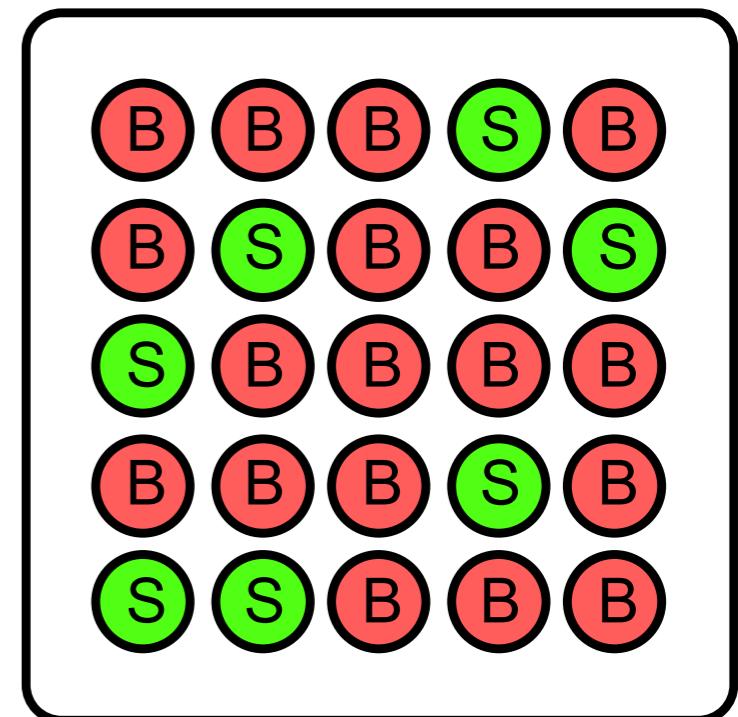
$$f_{\text{full}} = \operatorname{argmin}_{f': \mathbb{R}^n \rightarrow \{0,1\}} \sum_{i=1}^N \ell(f'(x_i) - t_i)$$

loss fcn. *labels*

Mixed Sample 1



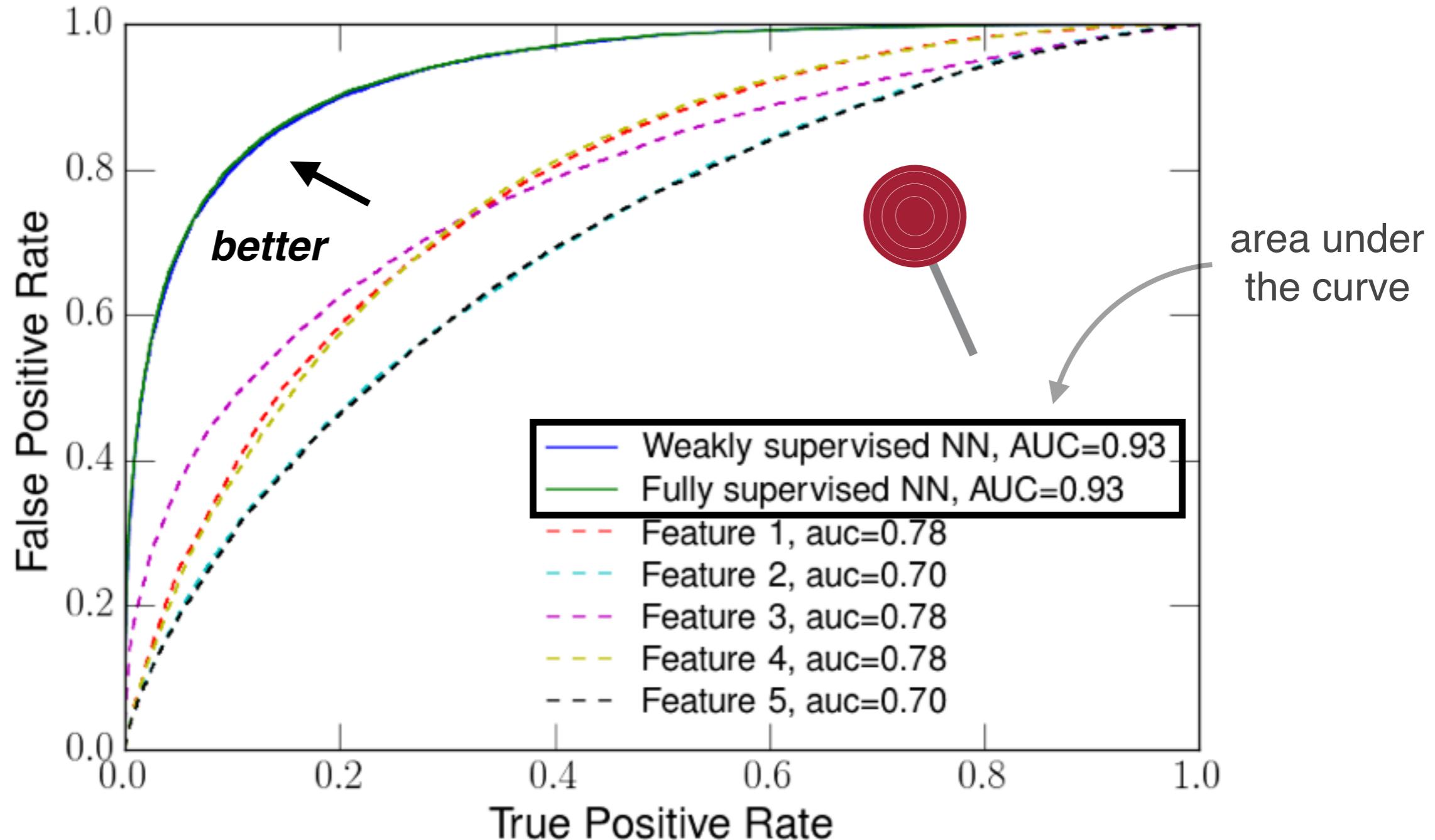
Mixed Sample 2



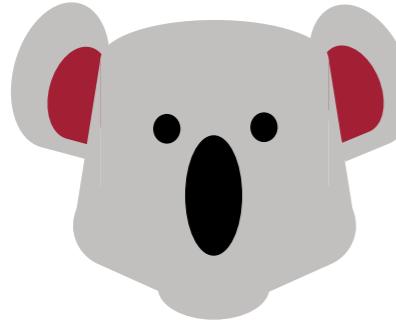
$$f_{\text{weak}} = \operatorname{argmin}_{f': \mathbb{R}^n \rightarrow [0,1]} \ell \left(\sum_{i=1}^N \frac{f'(x_i)}{N} - y \right)$$

proportions

Works!

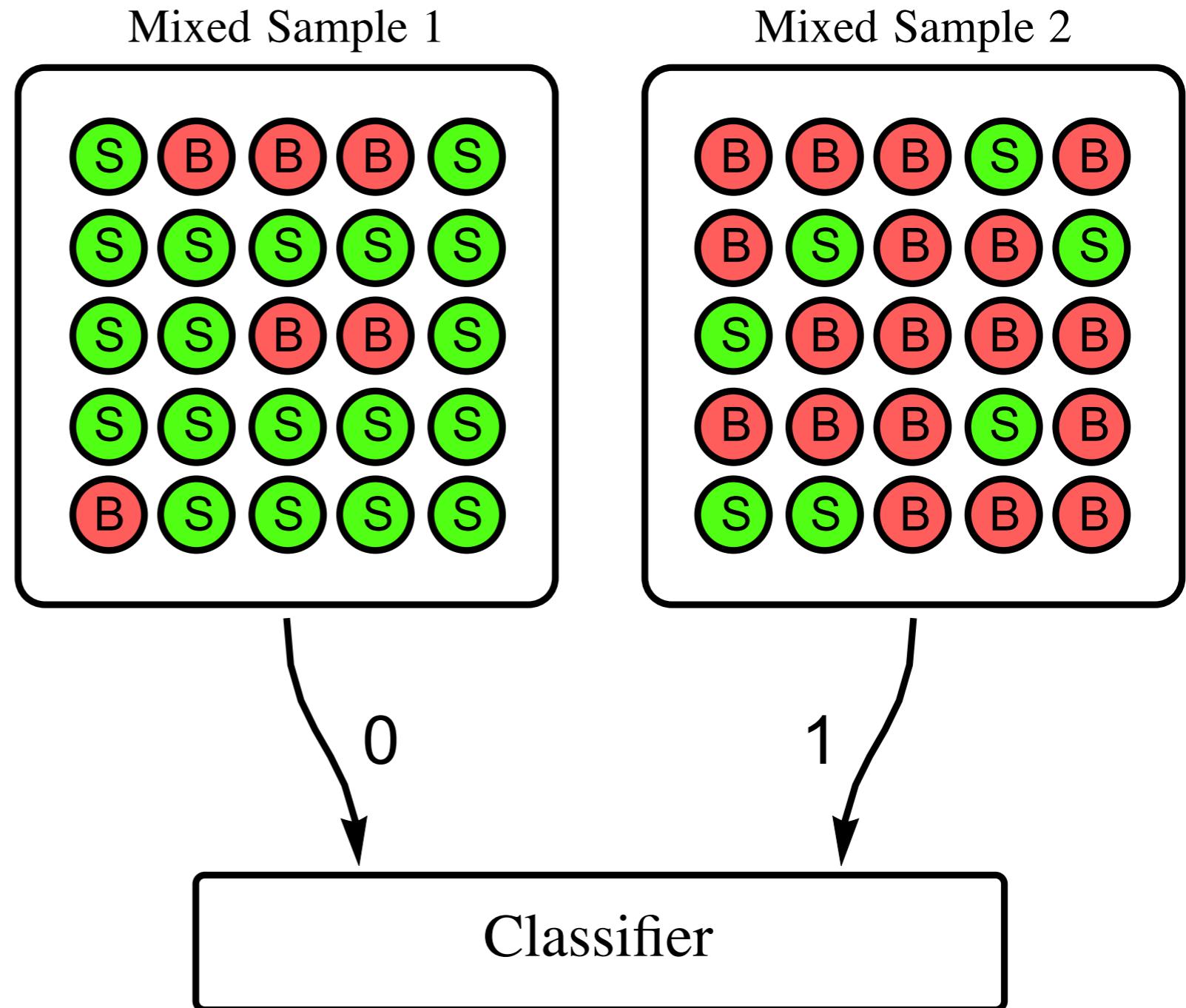


Method 2: Learning without Proportions



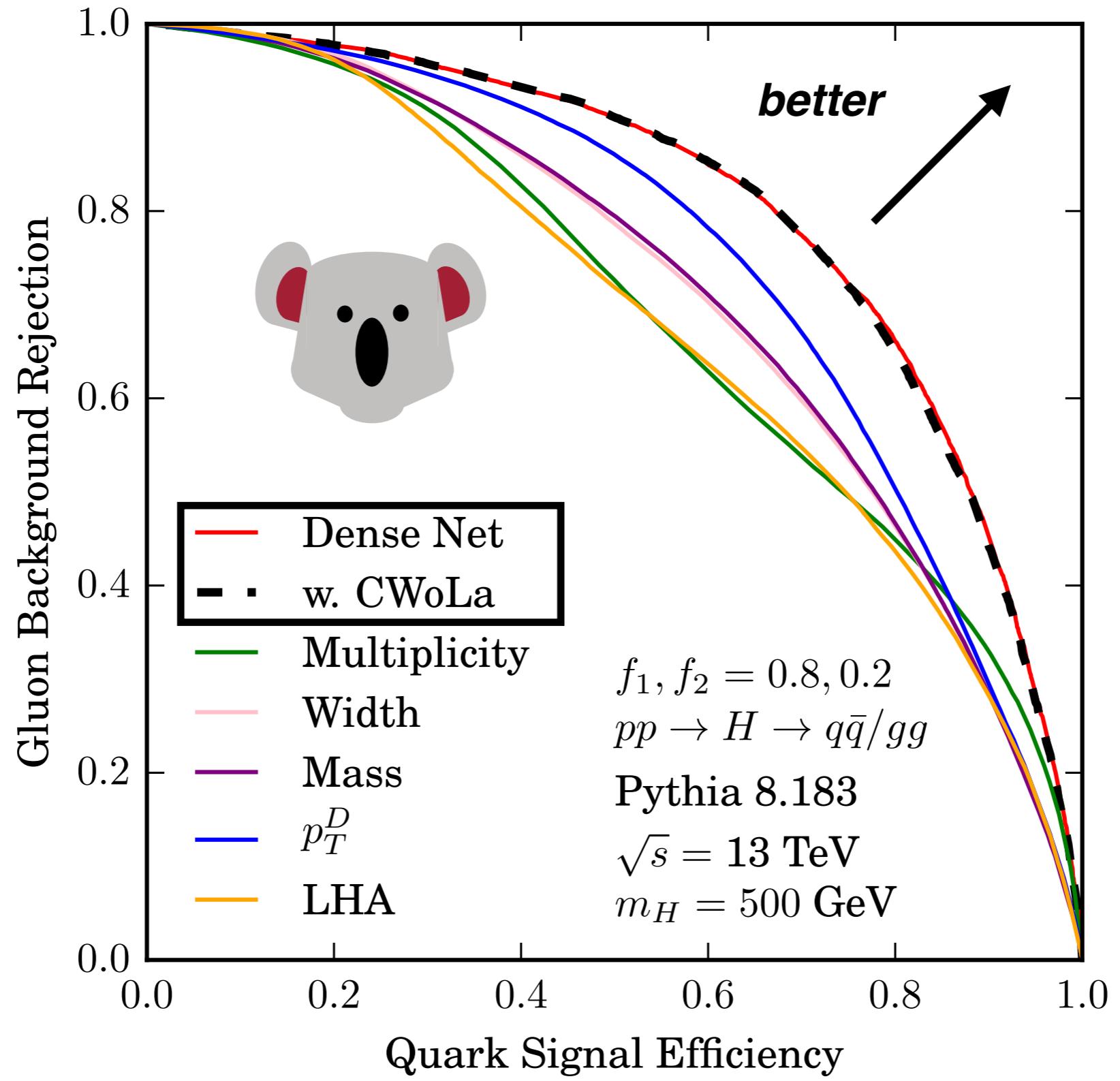
CWoLa
*Classification
Without Labels*

Solution: Train
directly on data using
mixed samples

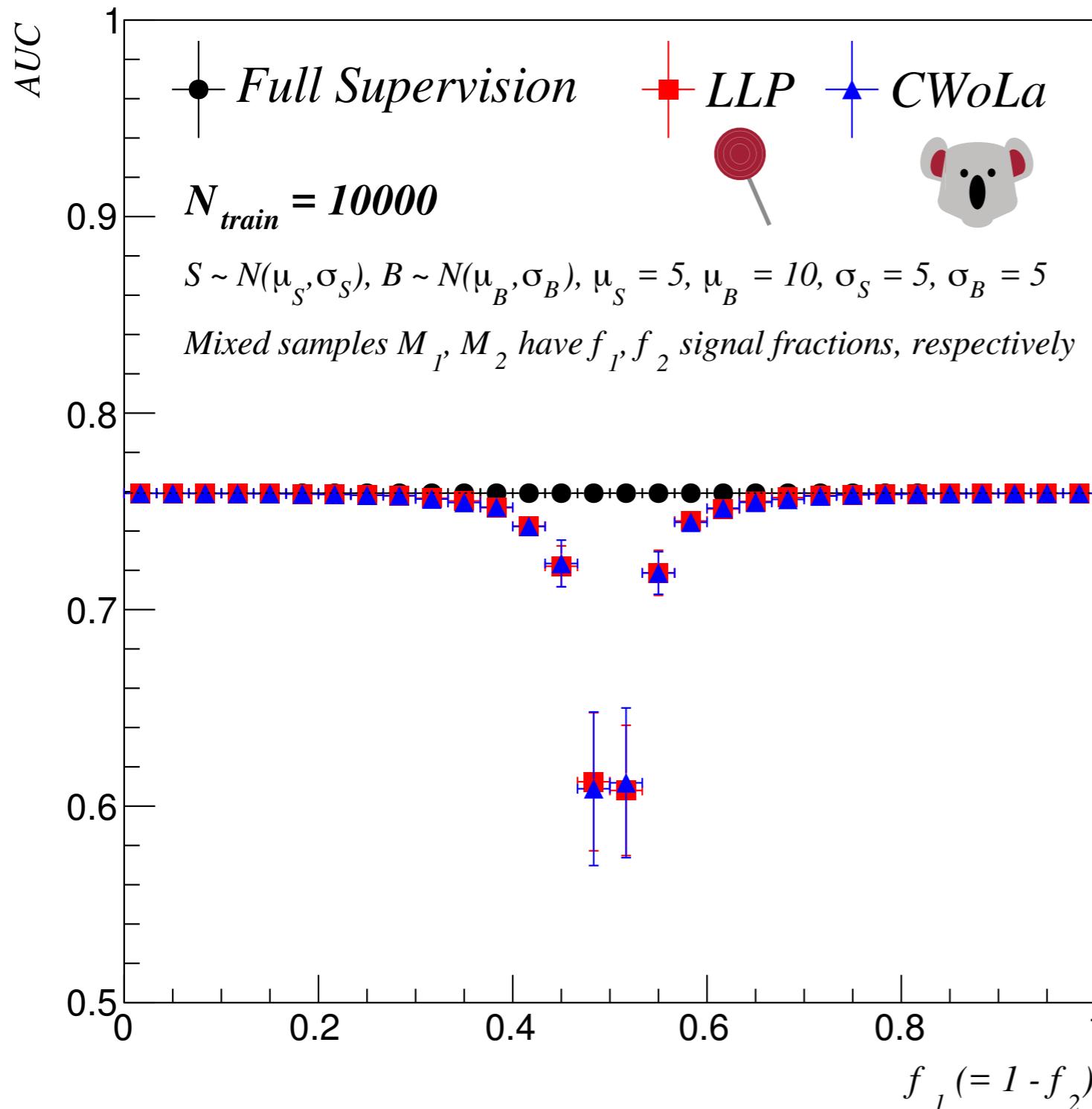


Works!

30



A note about training statistics



signal fraction of mixed sample 1

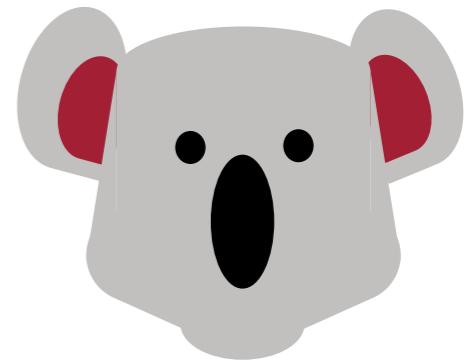
Can't learn when
the two proportions
are the same.

The more similar
they are, the worse
the performance.
(lower effective
statistics)

CWoLa Hunting for BSM

- Why model agnostic searches?

Intermezzo: Learning directly from data



- Why it is important to learn directly from data
 - How to learn without labels
- BSM Searches** background / signal model independent

Most searches at the LHC
(with or without ML)

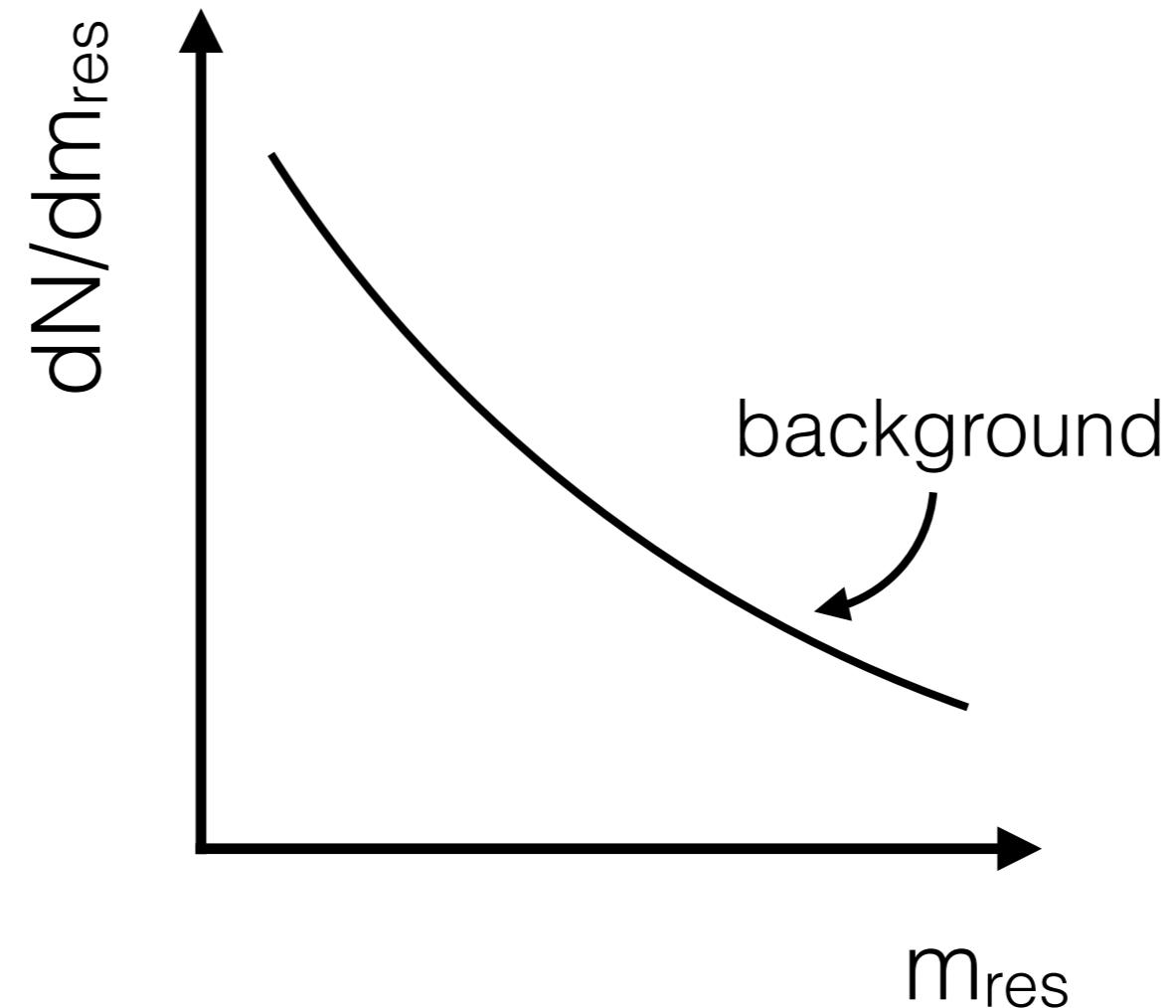
Data versus SM simulation
("general search")

CWoLa Hunting
(focus of this talk)

Autoencoders

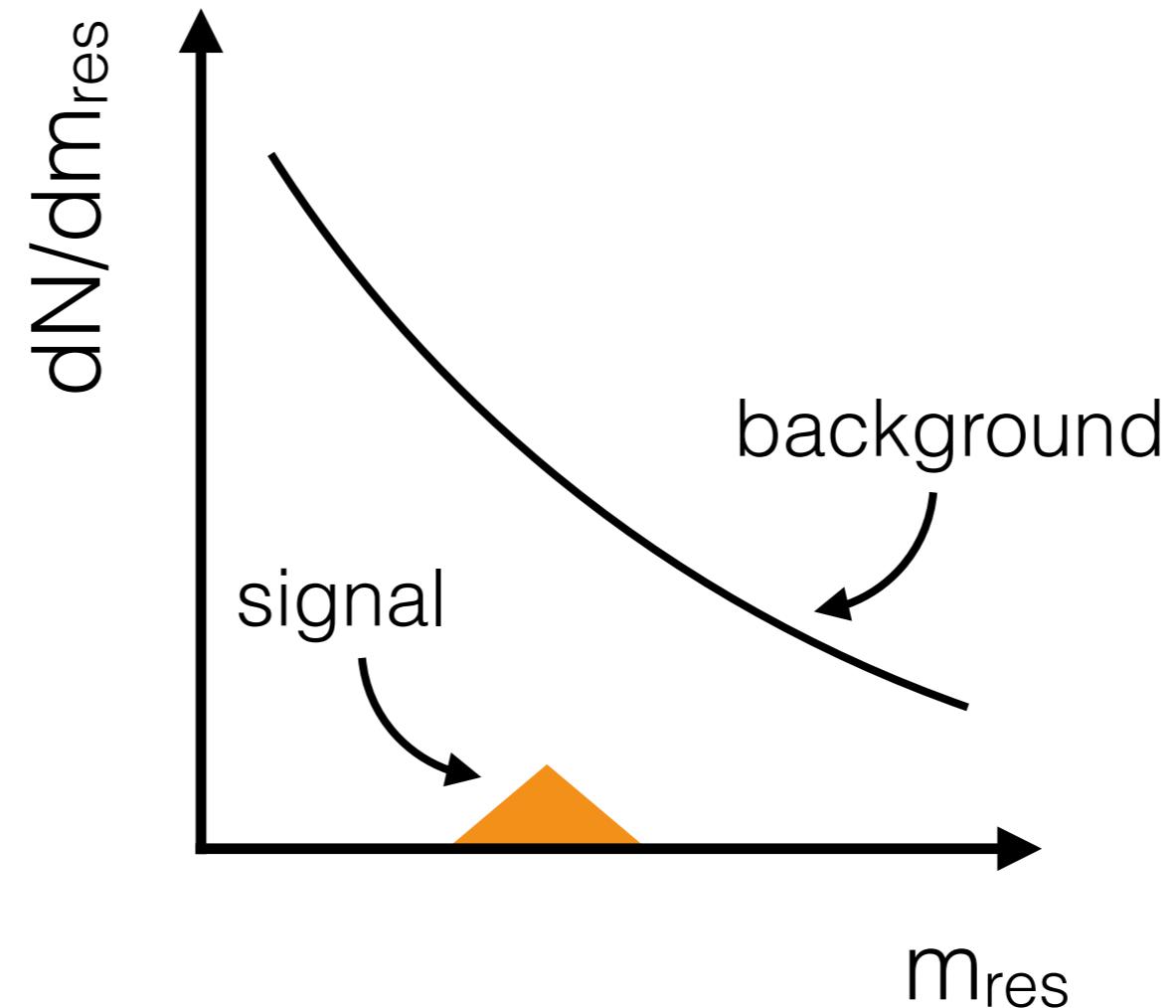
- The future

CWoLa Hunting for BSM



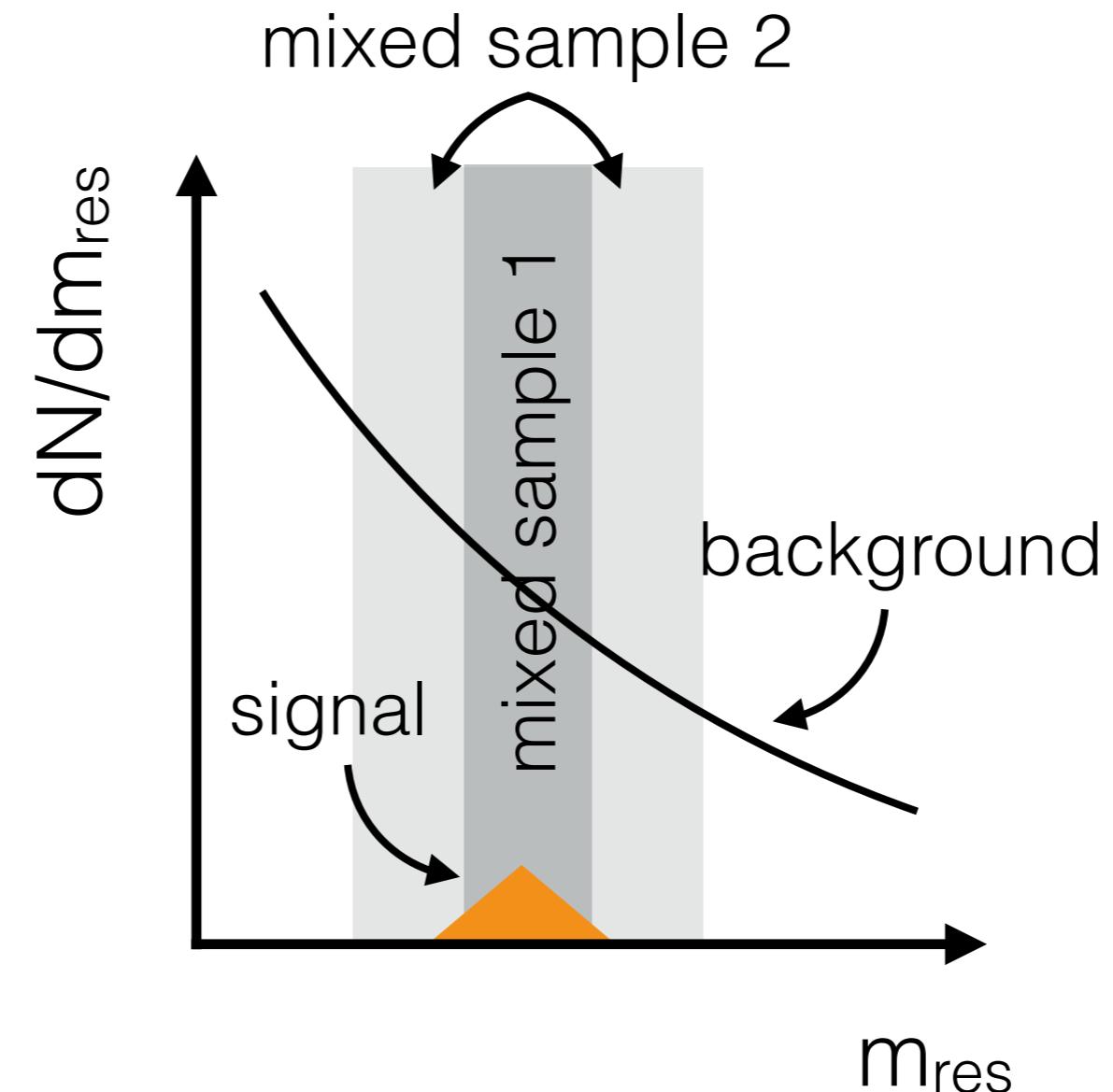
Assumption: there is a feature that we know about where the background is smooth and the signal (if it exists) is localized.

CWoLa Hunting for BSM



Assumption: there is a feature that we know about where the background is smooth **and the signal (if it exists) is localized**.

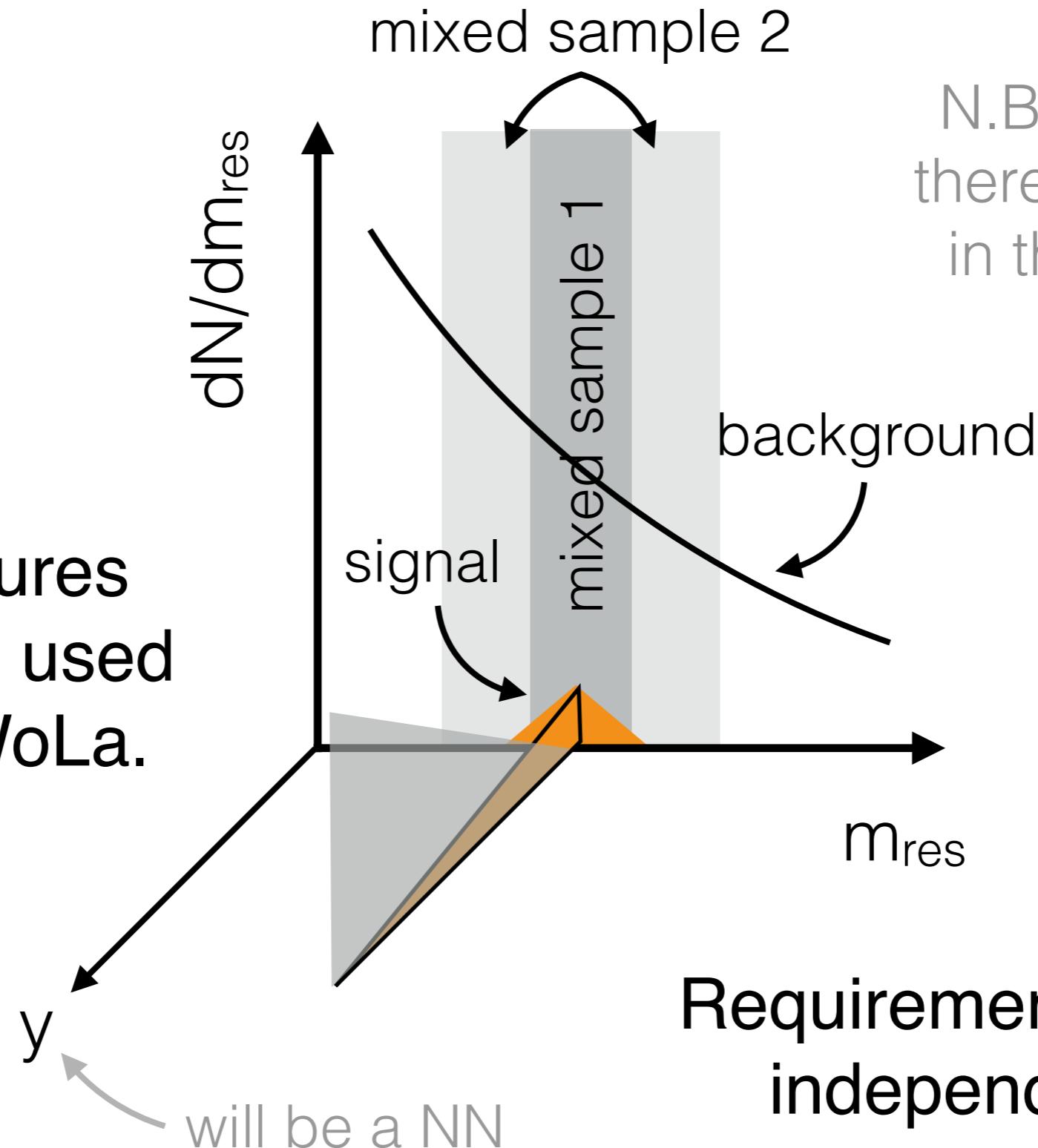
CWoLa Hunting for BSM



We don't know where the signal is, but for a given hypothesis, we can make signal windows and sidebands.

CWoLa Hunting for BSM

Other features
can then be used
to train CWoLa.



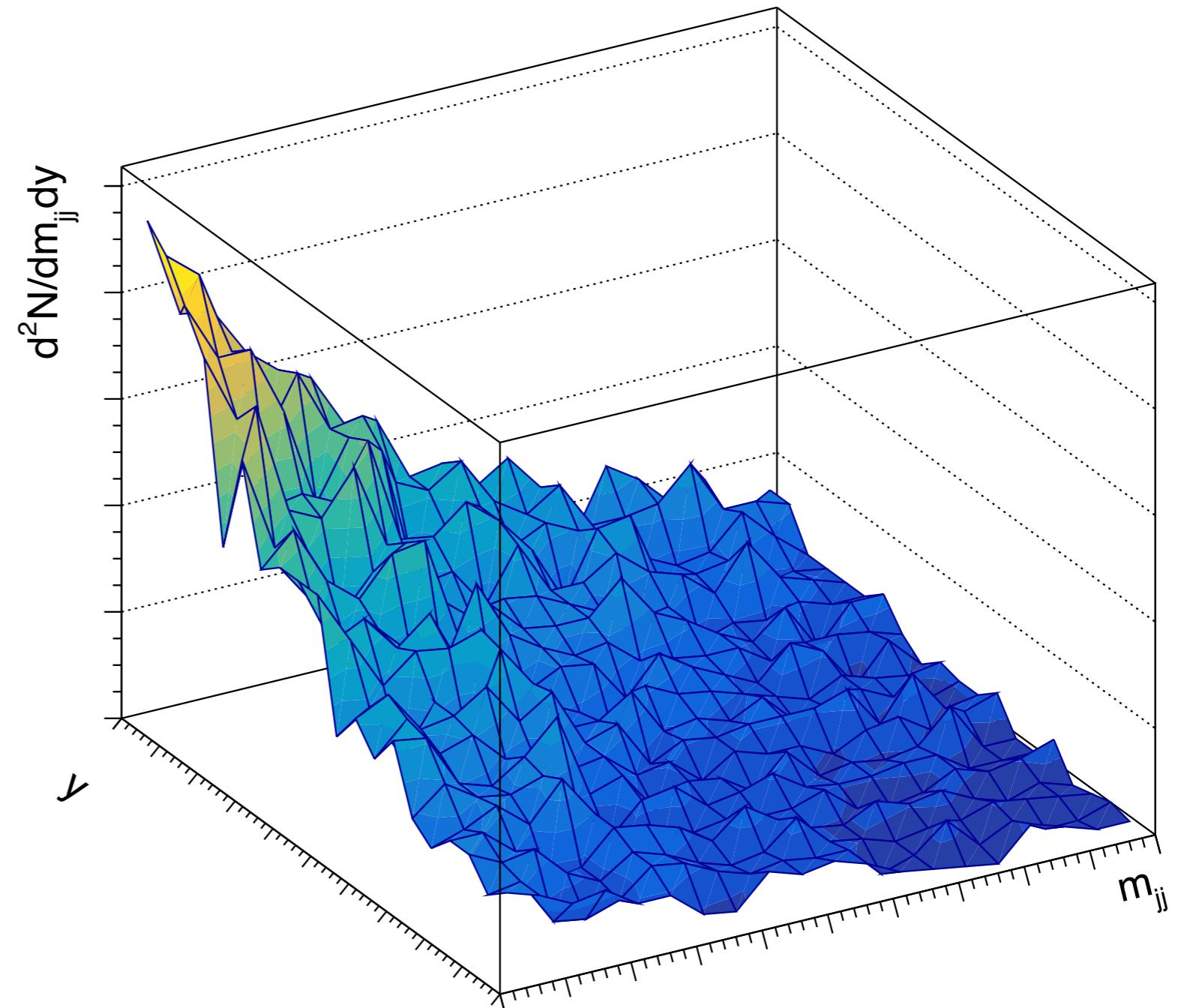
N.B. it is okay that
there is some signal
in the sidebands!

Requirement: y is (nearly)
independent of m_{res} .

Overtraining & Look Elsewhere Effect

Naively, pay a huge LEE penalty because y can be high-dimensional.

i.e. you will sculpt lots of bumps!

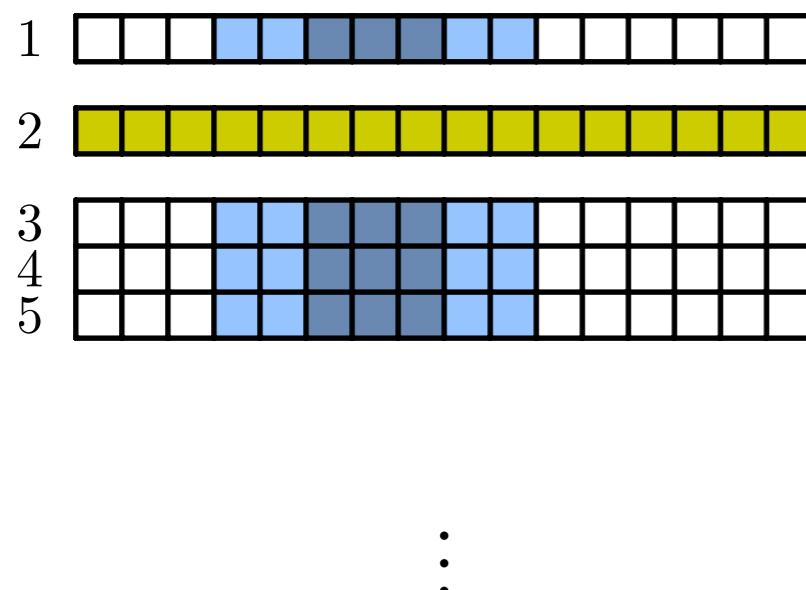
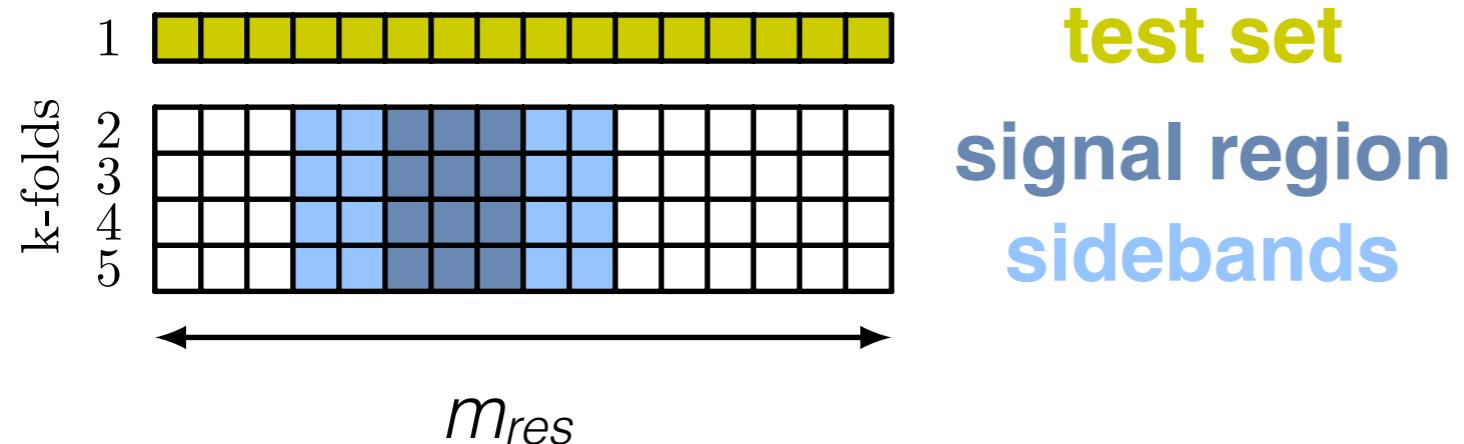


Solution: **(nested) cross-training**

Nested cross-training

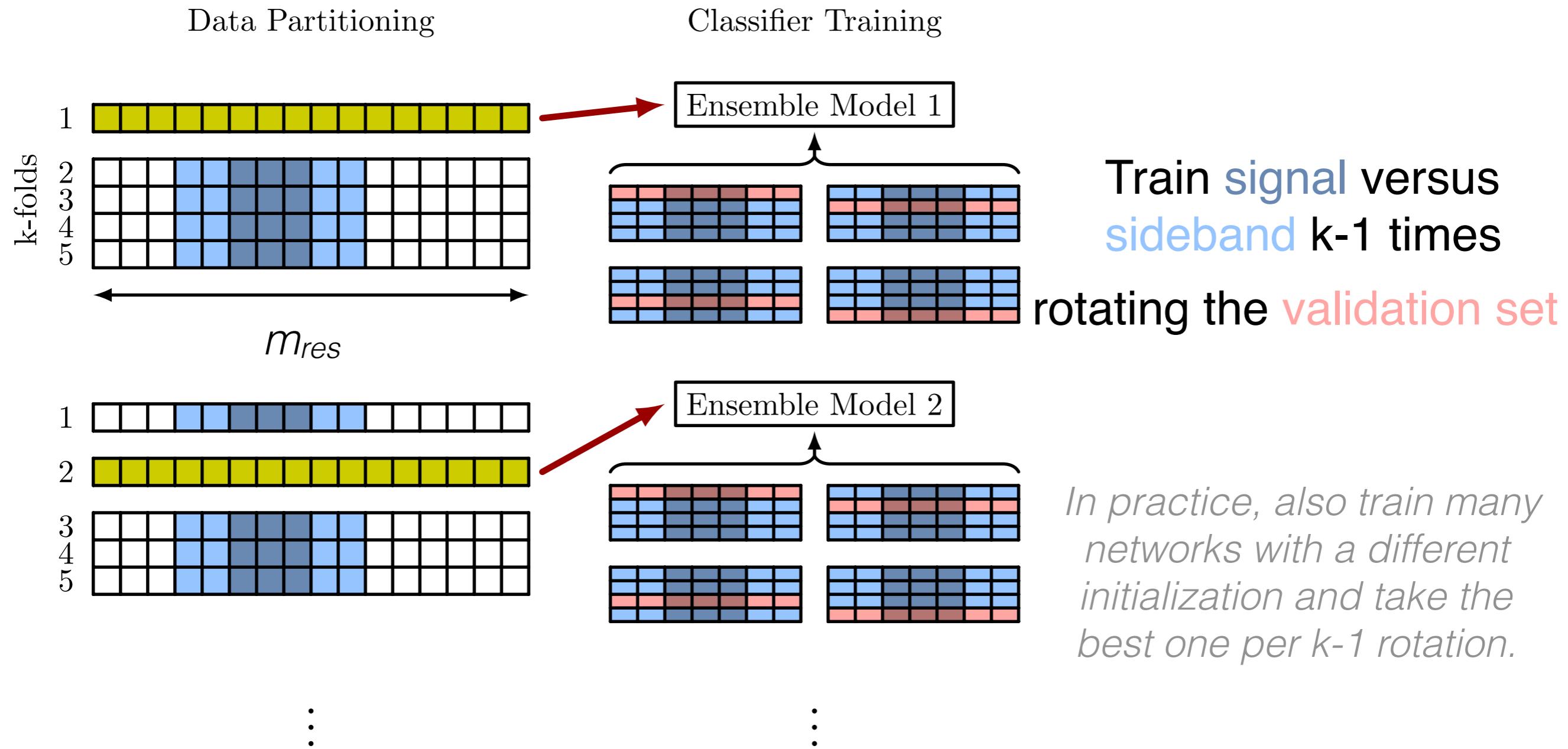
(1) Divide the entire dataset into k-folds.

Data Partitioning



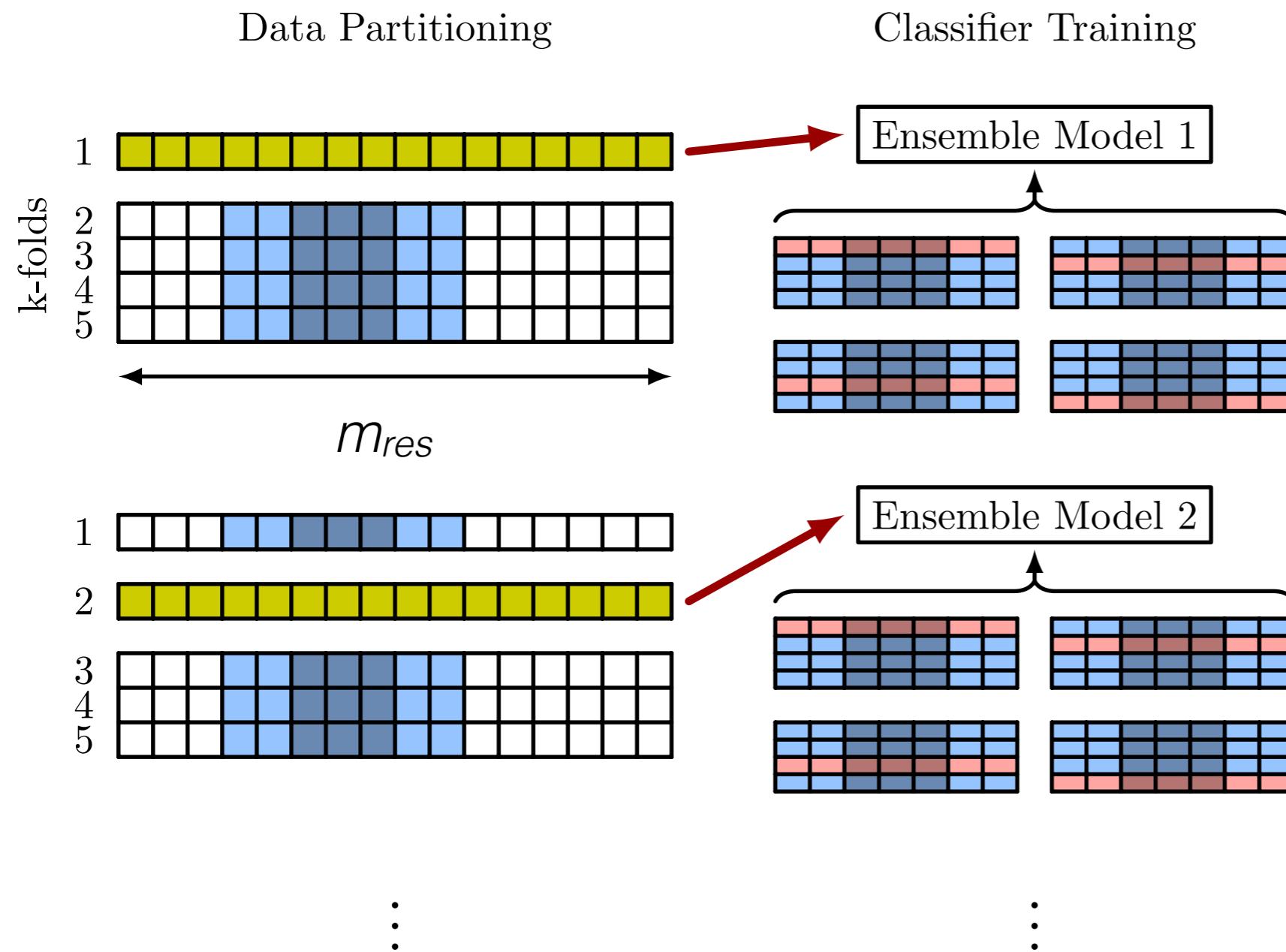
Nested cross-training

(2) Train CWoLa classifiers.



Nested cross-training

(2) Train CWoLa classifiers.

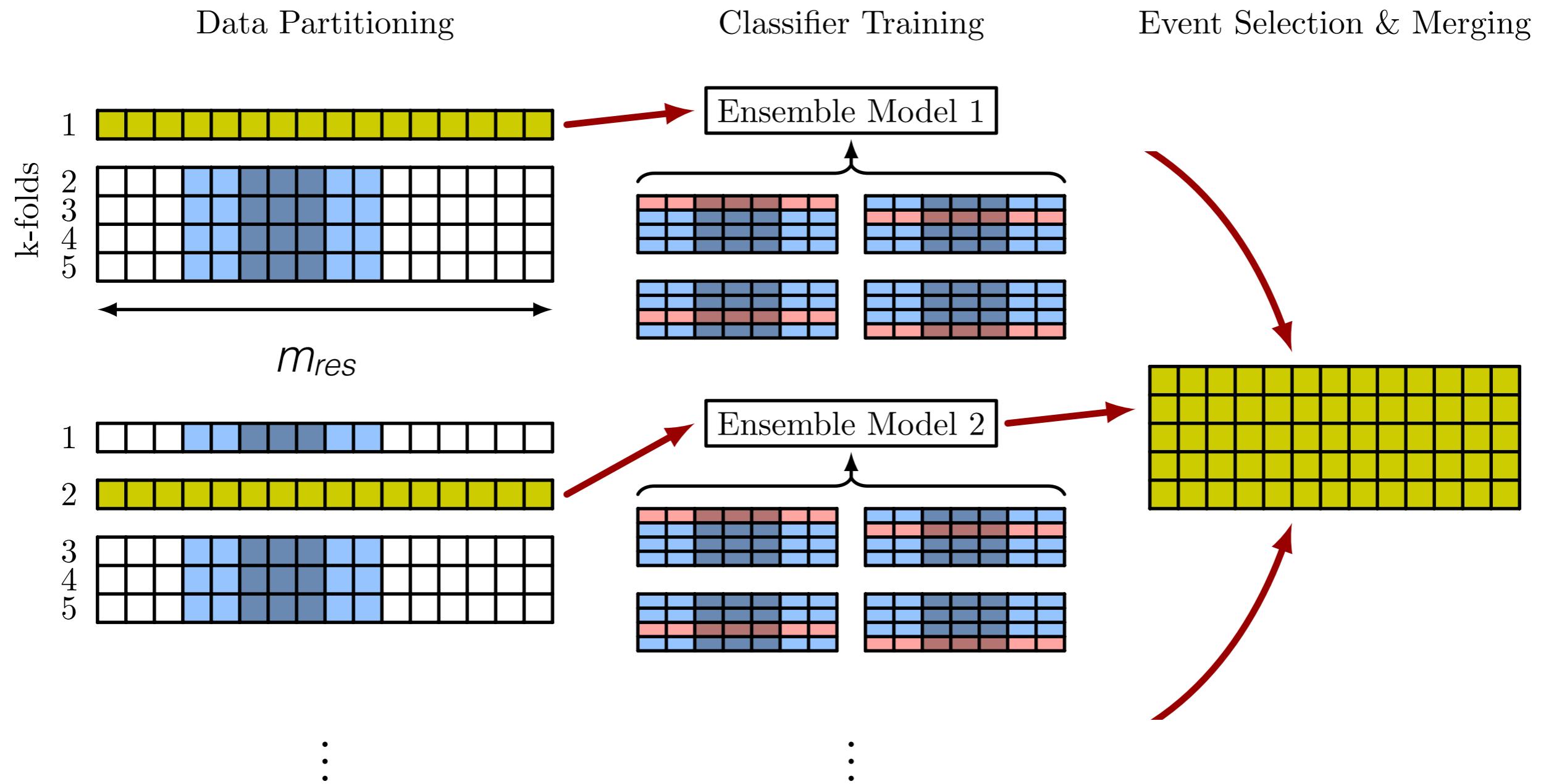


The **Ensemble Model** is just the average of the four networks.

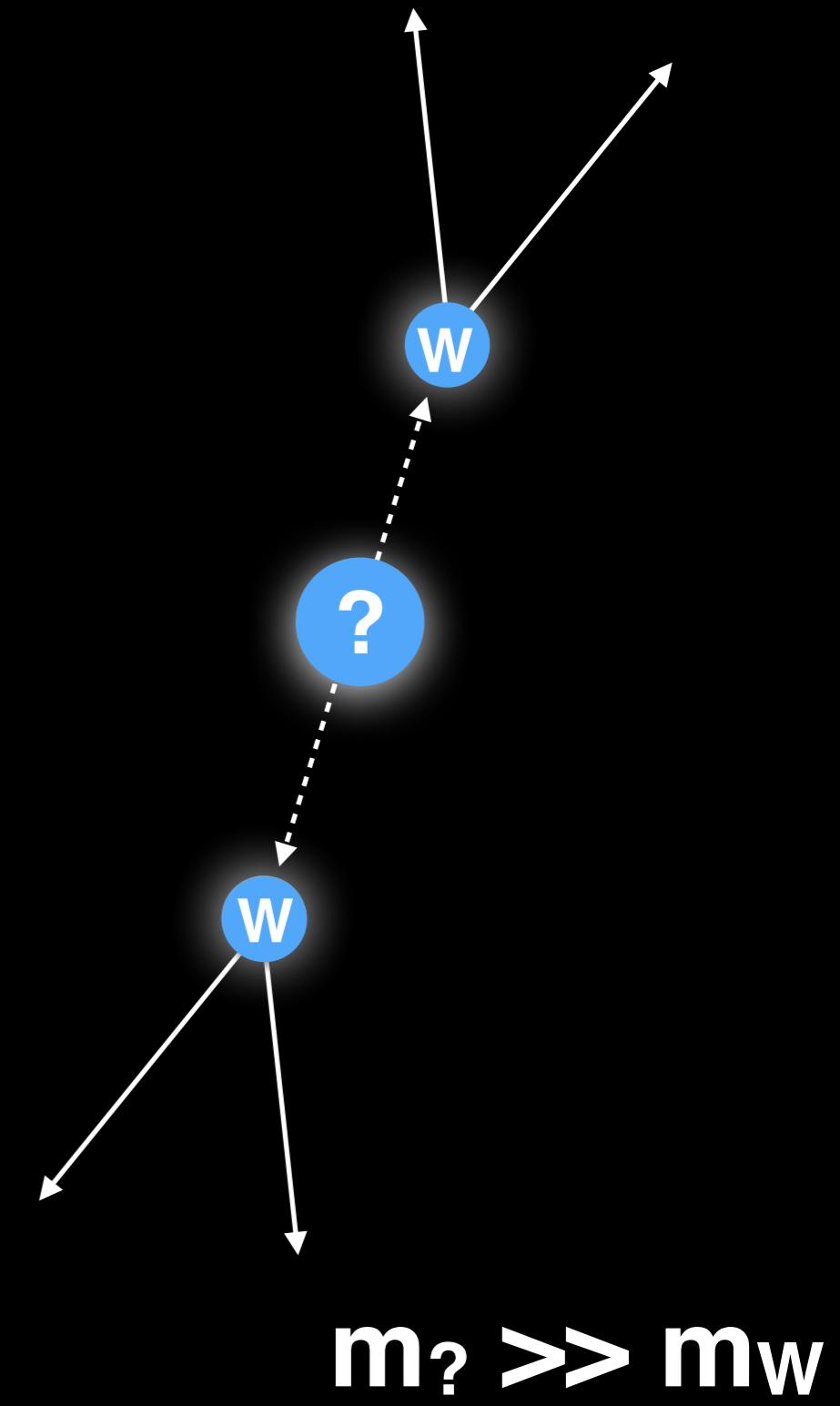
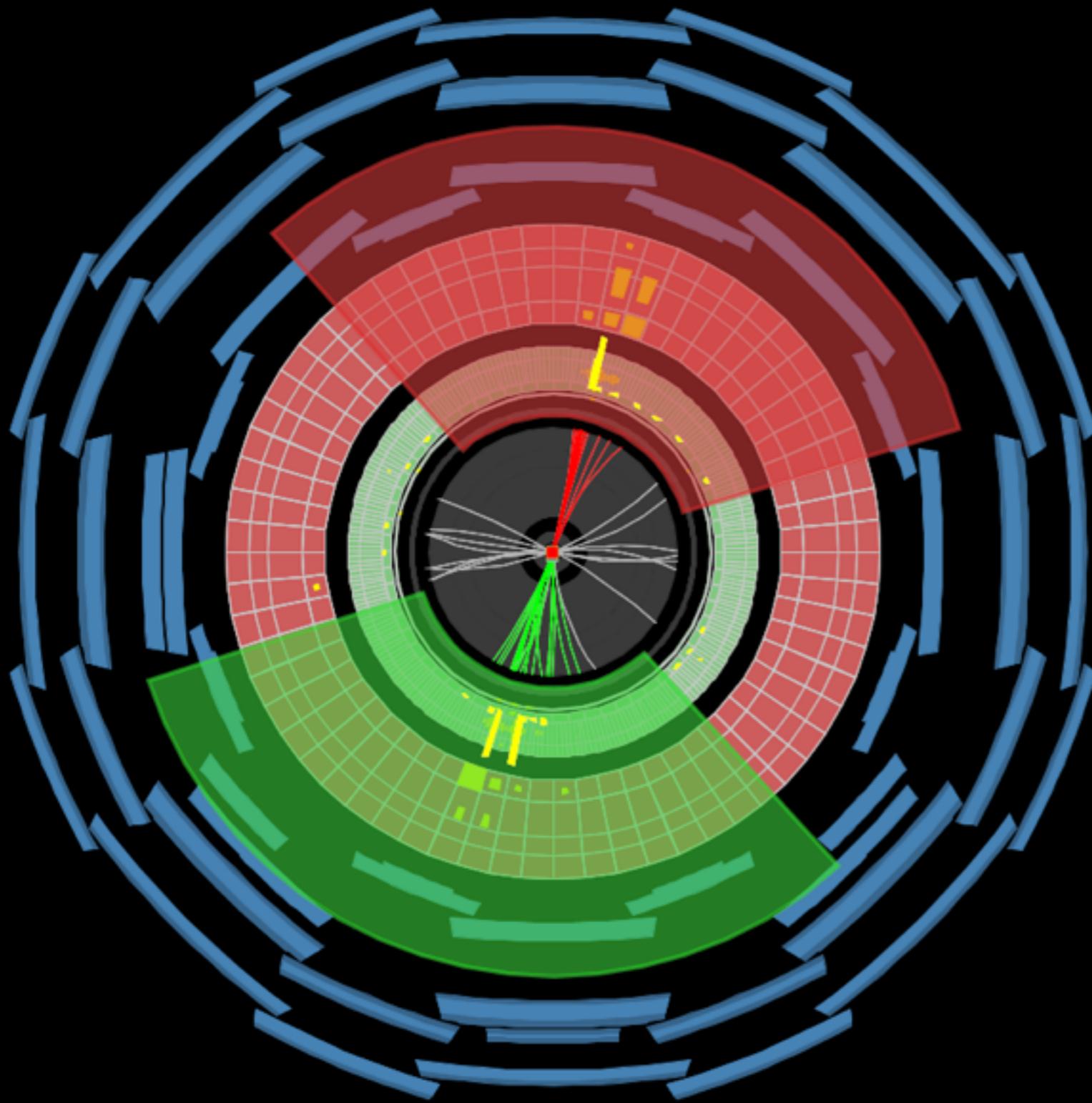
Data fluctuations will cancel destructively while signal interferes constructively.

Nested cross-training

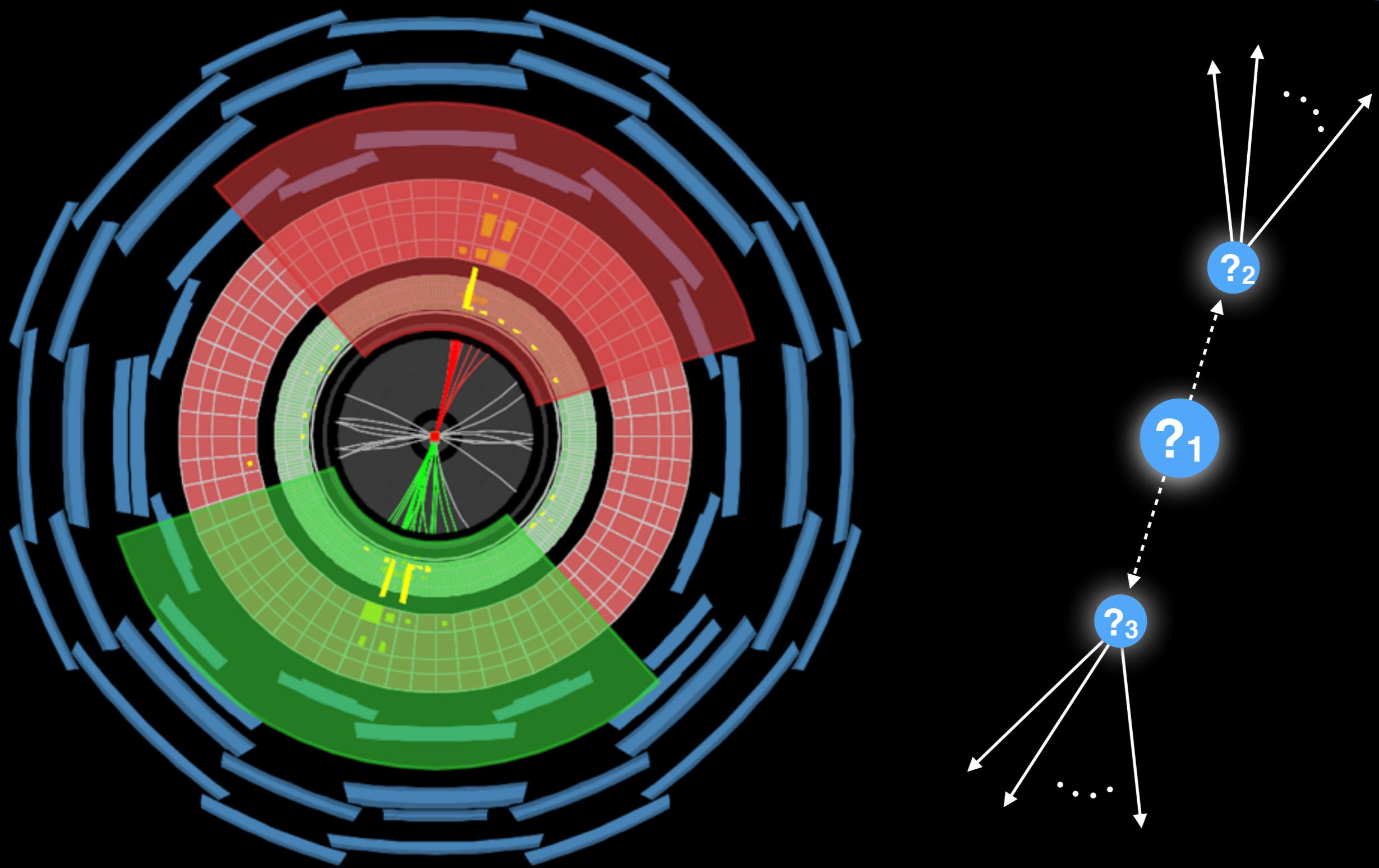
(3) Apply classifiers to holdout test sets and sum.



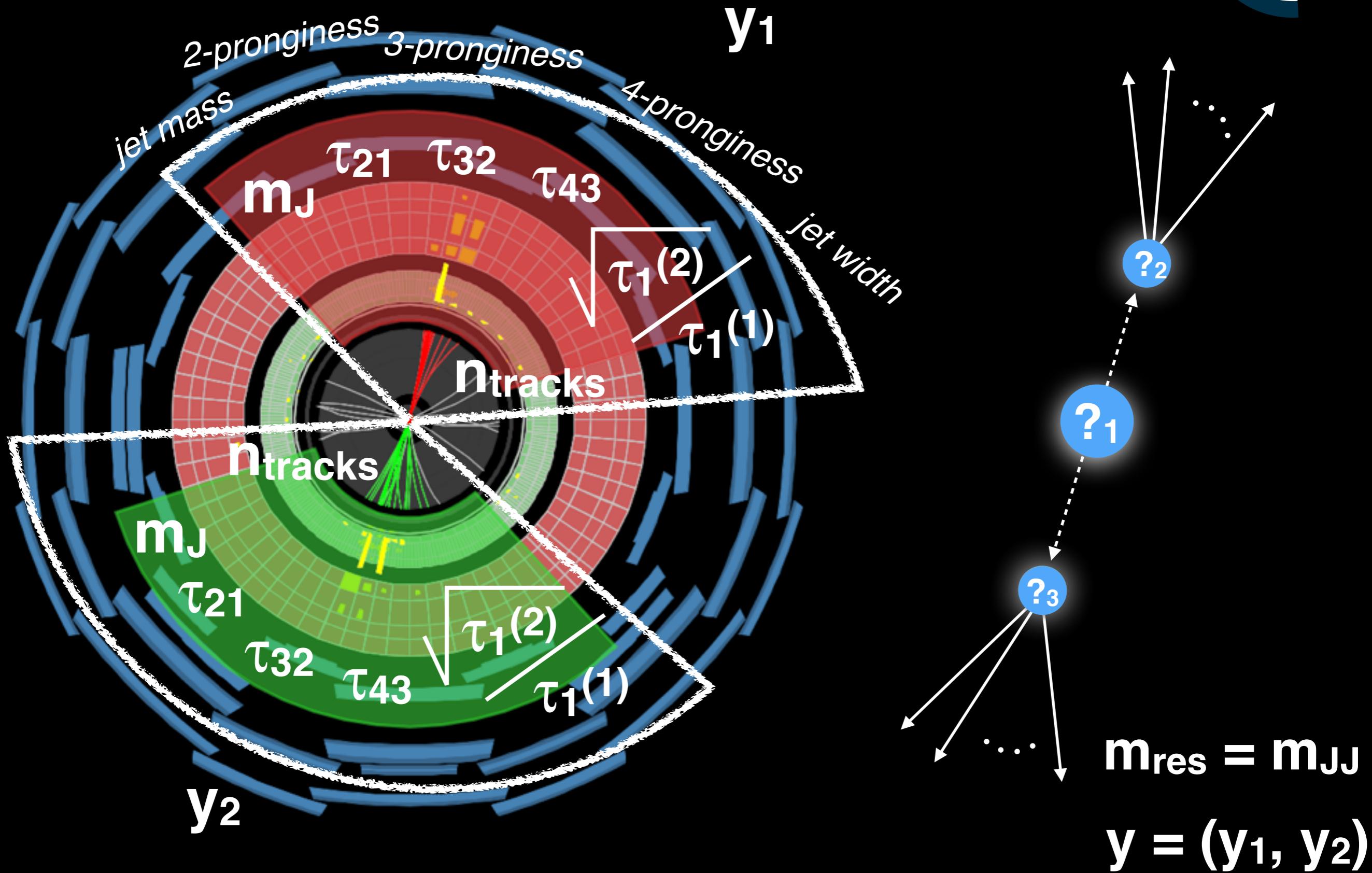
Application to dijet resonance search



Application to dijet resonance search



Application to dijet resonance search



Note about correlations

Why? $\sqrt{\frac{\tau_1(2)}{\tau_1(1)}}$

*Remember: if NN can learn m_{JJ} ,
violate CWoLa assumptions.*

$$m_{JJ} \sim p_{T,J_1} p_{T,J_2} \Delta R_{JJ}^2 \Rightarrow \text{Can't learn jet } p_T.$$

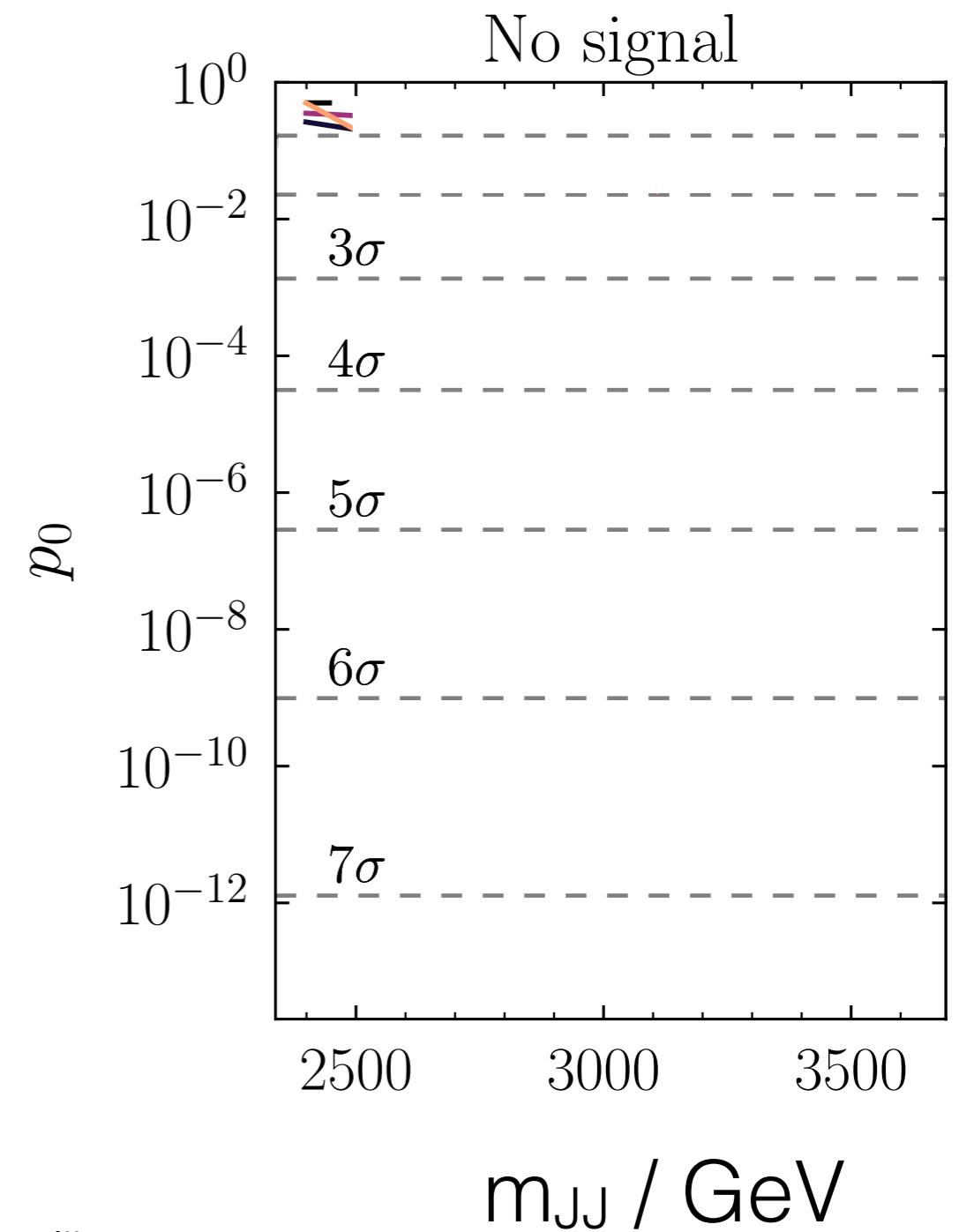
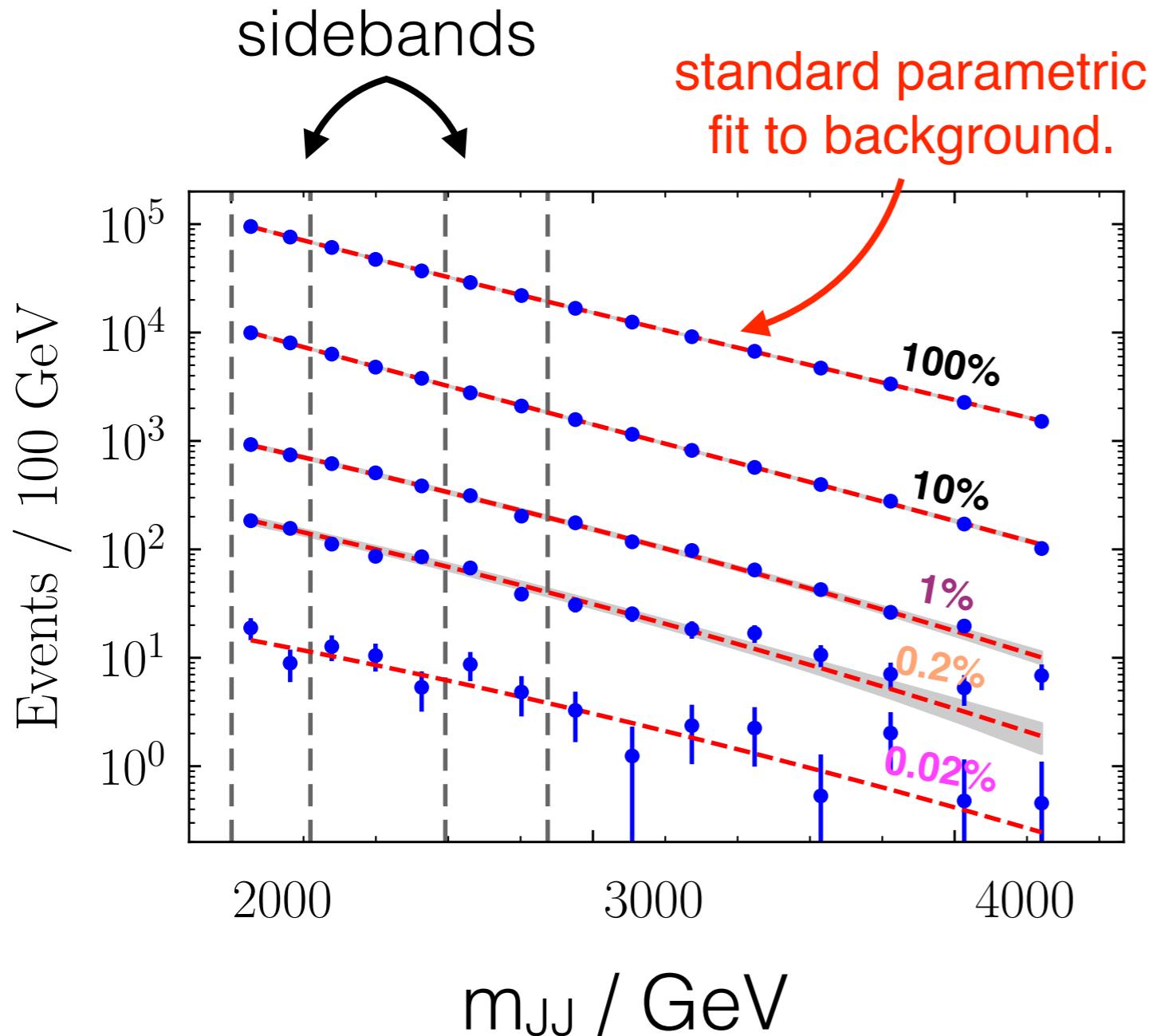
$$\tau_1^{(j)} \sim \frac{1}{p_{T,\text{jet}}} \sum_{i \in \text{jet}} p_{T,i} \Delta R(\text{jet}, i)^j \quad \Delta R \sim 1/\gamma \sim m/p_T$$

\Rightarrow can't have both m and $\tau_1^{(2)}$

\Rightarrow since $\tau_1^{(2)} \sim 1/\gamma^2$ and $\tau_1^{(1)} \sim 1/\gamma$, use $\sqrt{\tau_1^{(2)}} / \tau_1^{(1)}$

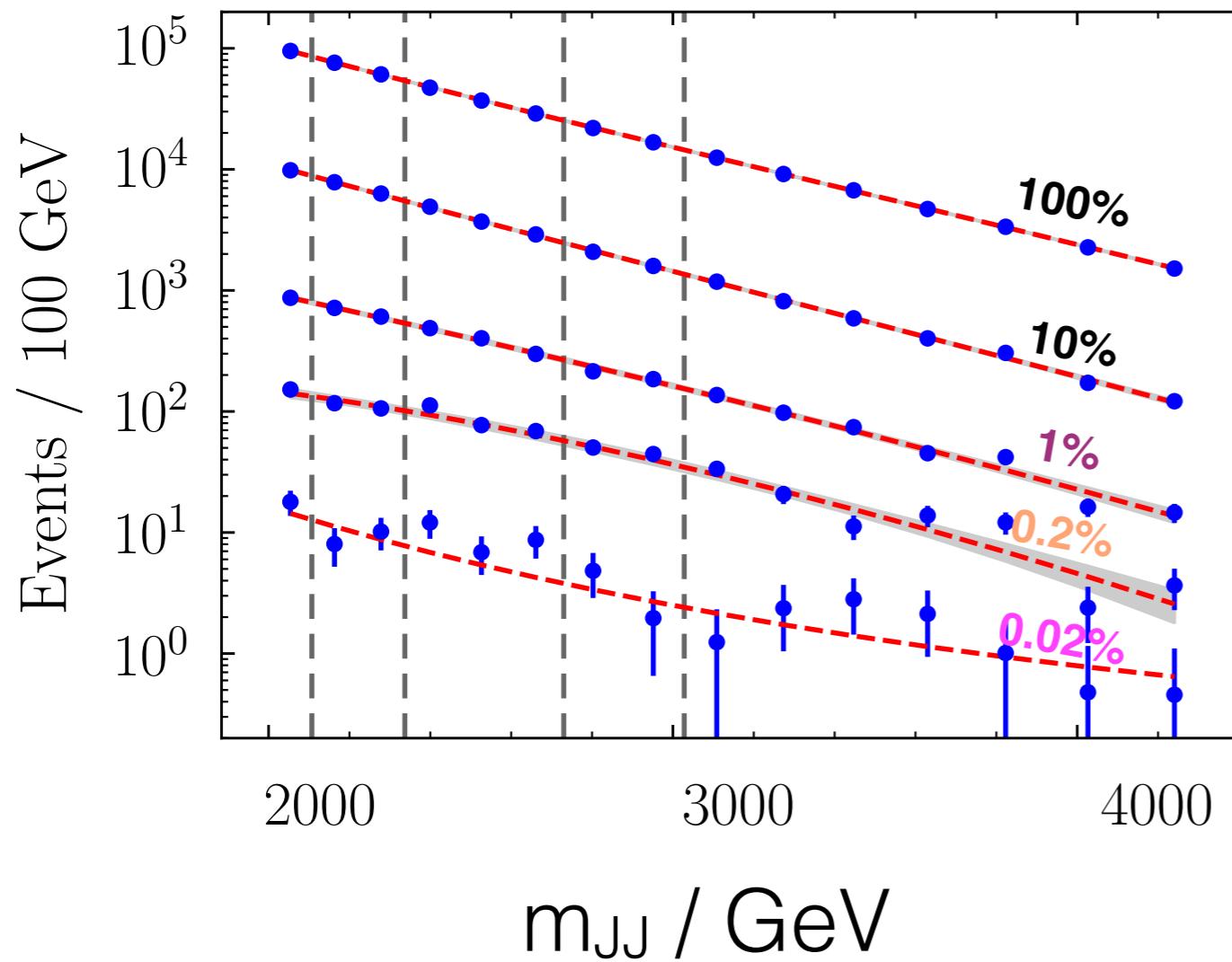
Could also imagine doing this automatically with ML - see
C. Shimmin, P. Sadowski, P. Baldi, E. Weik, D. Whiteson et al. 1703.03507.

Application to dijets



- no cut on NN
- most 10% signal-region-like
- most 1% signal-region-like
- most 0.2% signal-region-like
- most 0.02% signal-region-like

Application to dijets

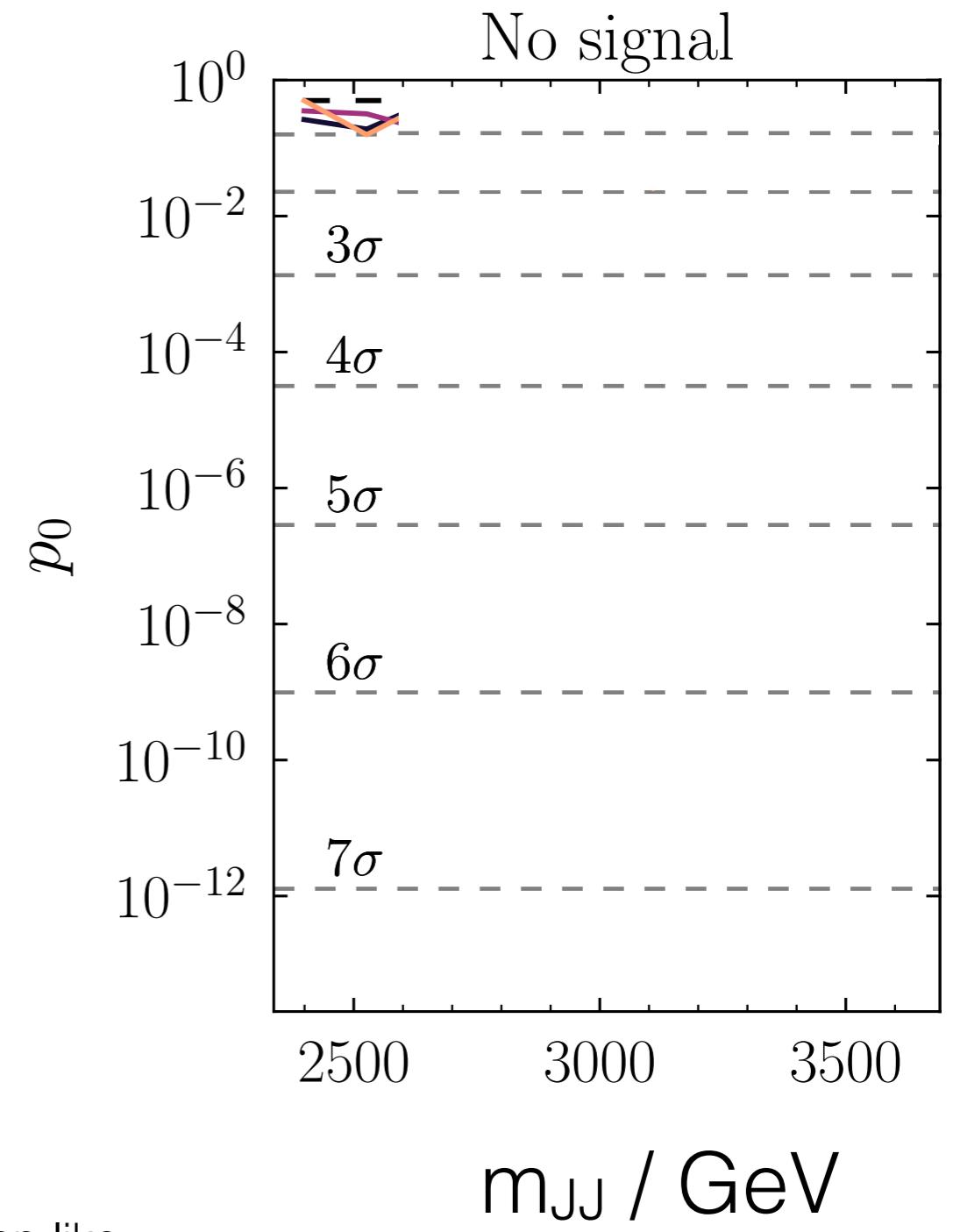


---- no cut on NN

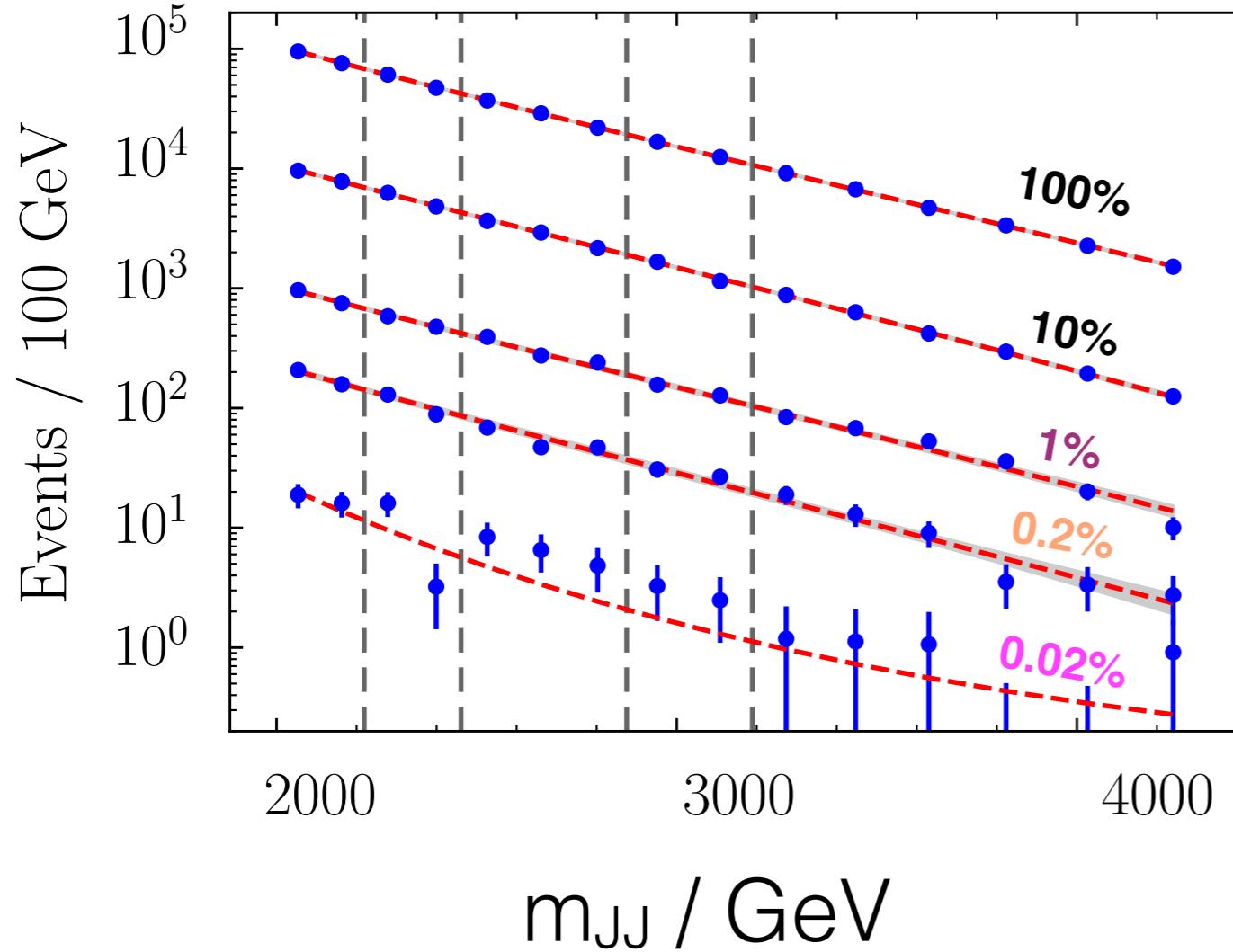
— most 10% signal-region-like

— most 1% signal-region-like

— most 0.2% signal-region-like



Application to dijets

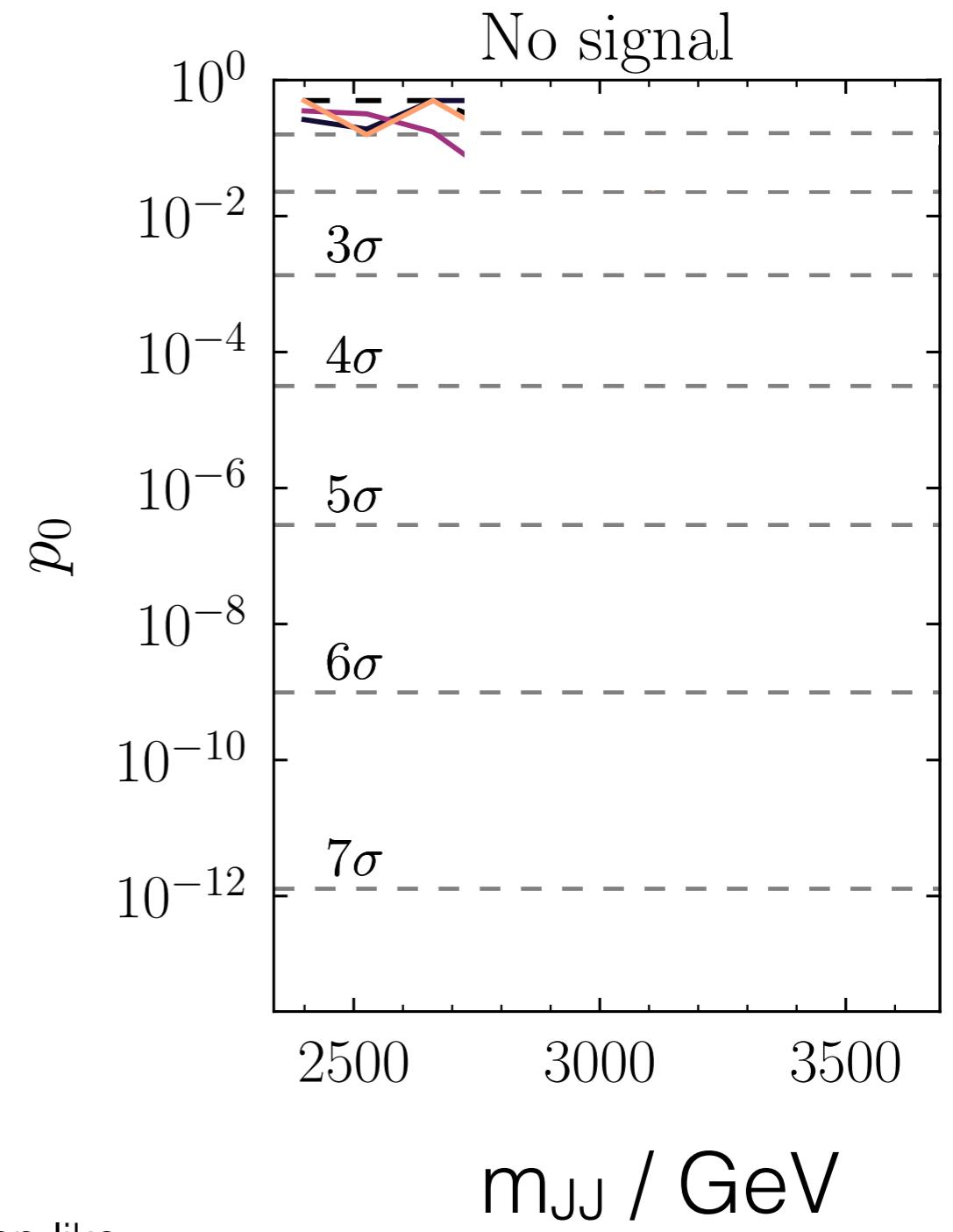


---- no cut on NN

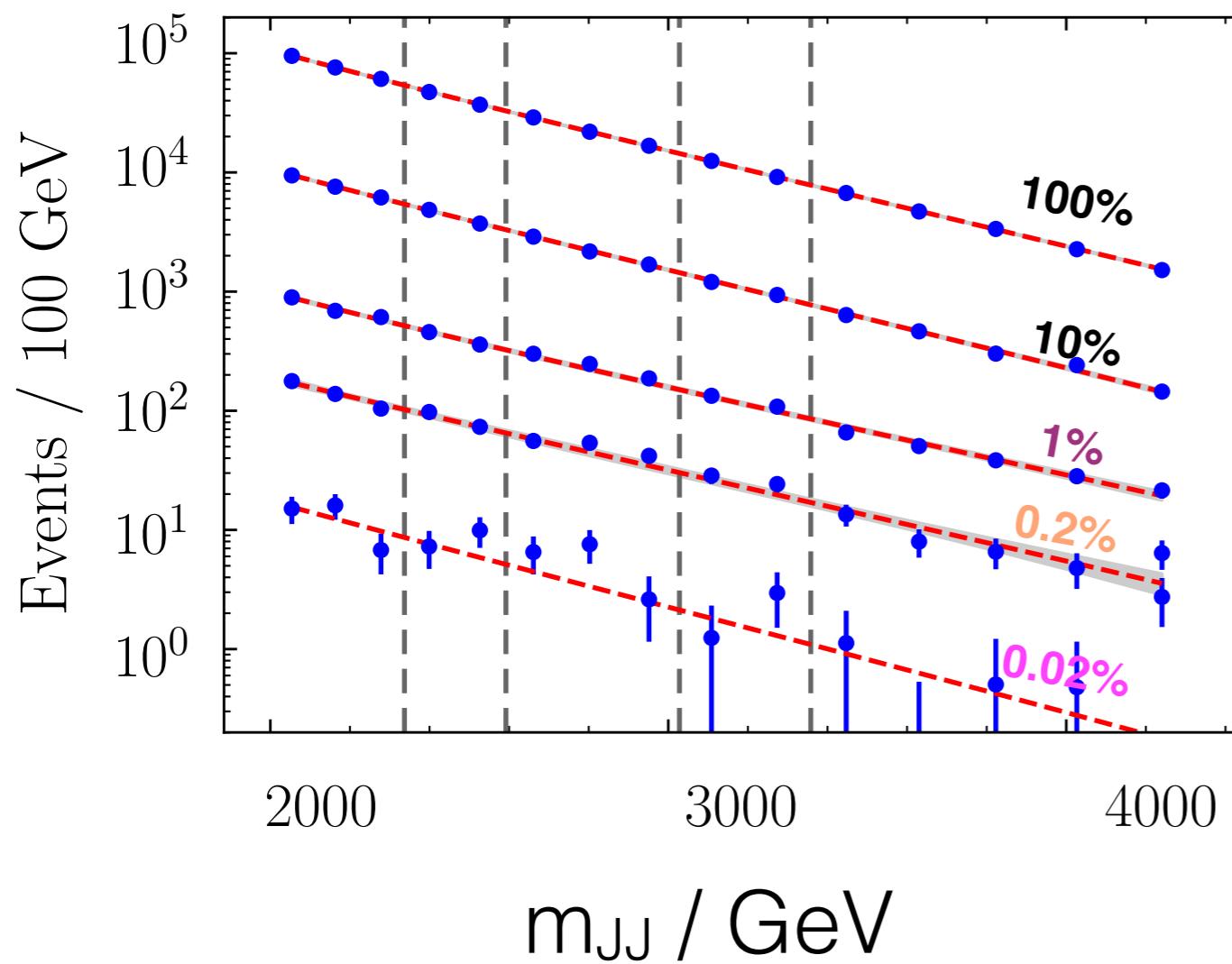
— most 10% signal-region-like

— most 1% signal-region-like

— most 0.2% signal-region-like



Application to dijets

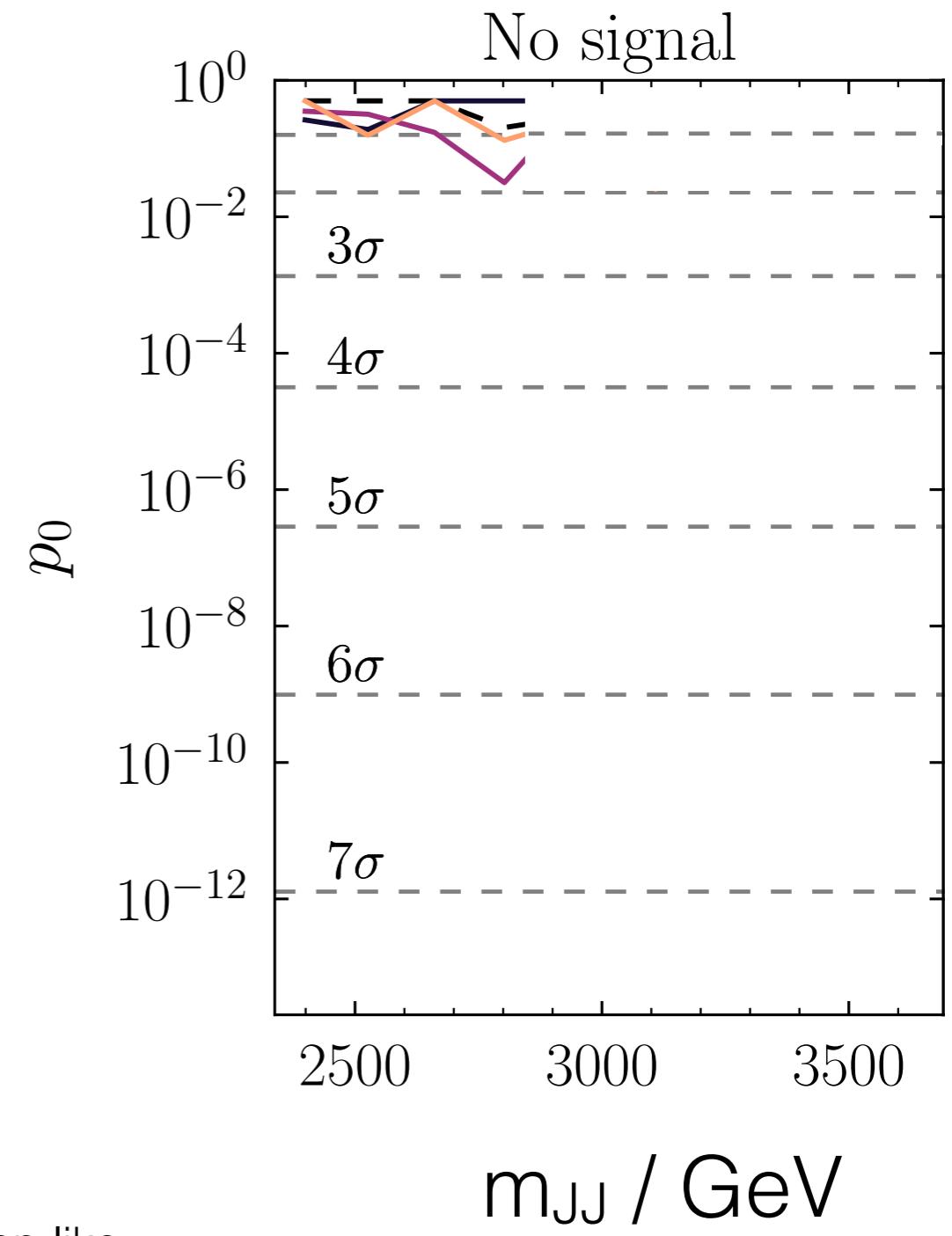


---- no cut on NN

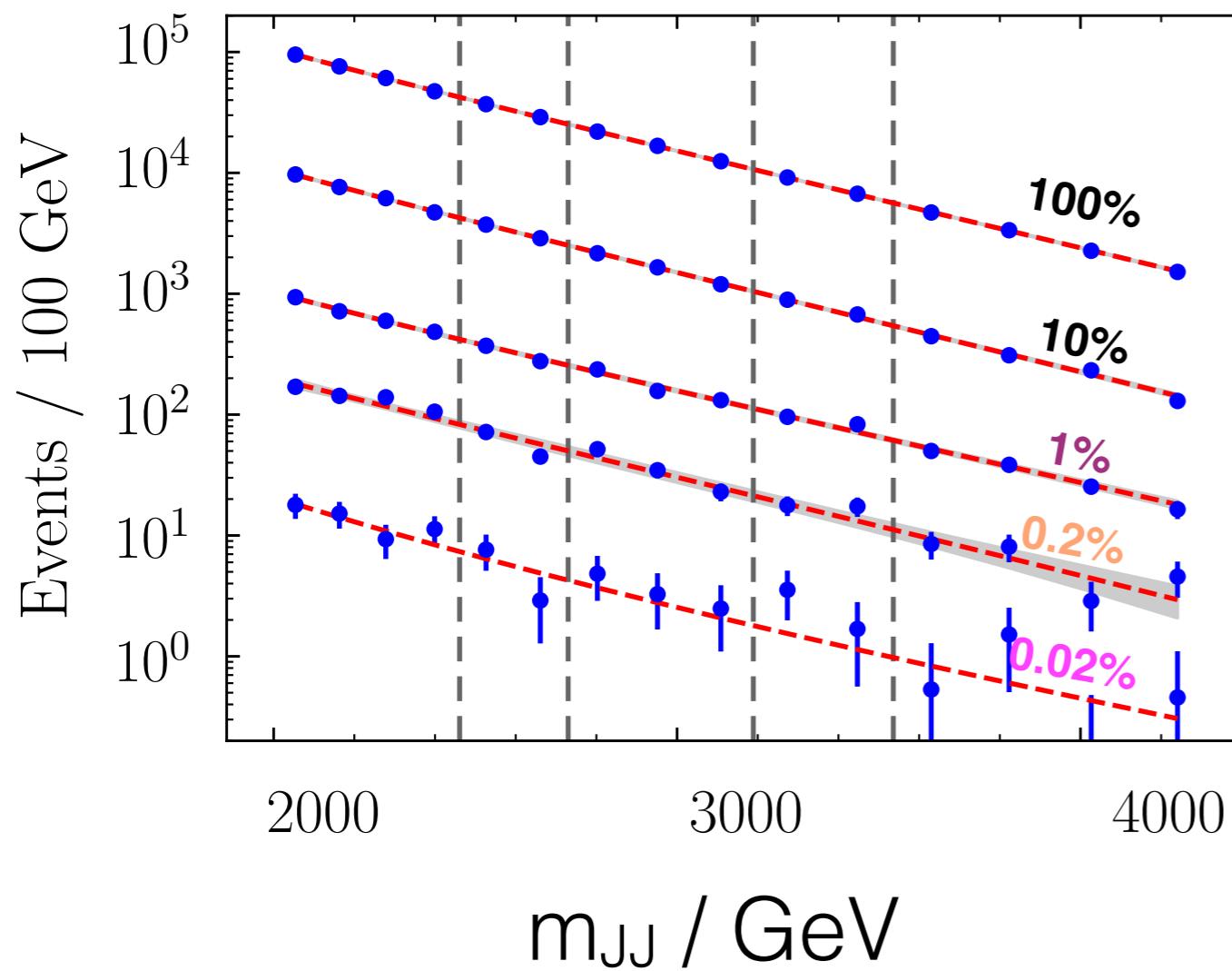
— most 10% signal-region-like

— most 1% signal-region-like

— most 0.2% signal-region-like

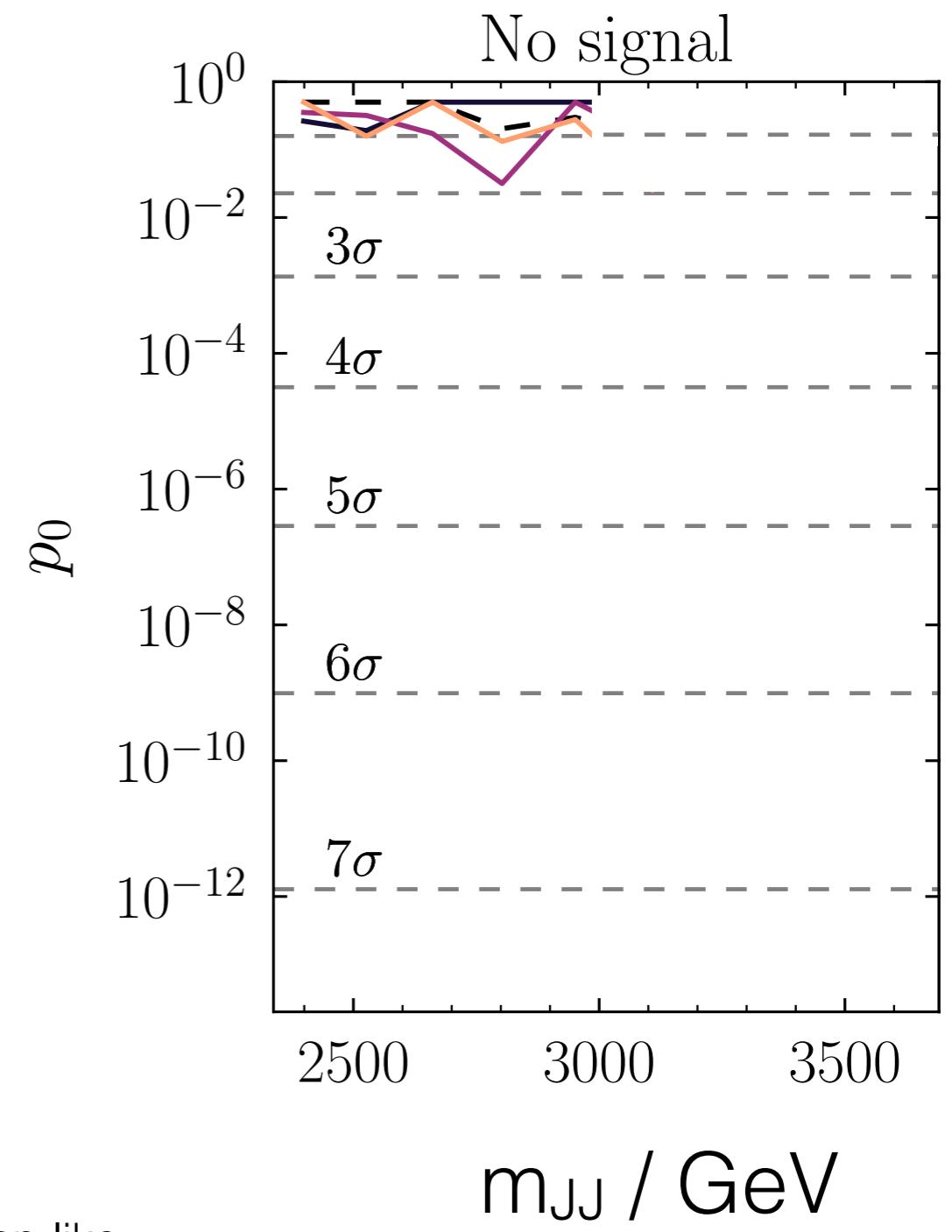


Application to dijets



Legend:

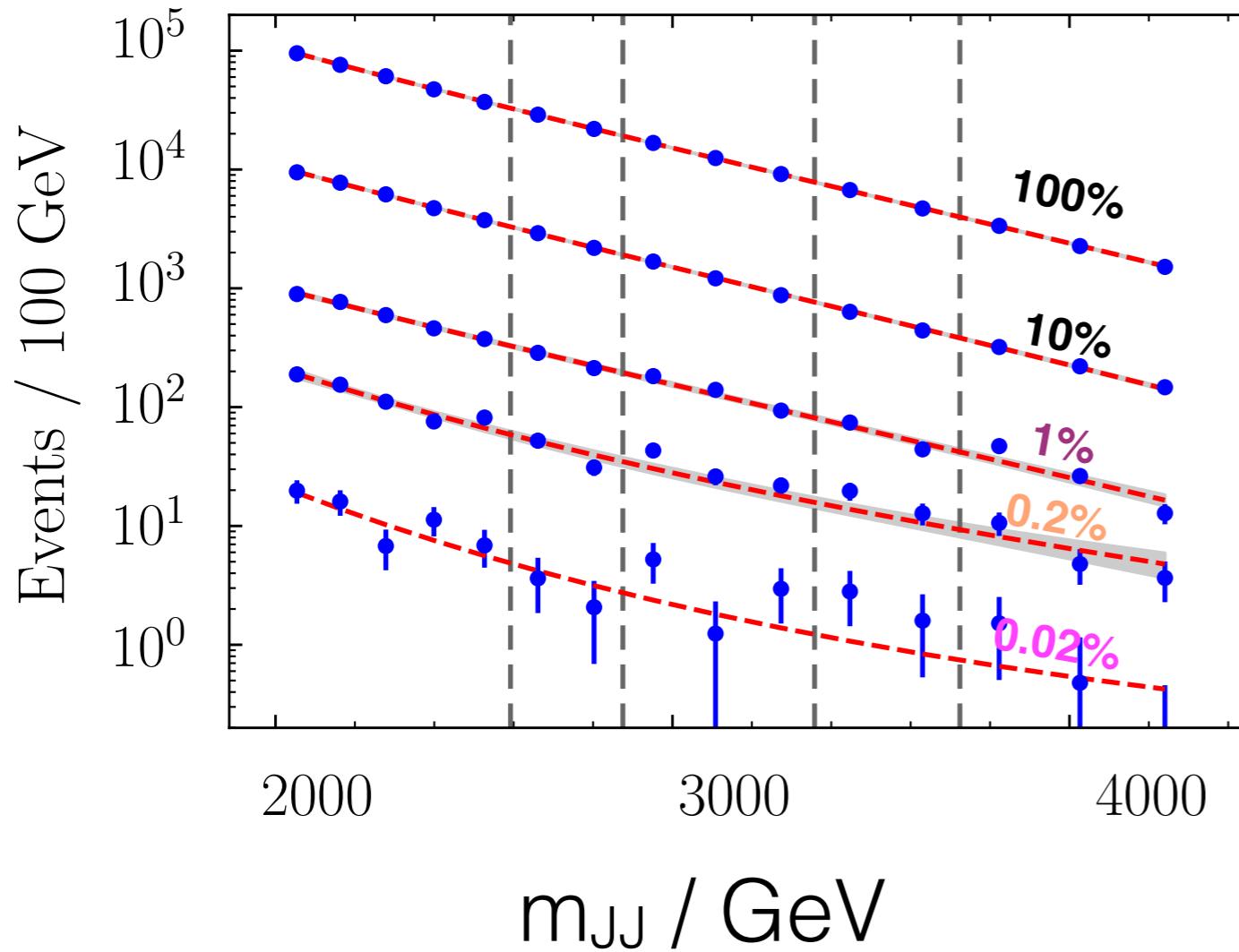
- no cut on NN
- most 10% signal-region-like
- most 1% signal-region-like
- most 0.2% signal-region-like
- most 0.02% signal-region-like



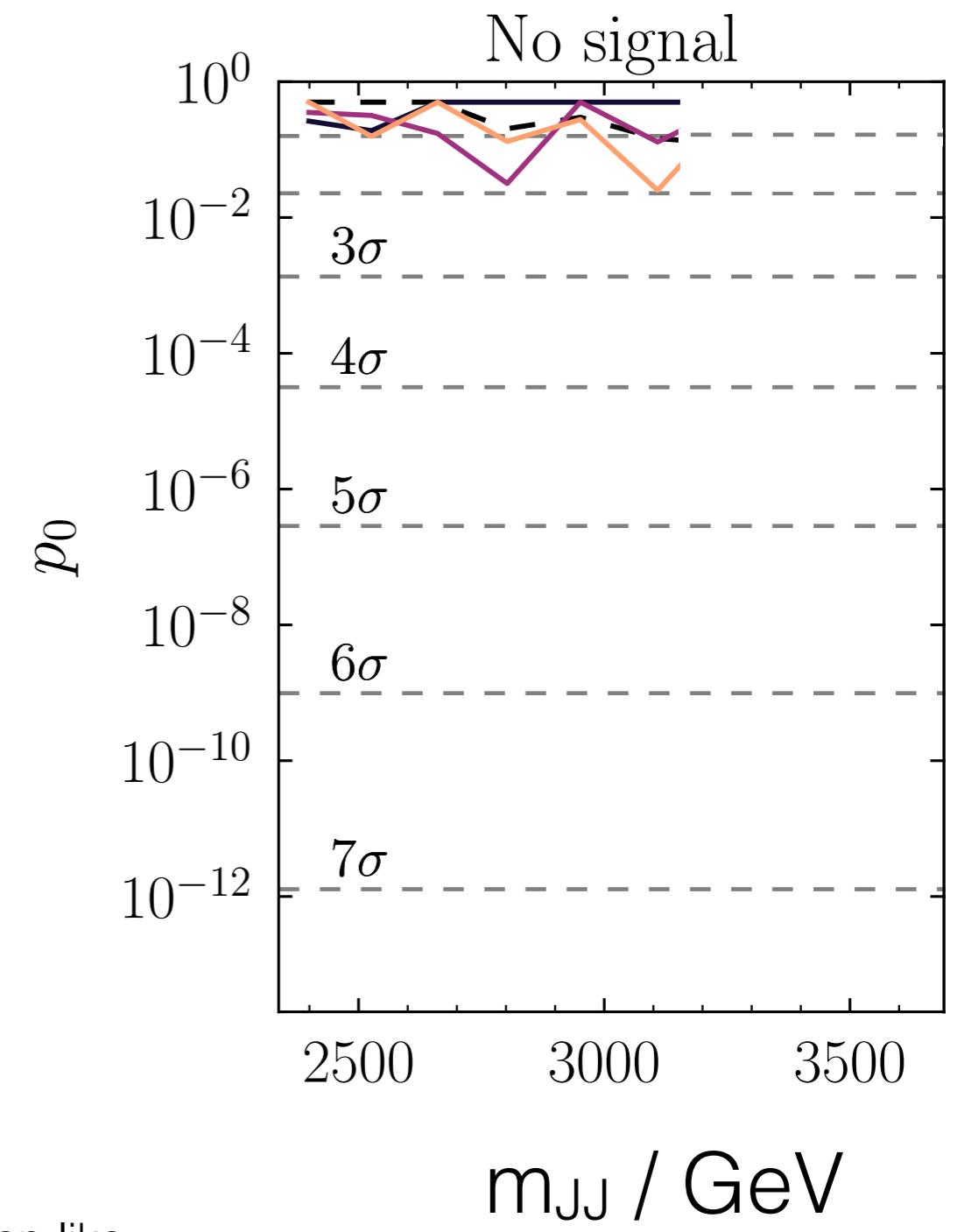
Legend:

- no cut on NN
- most 10% signal-region-like
- most 1% signal-region-like
- most 0.2% signal-region-like
- most 0.02% signal-region-like

Application to dijets

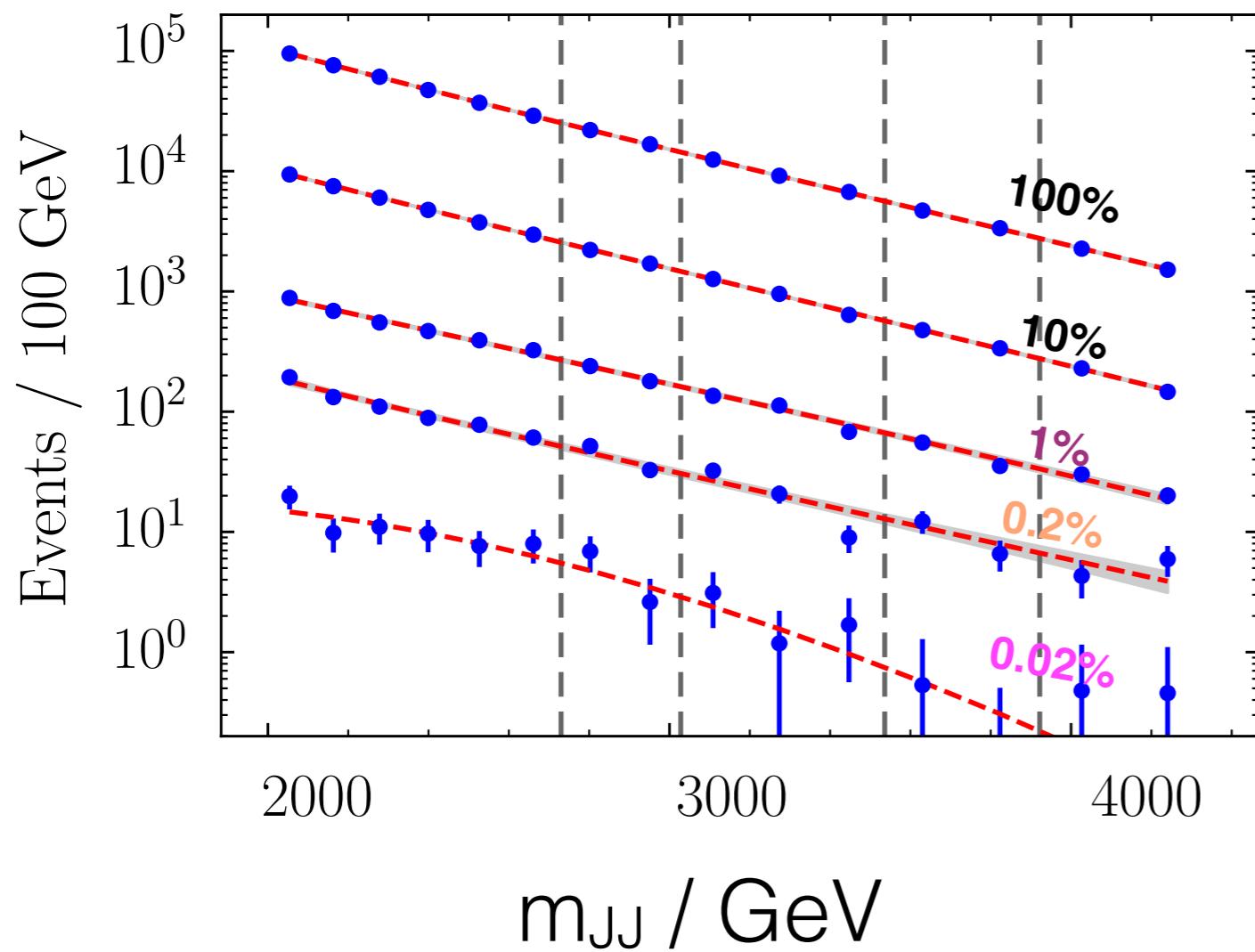


--- no cut on NN
— most 10% signal-region-like



— most 1% signal-region-like
— most 0.2% signal-region-like

Application to dijets

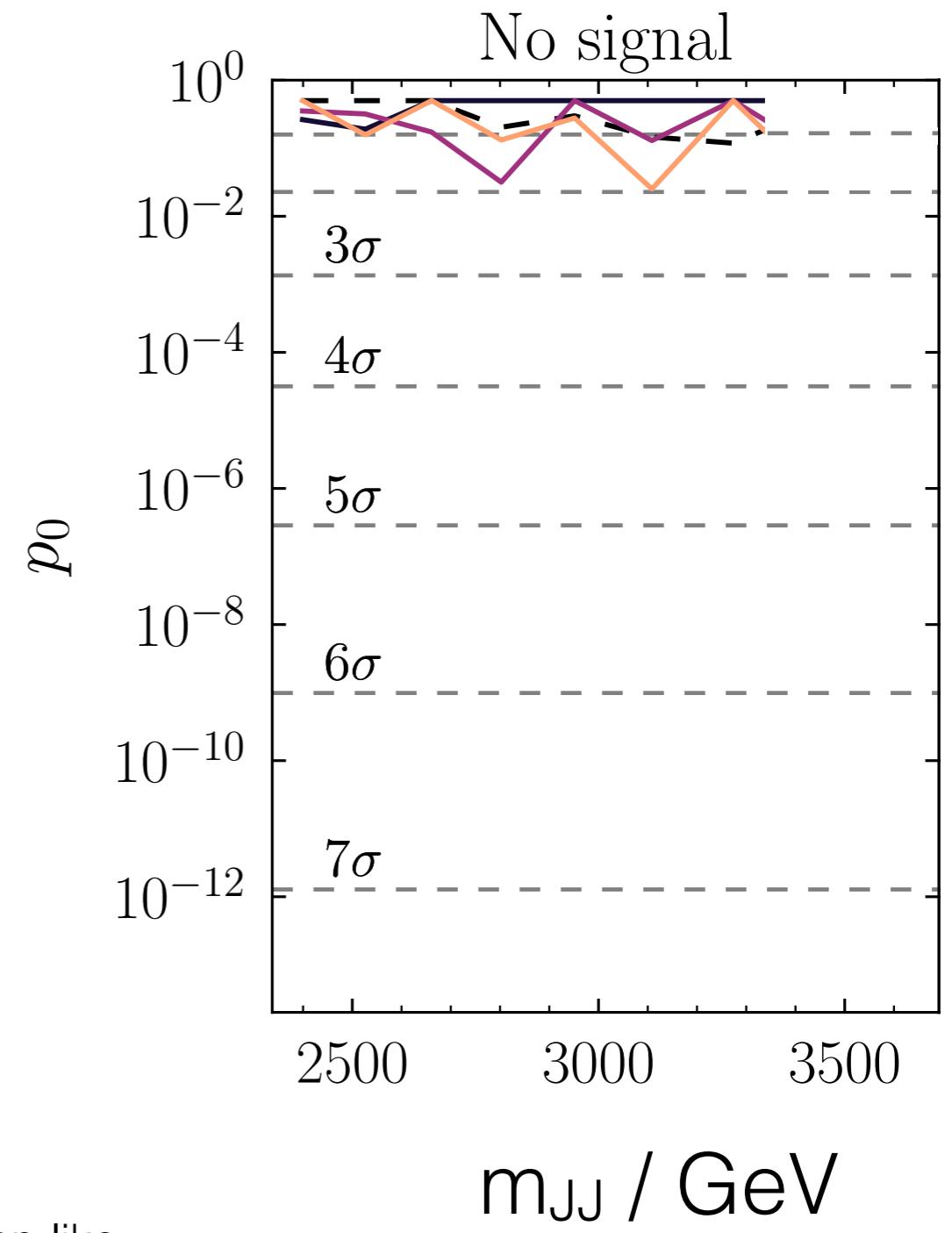


---- no cut on NN

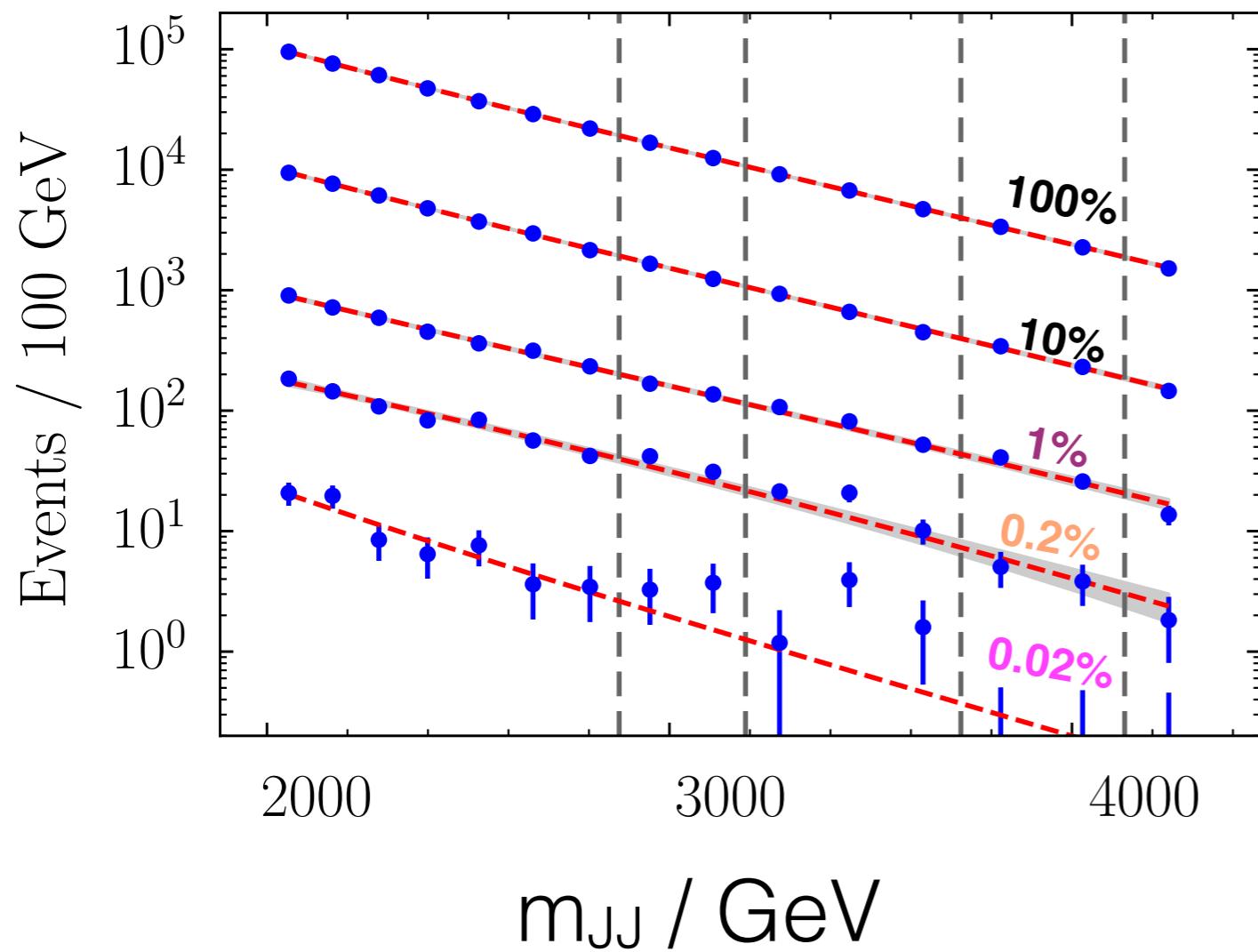
— most 10% signal-region-like

— most 1% signal-region-like

— most 0.2% signal-region-like



Application to dijets

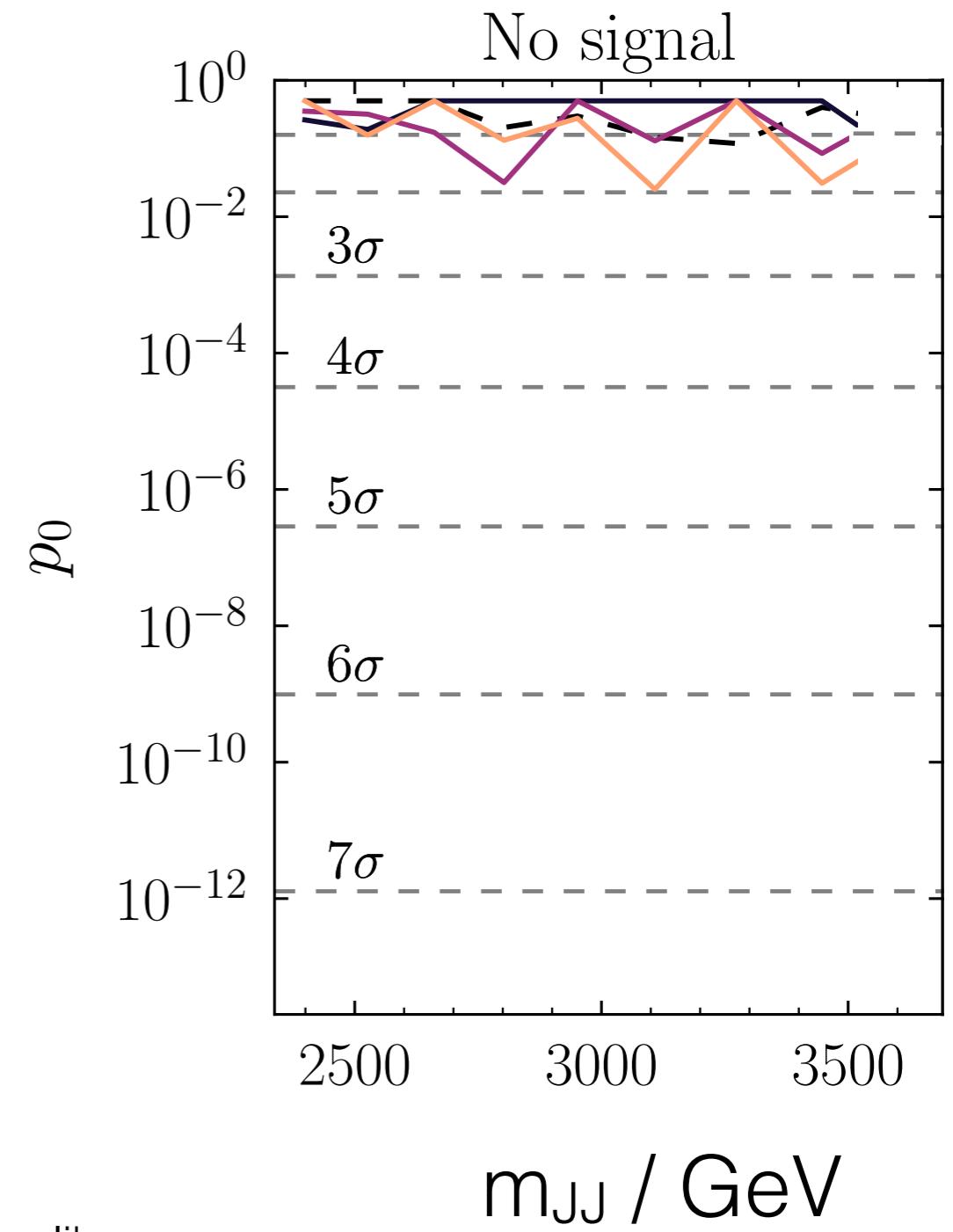


---- no cut on NN

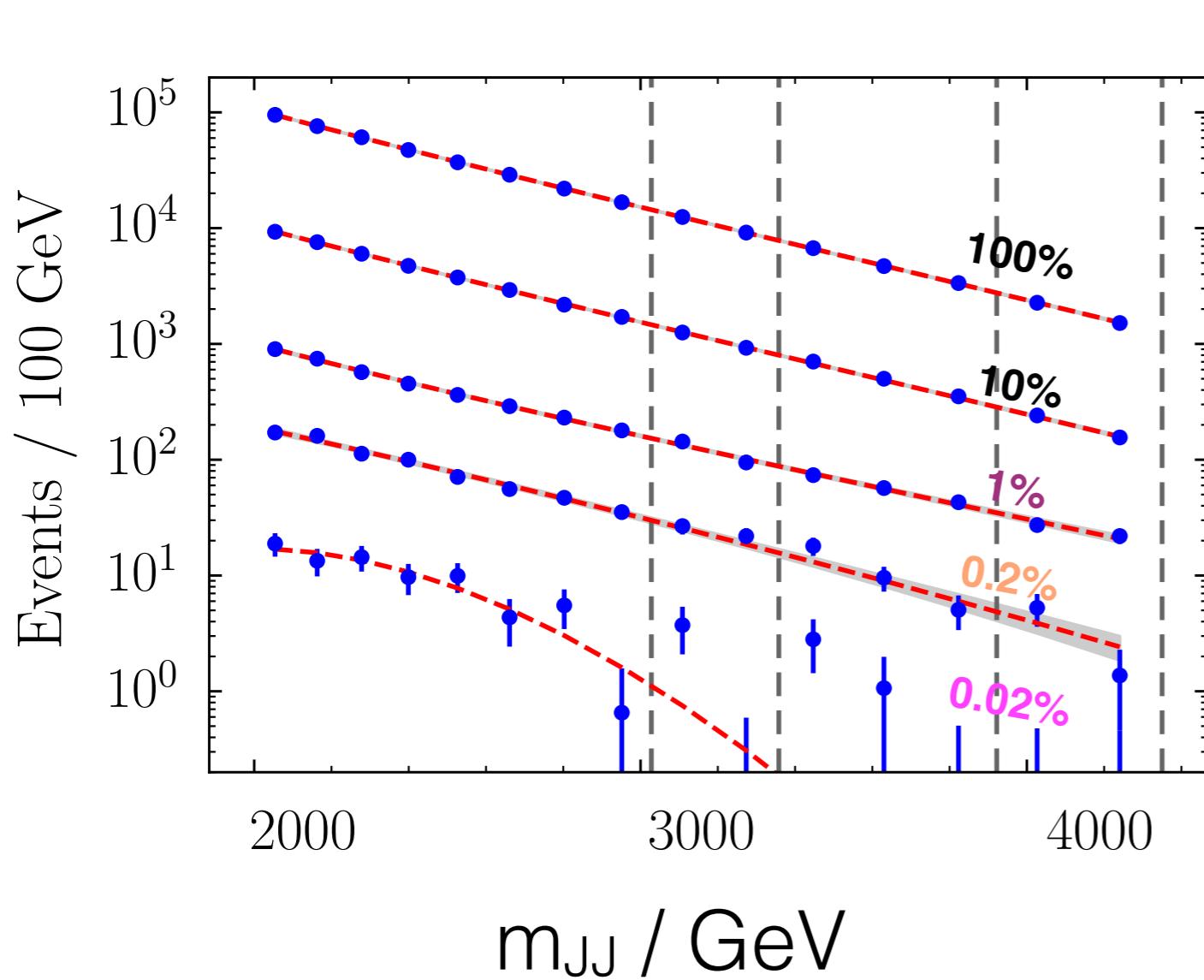
— most 10% signal-region-like

— most 1% signal-region-like

— most 0.2% signal-region-like



Application to dijets

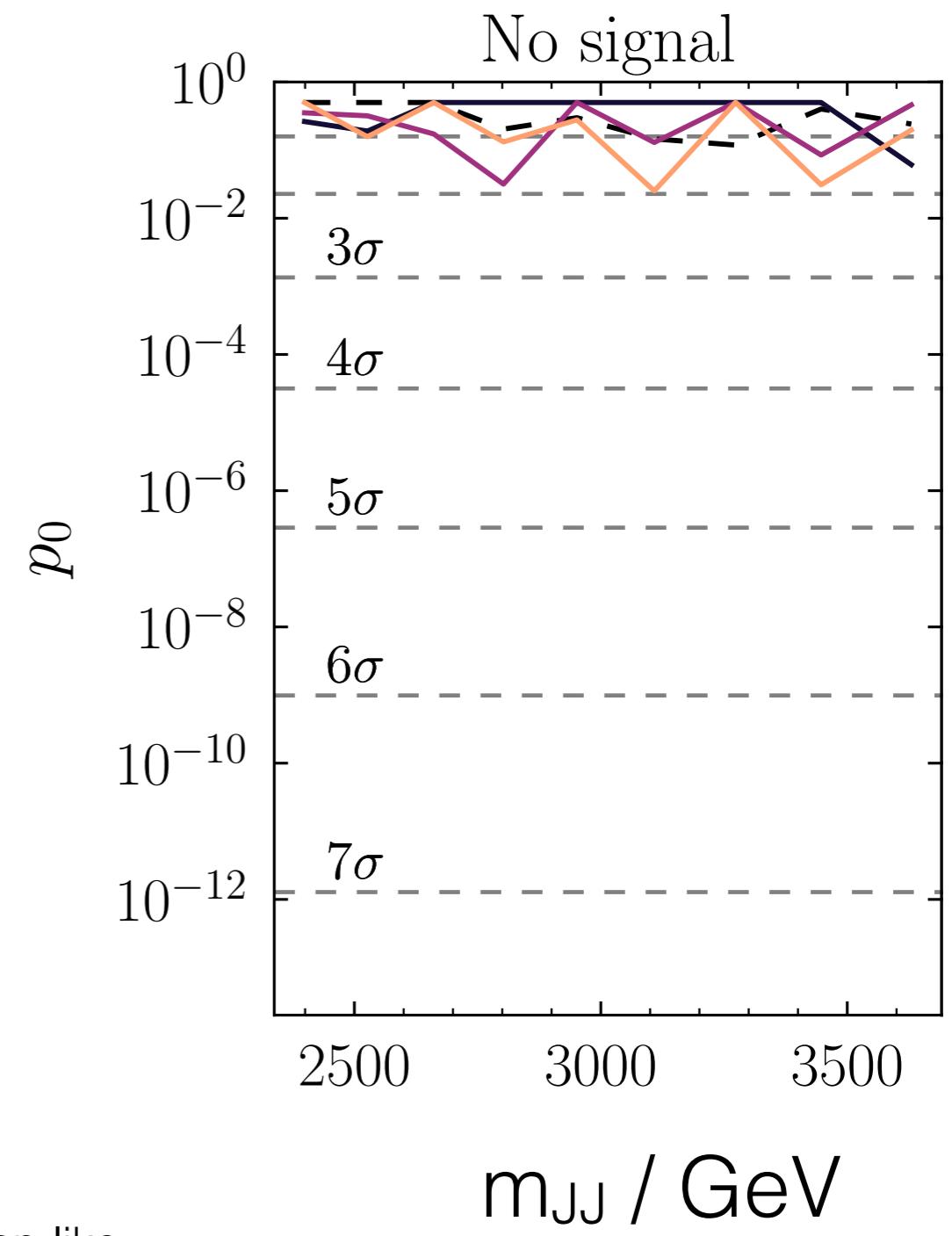


---- no cut on NN

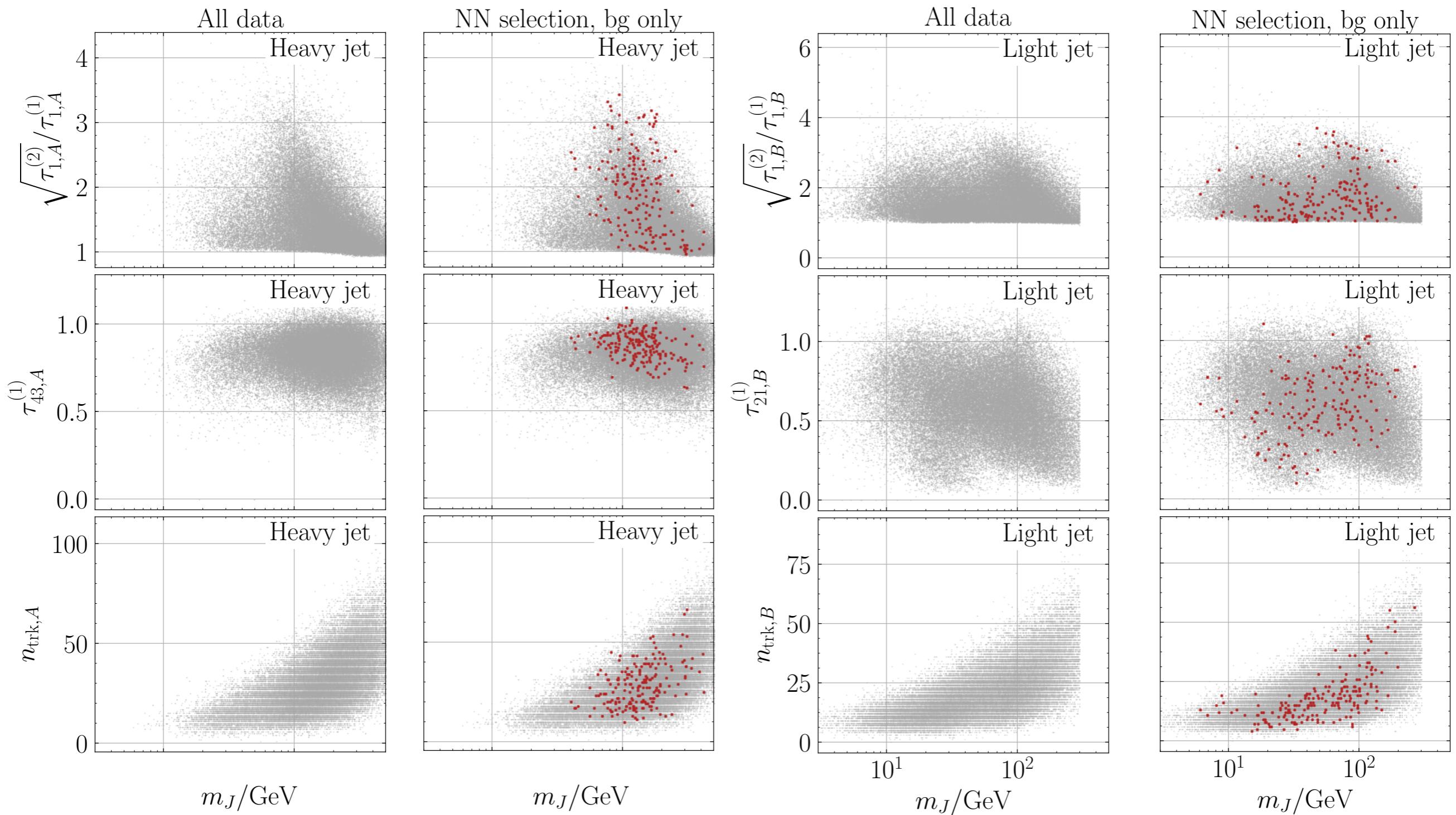
— most 10% signal-region-like

— most 1% signal-region-like

— most 0.2% signal-region-like

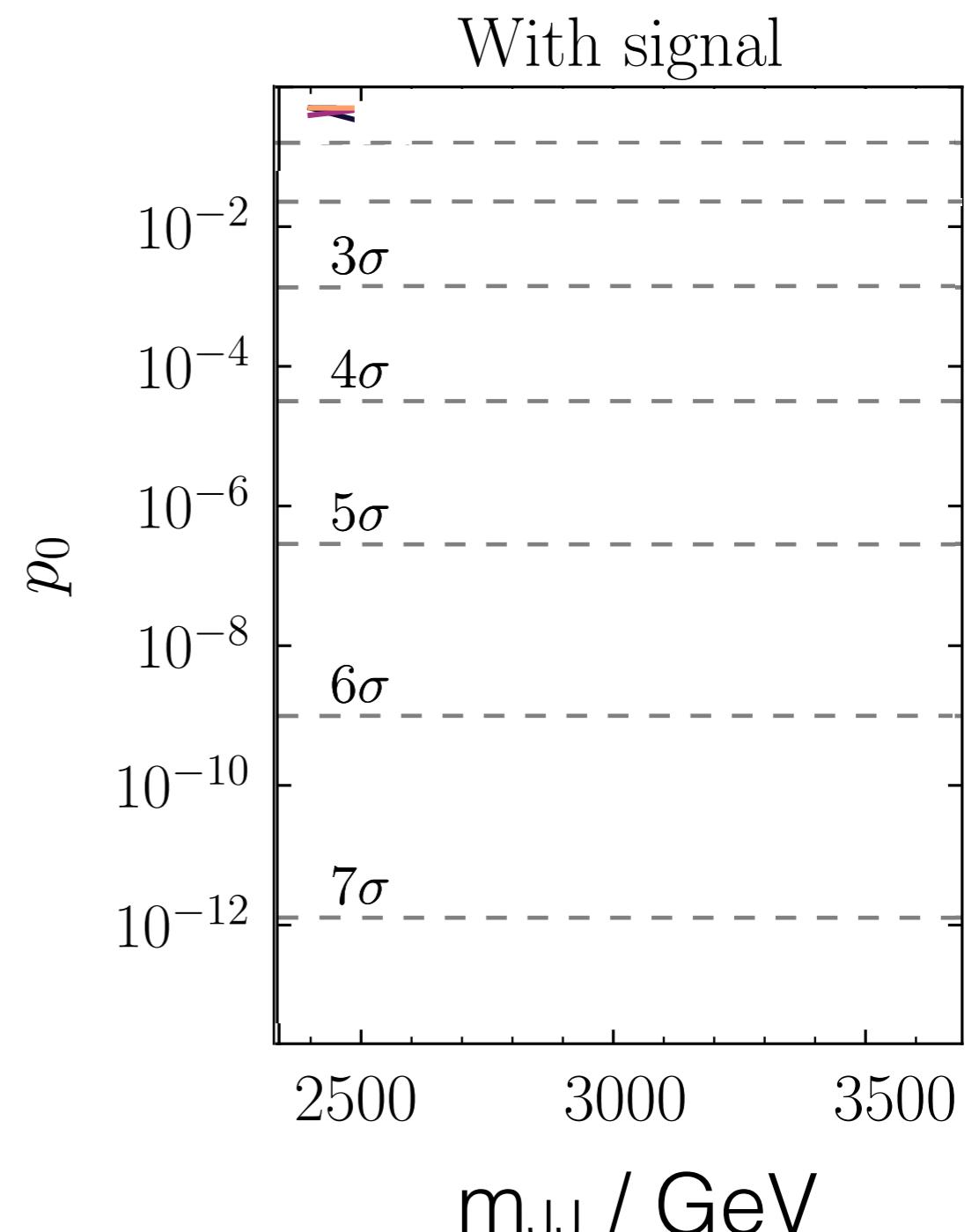
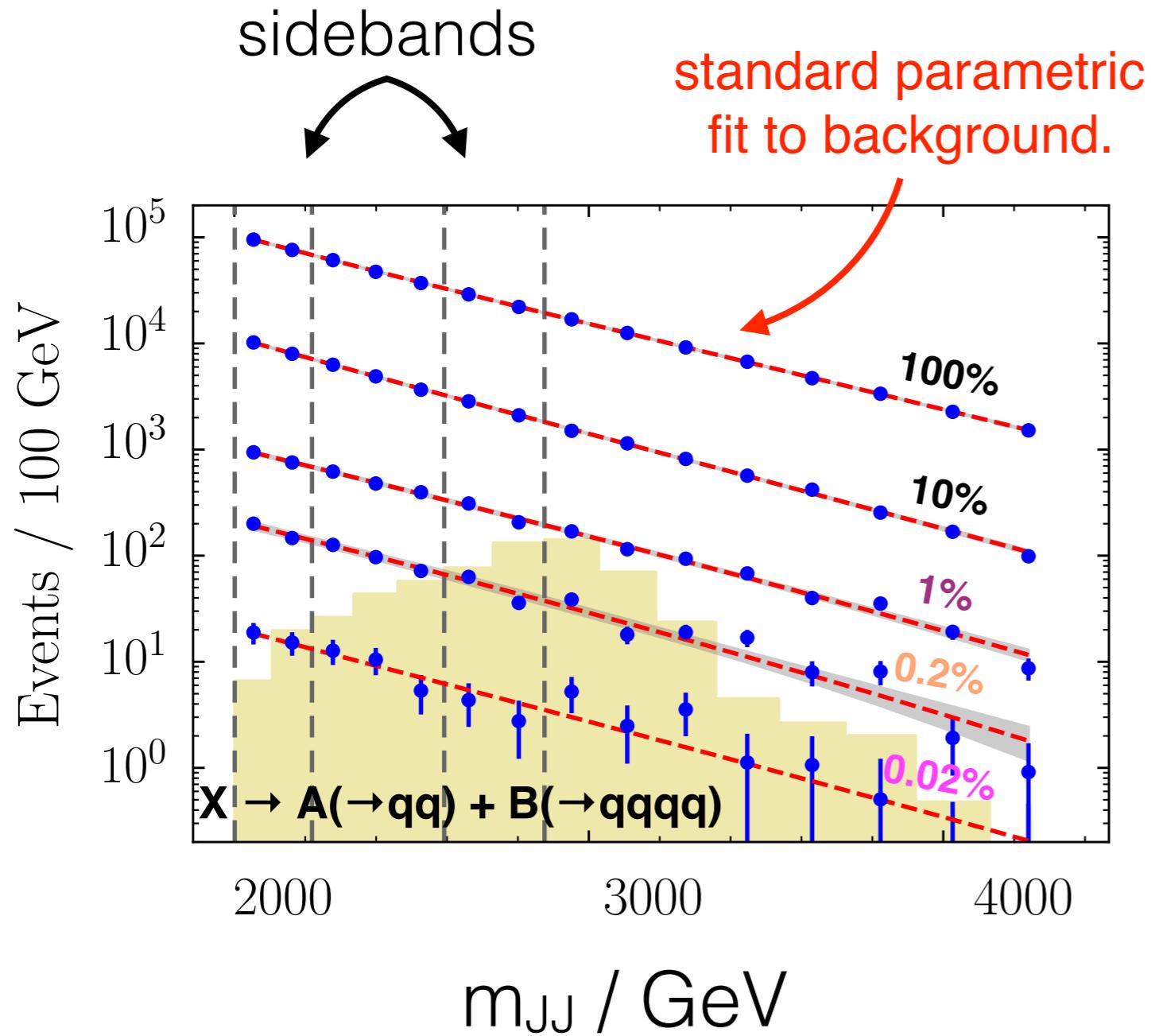


What is the network learning?



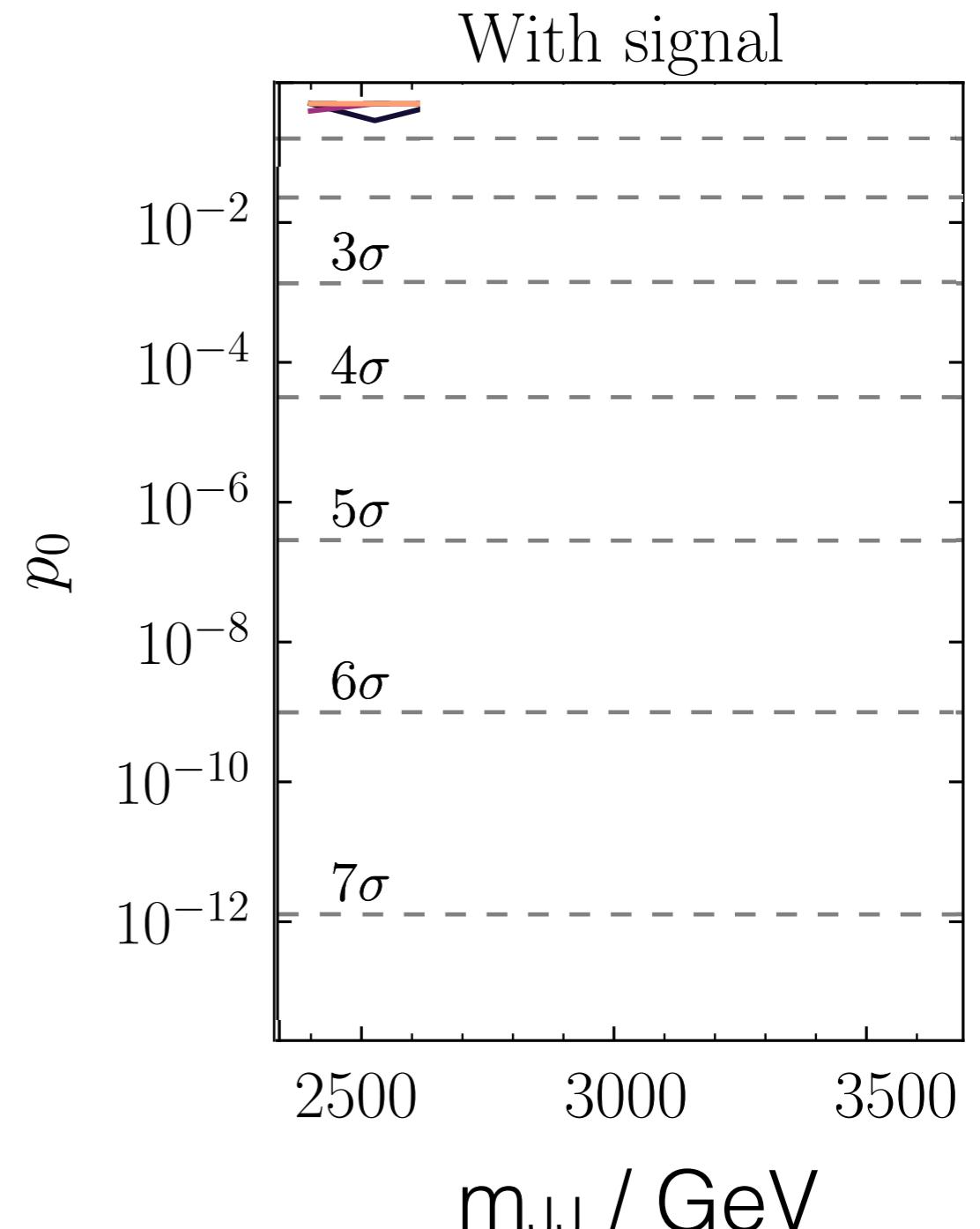
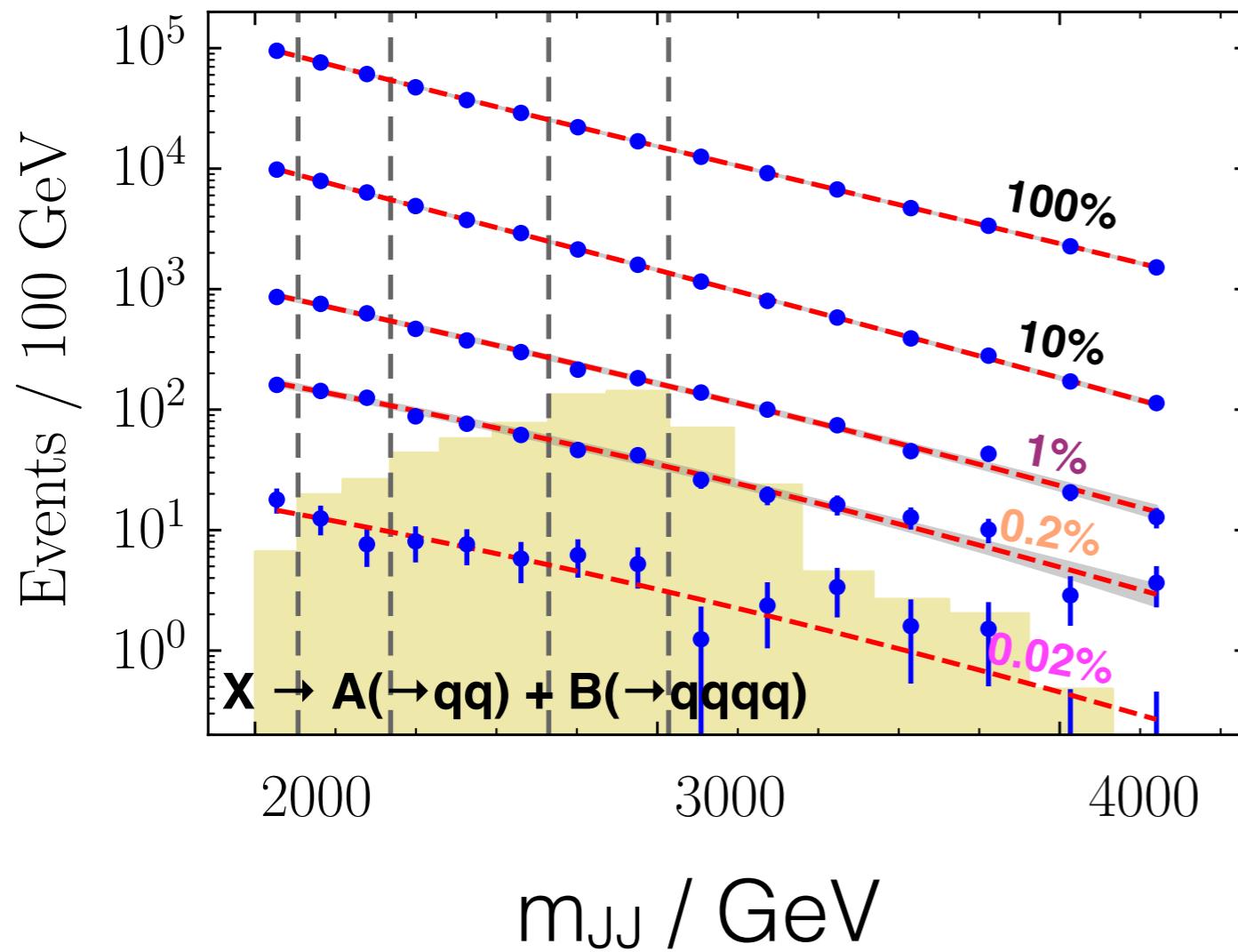
Learns nothing ... as desired !

...and when there is a signal?



- no cut on NN
- most 10% signal-region-like
- most 1% signal-region-like
- most 0.2% signal-region-like

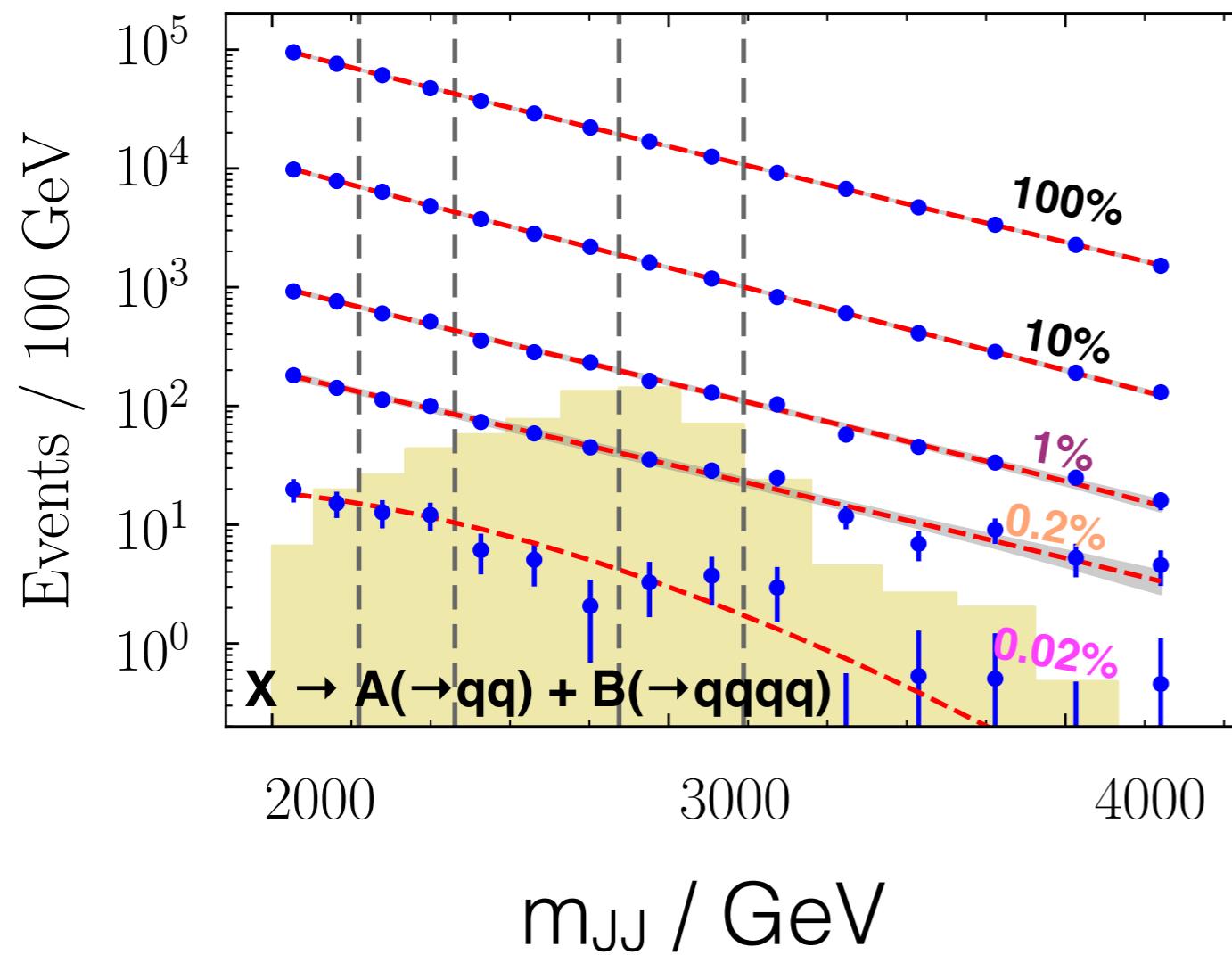
...and when there is a signal?



Legend:

- no cut on NN
- most 10% signal-region-like
- most 1% signal-region-like
- most 0.2% signal-region-like
- most 0.02% signal-region-like

...and when there is a signal?

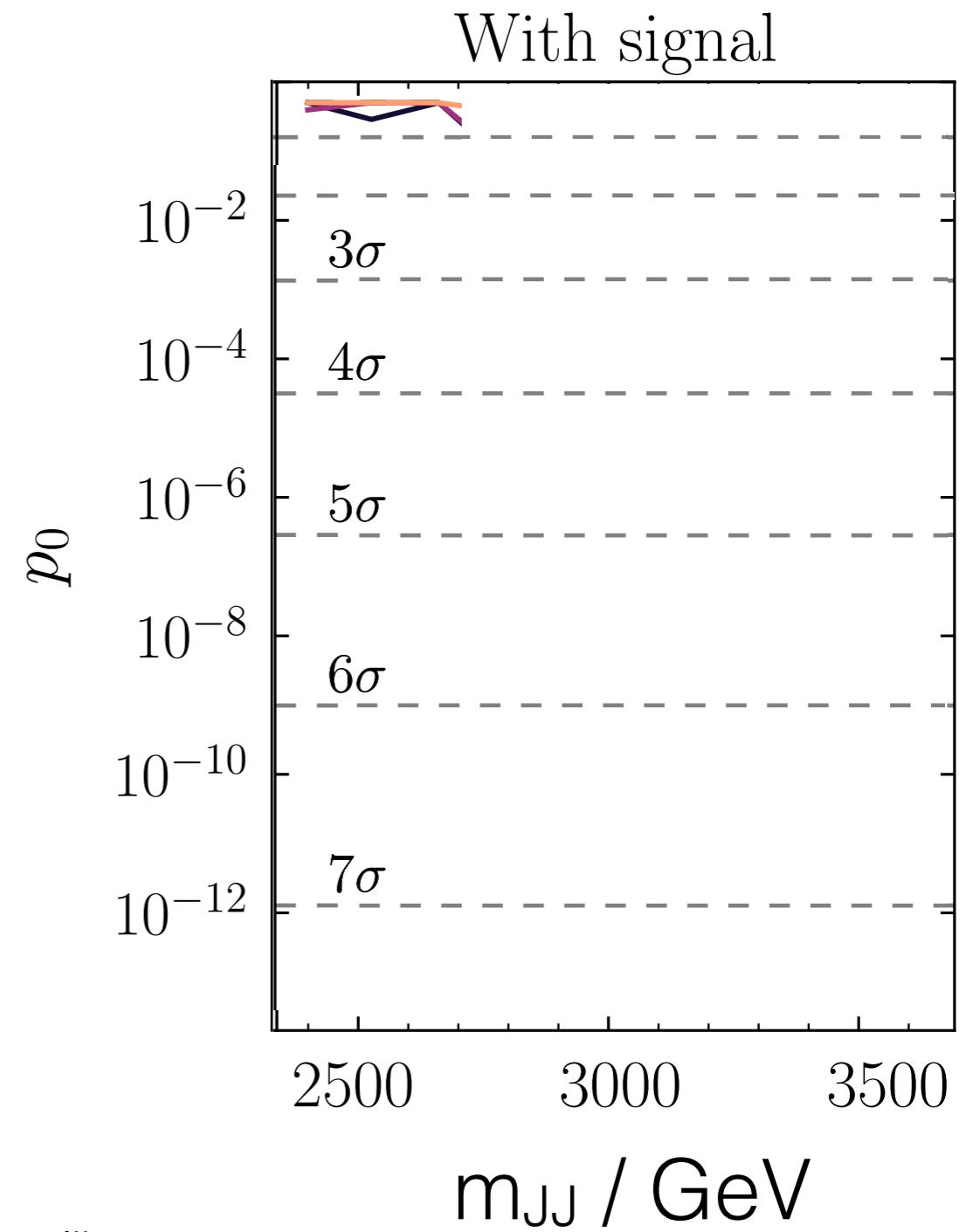


---- no cut on NN

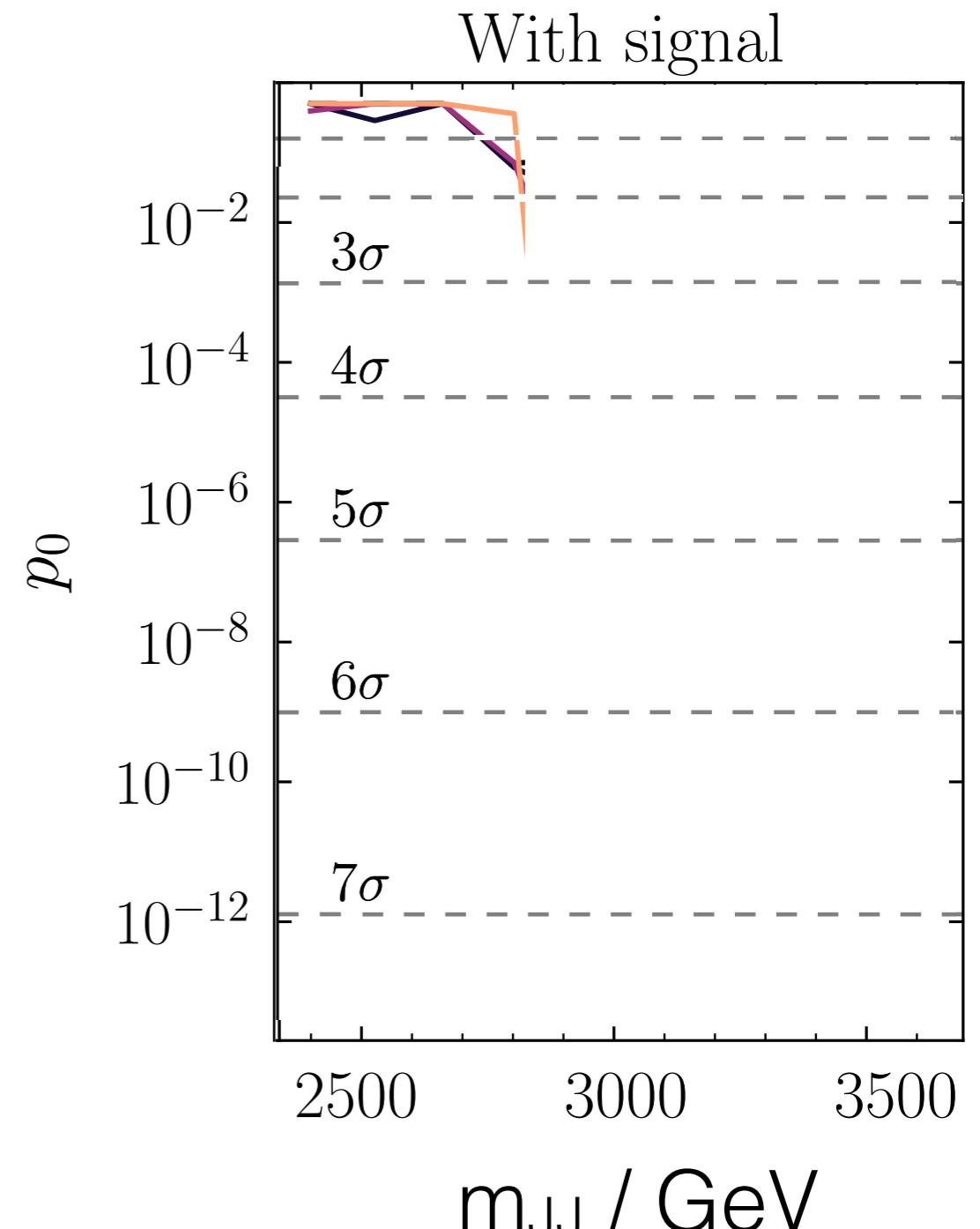
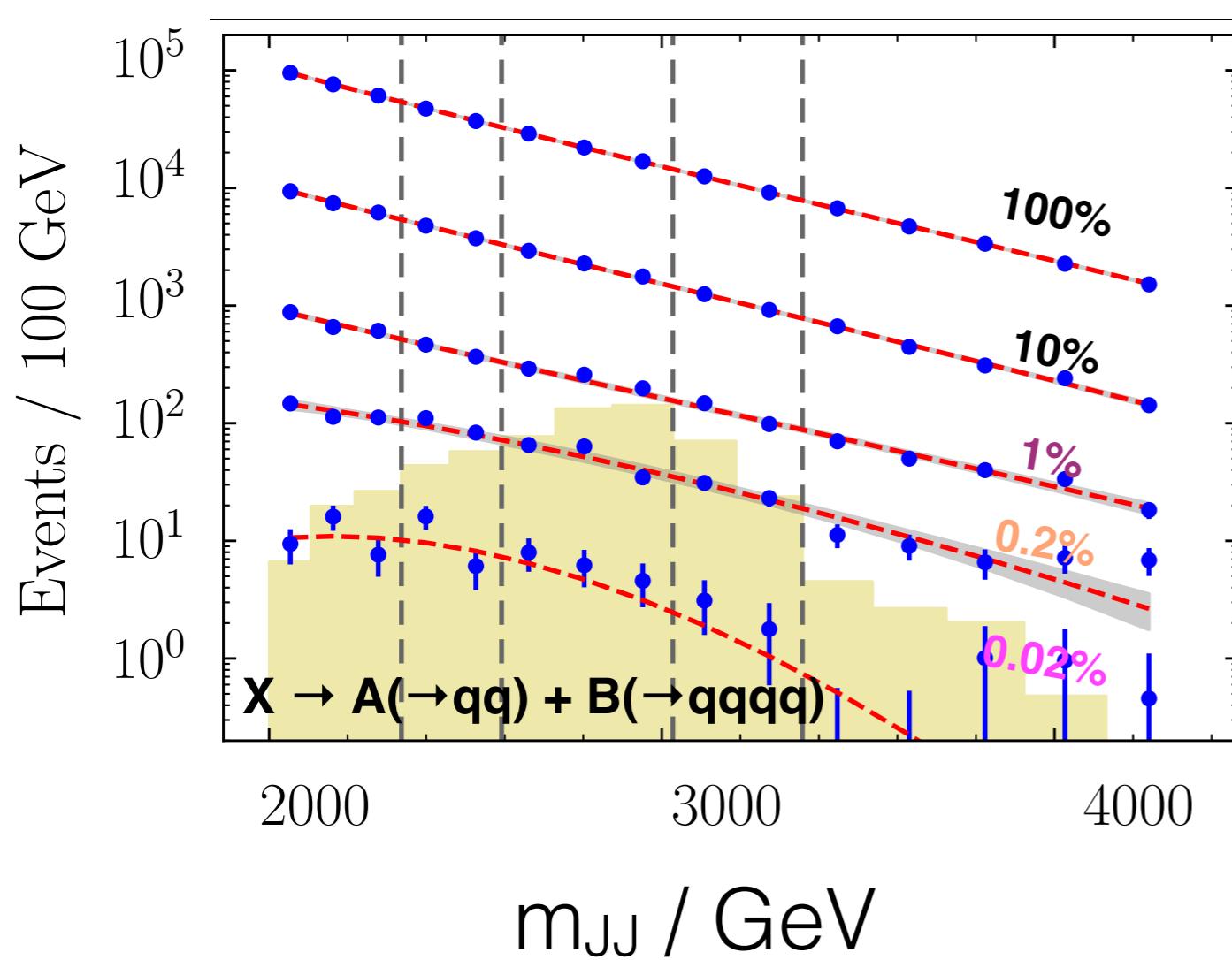
— most 10% signal-region-like

— most 1% signal-region-like

— most 0.2% signal-region-like

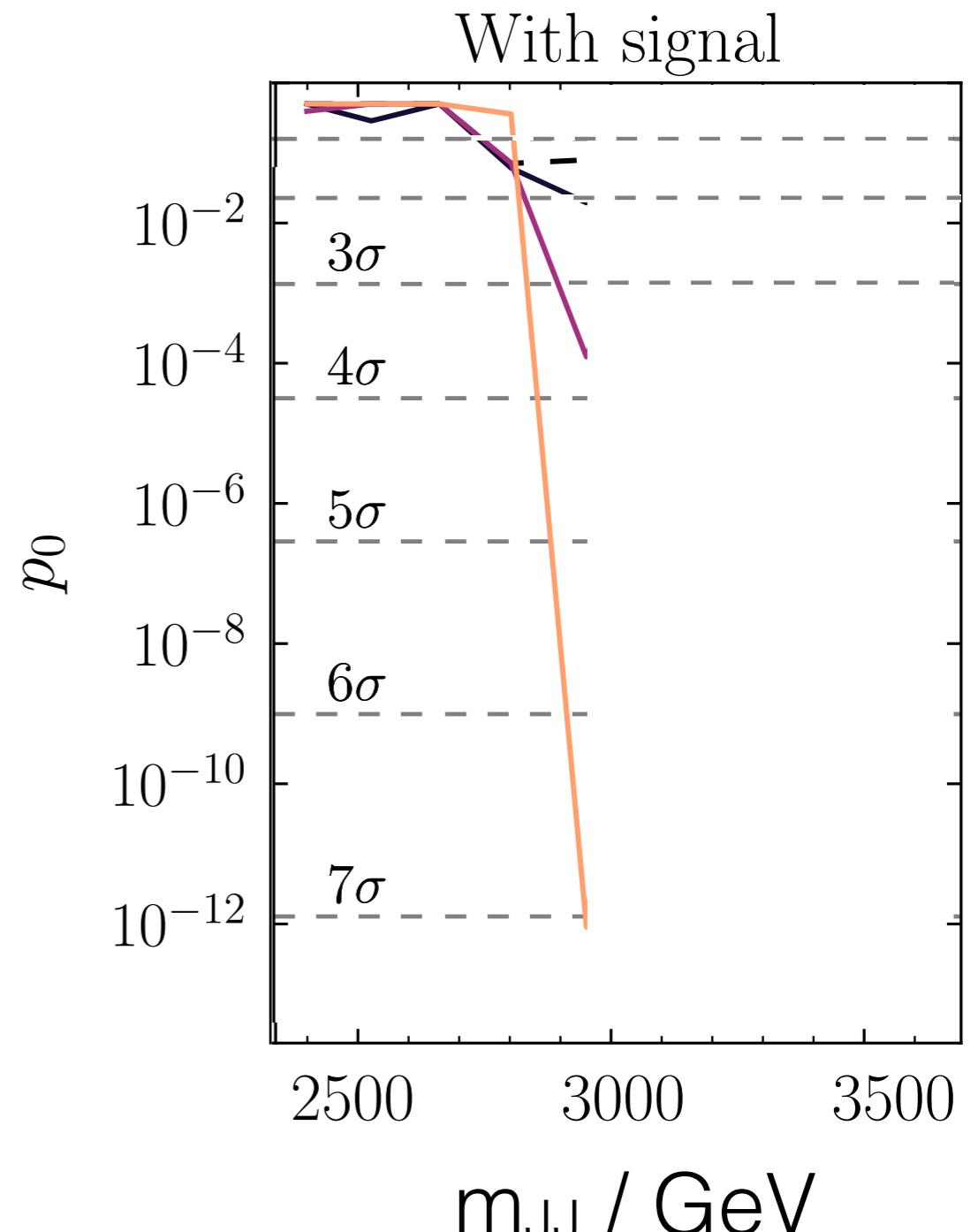
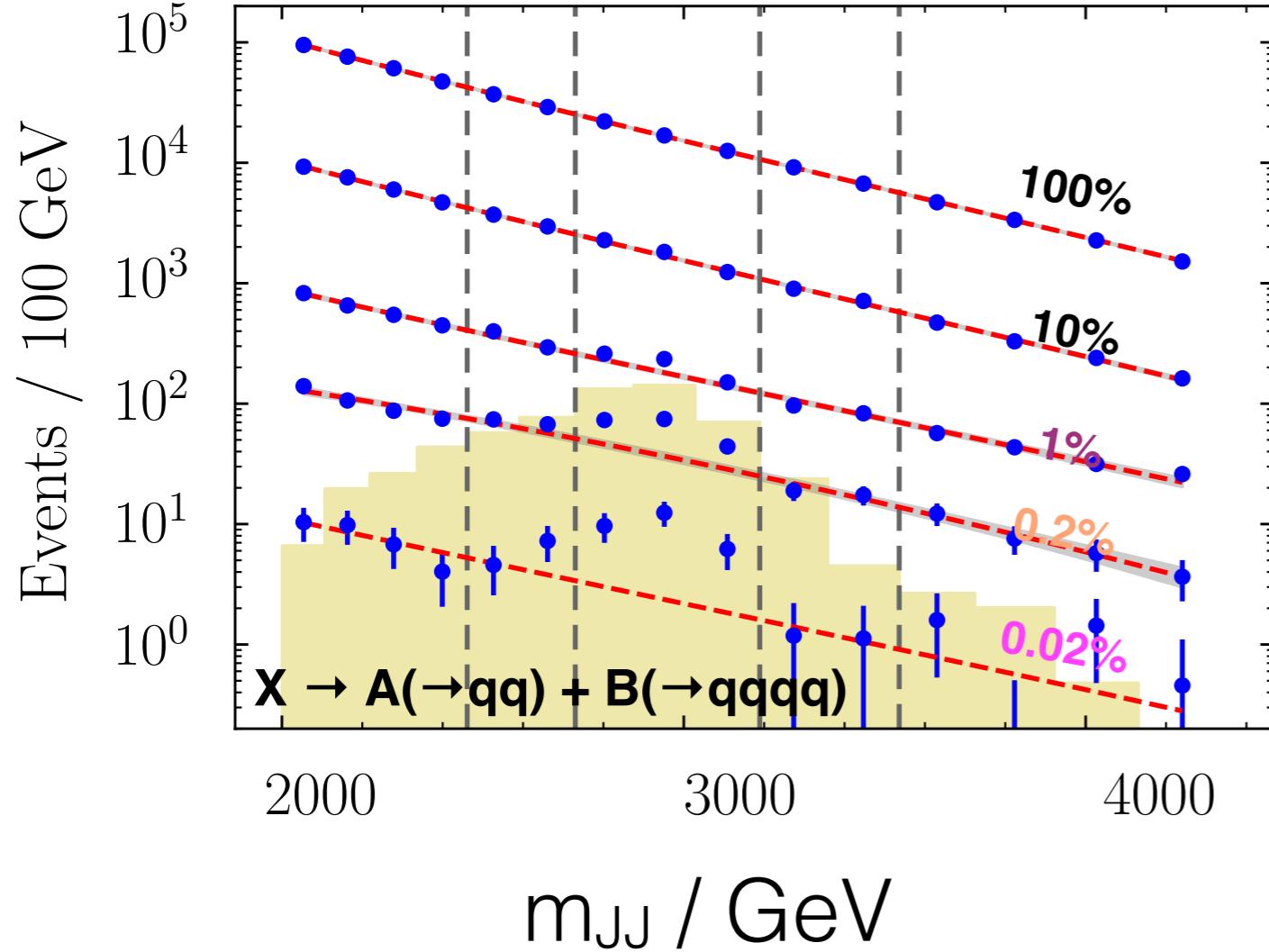


...and when there is a signal?



- no cut on NN
- most 10% signal-region-like
- most 1% signal-region-like
- most 0.2% signal-region-like

...and when there is a signal?



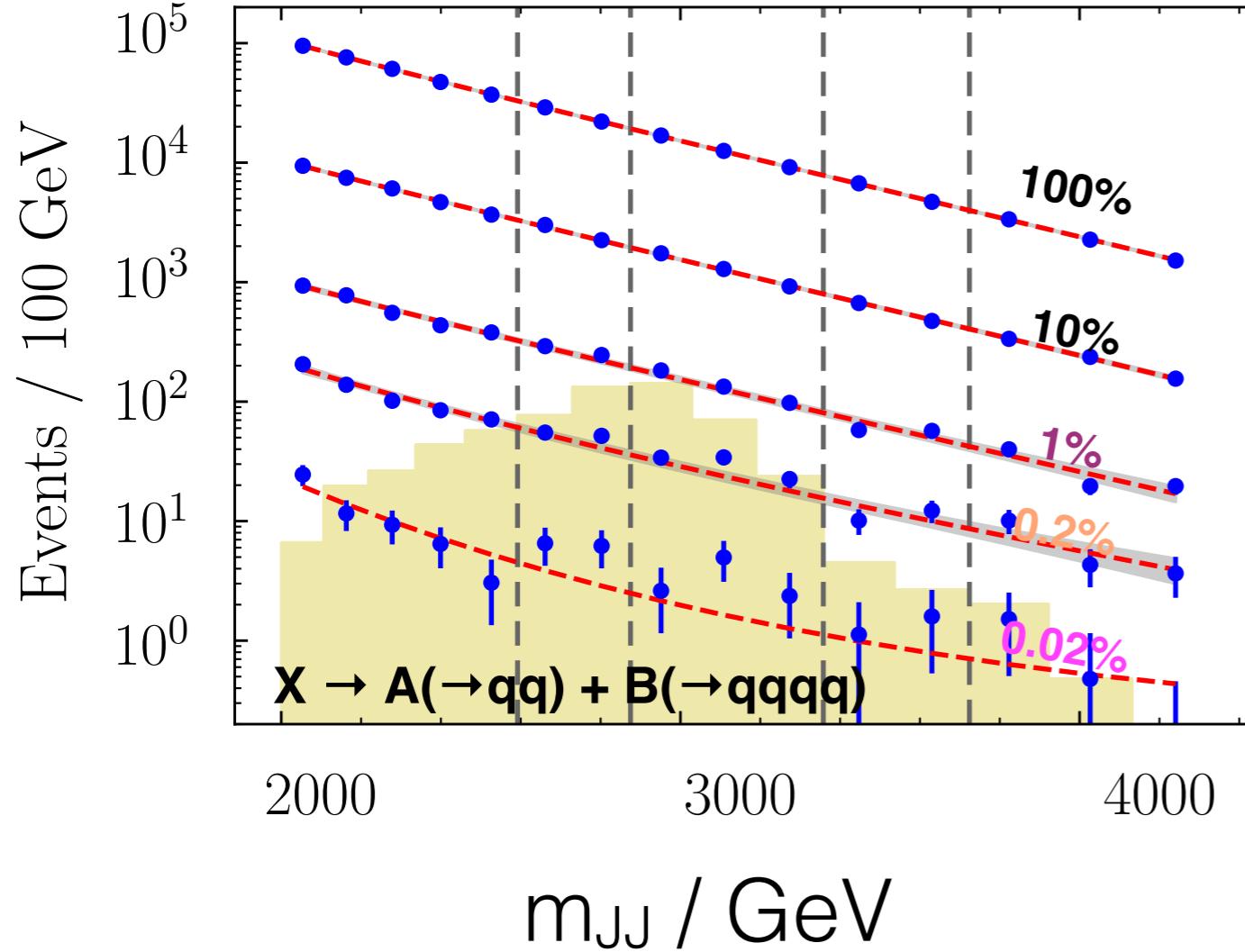
---- no cut on NN

— most 10% signal-region-like

— most 1% signal-region-like

— most 0.2% signal-region-like

...and when there is a signal?

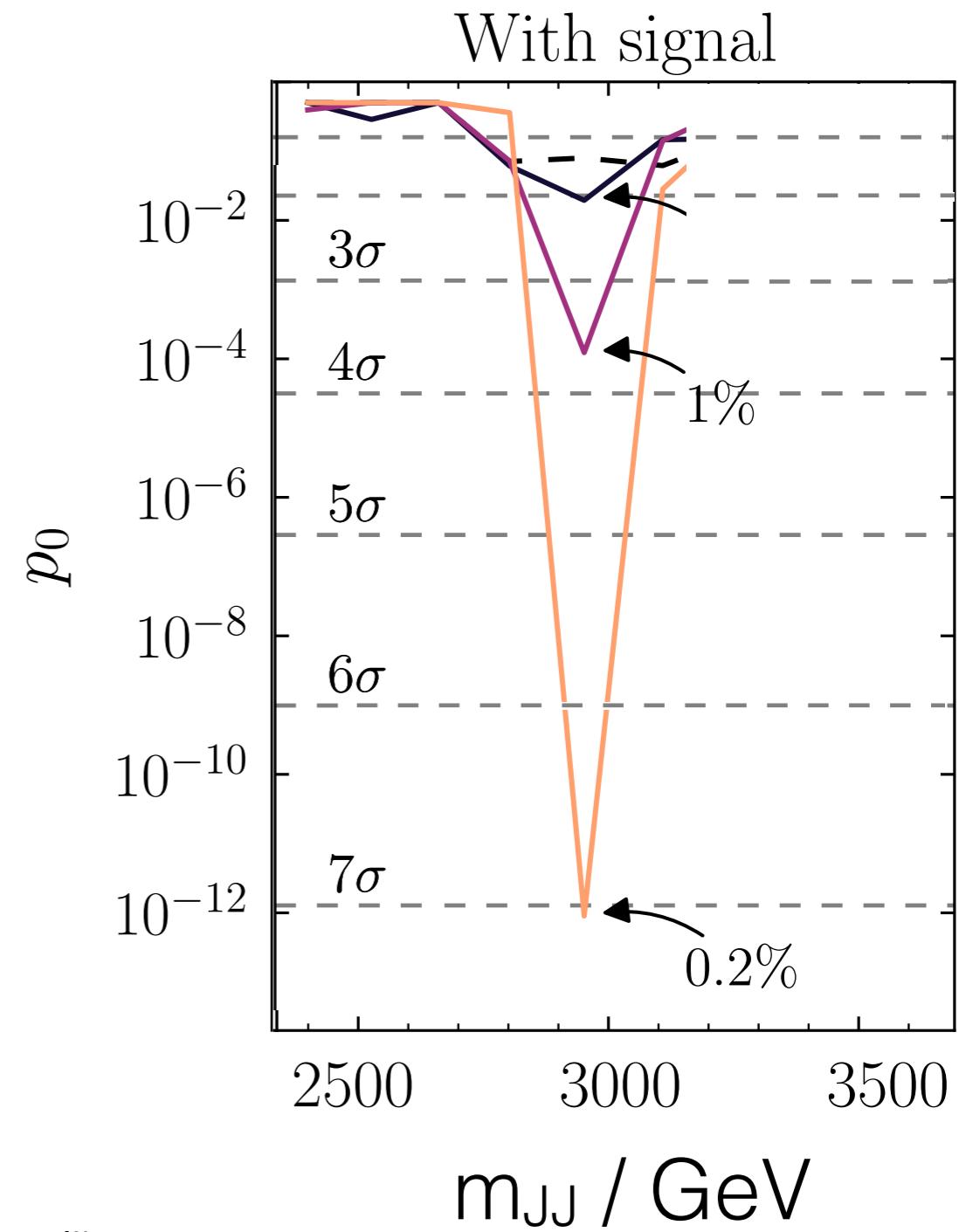


---- no cut on NN

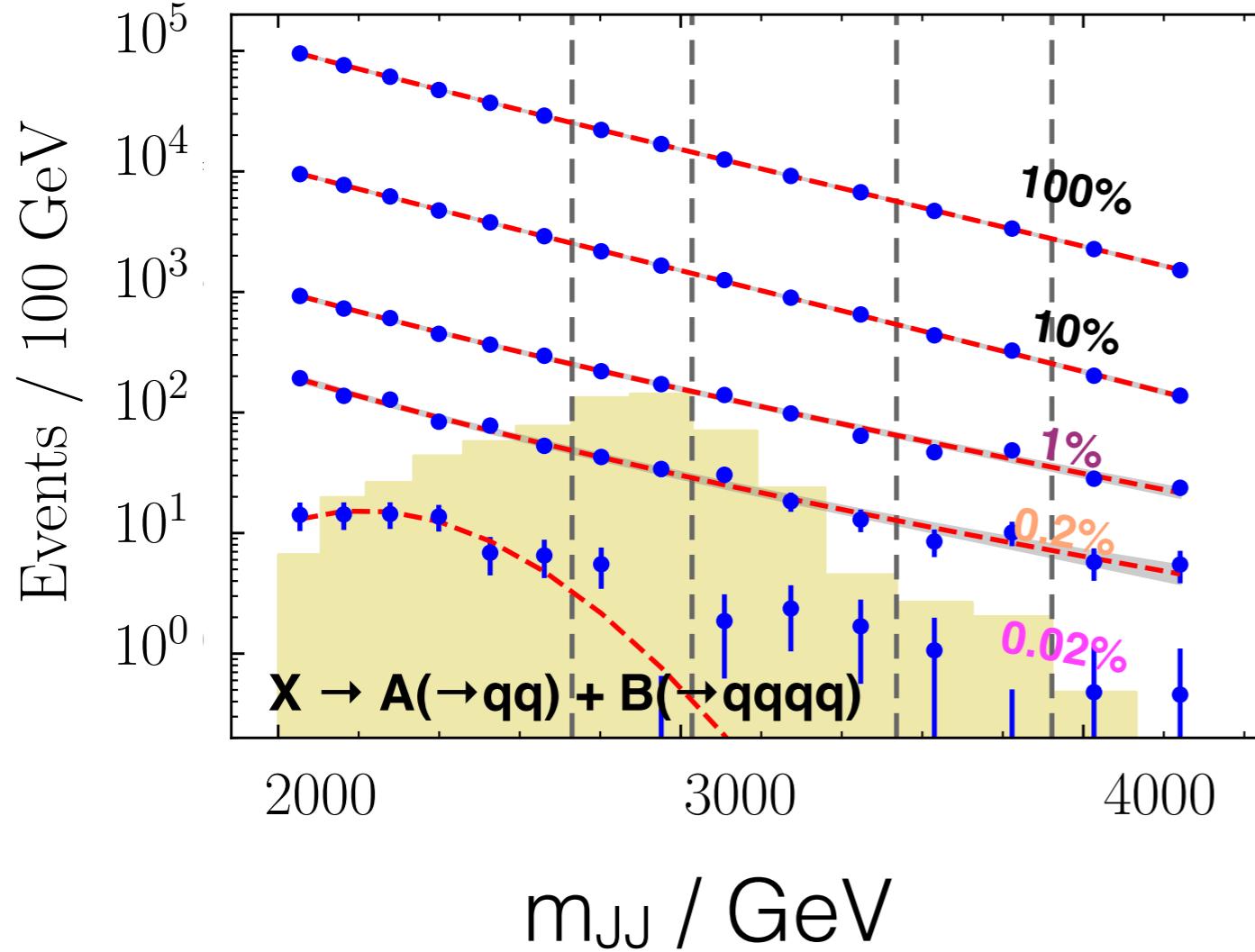
— most 10% signal-region-like

— most 1% signal-region-like

— most 0.2% signal-region-like



...and when there is a signal?

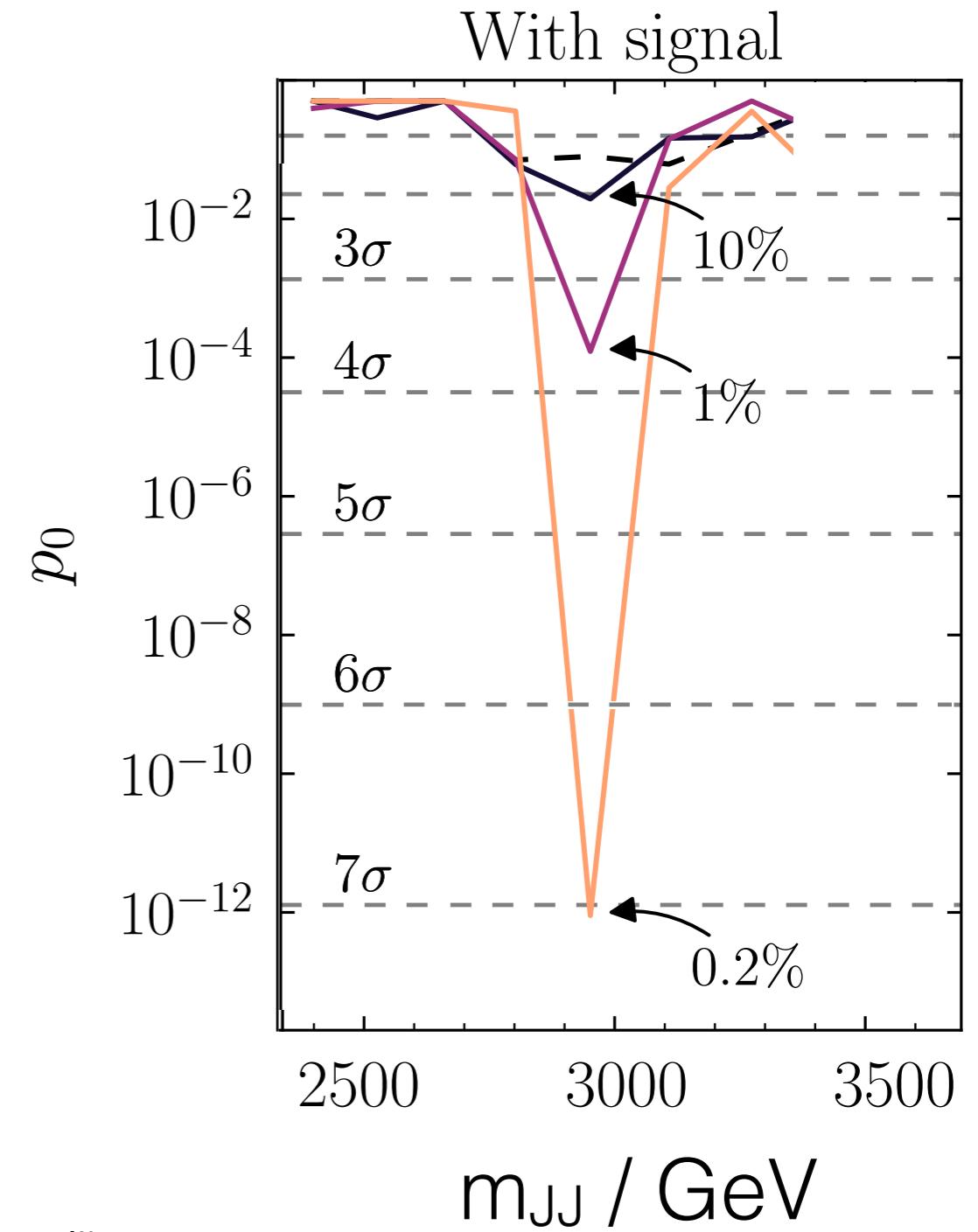


---- no cut on NN

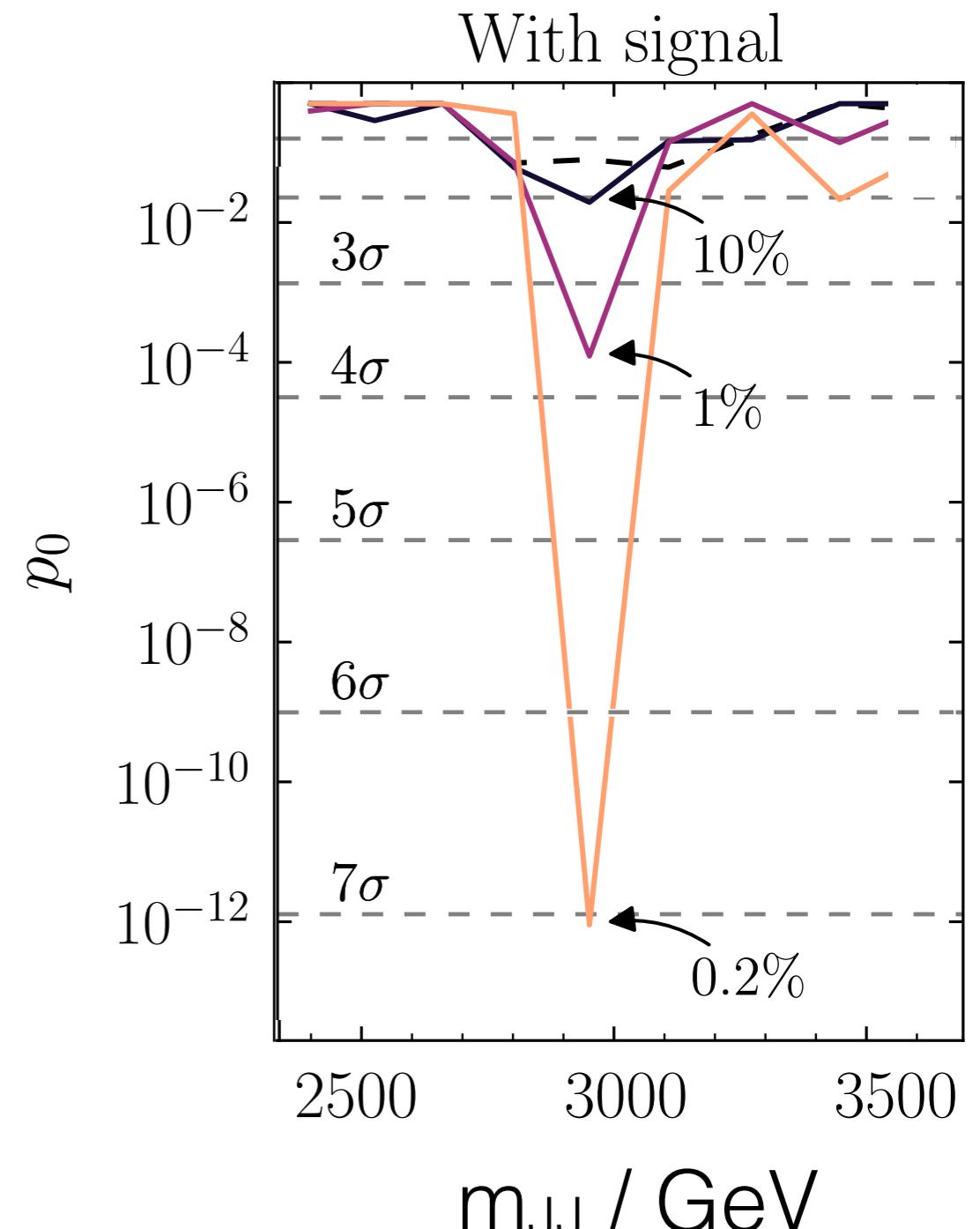
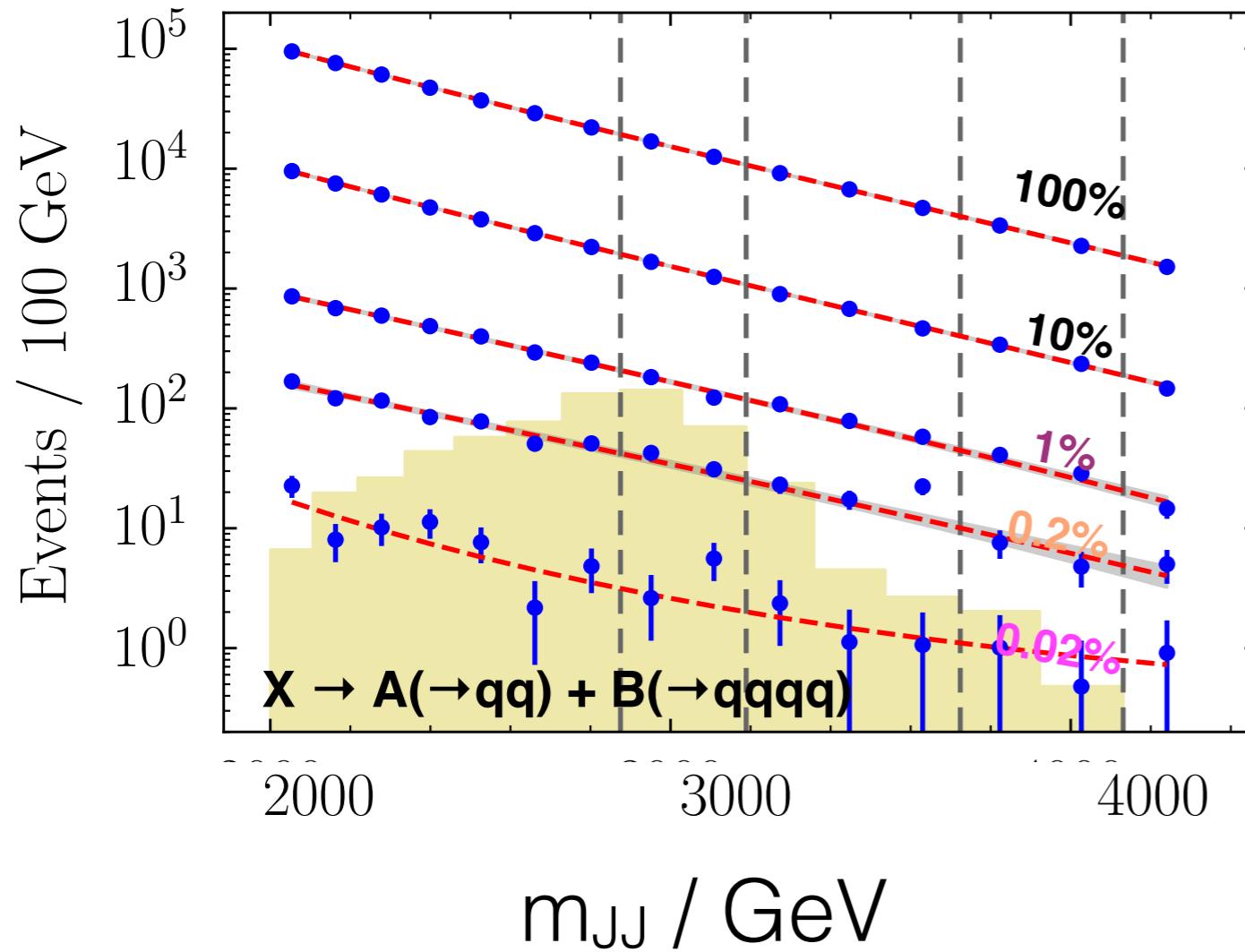
— most 10% signal-region-like

— most 1% signal-region-like

— most 0.2% signal-region-like



...and when there is a signal?



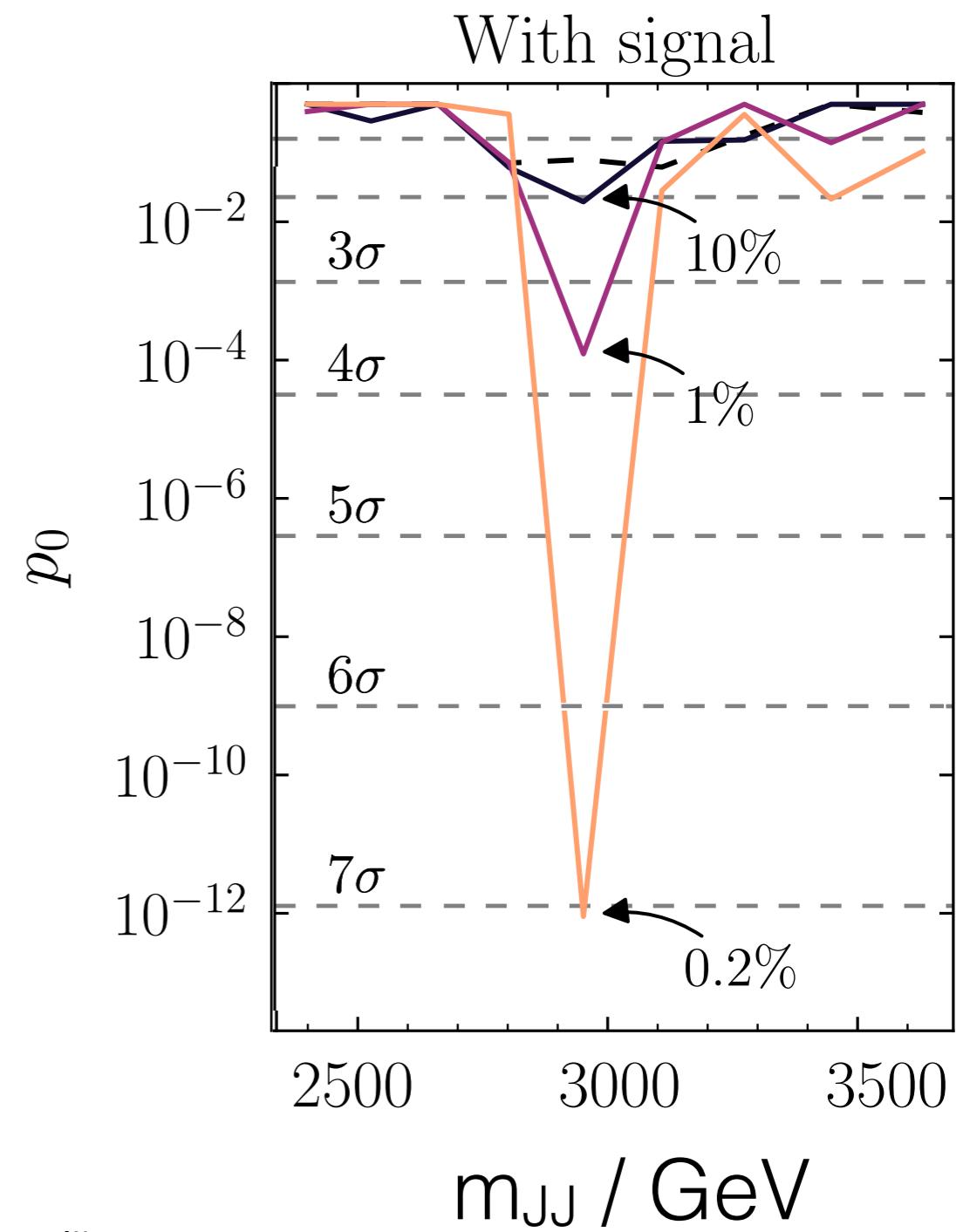
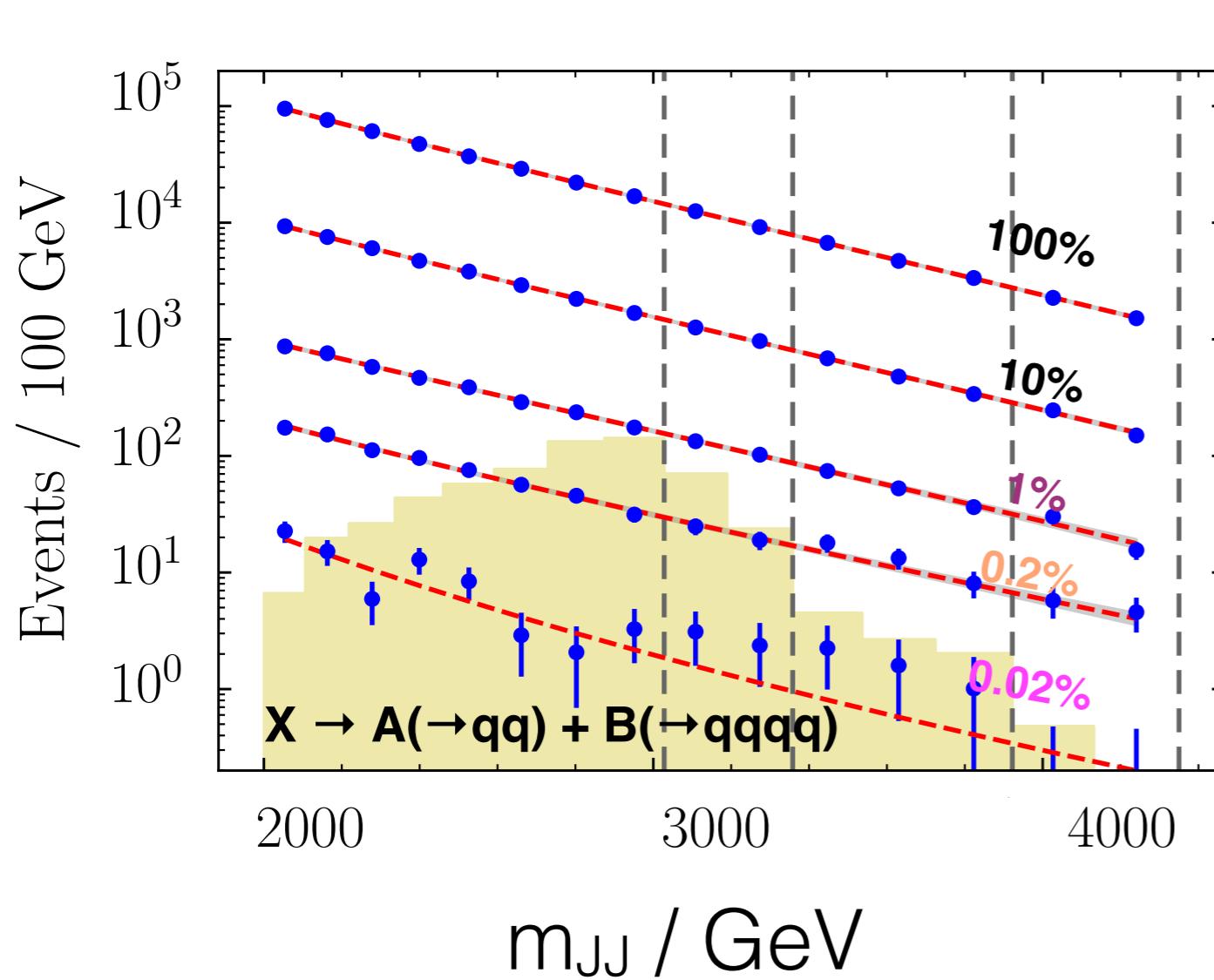
---- no cut on NN

— most 10% signal-region-like

— most 1% signal-region-like

— most 0.2% signal-region-like

...and when there is a signal?



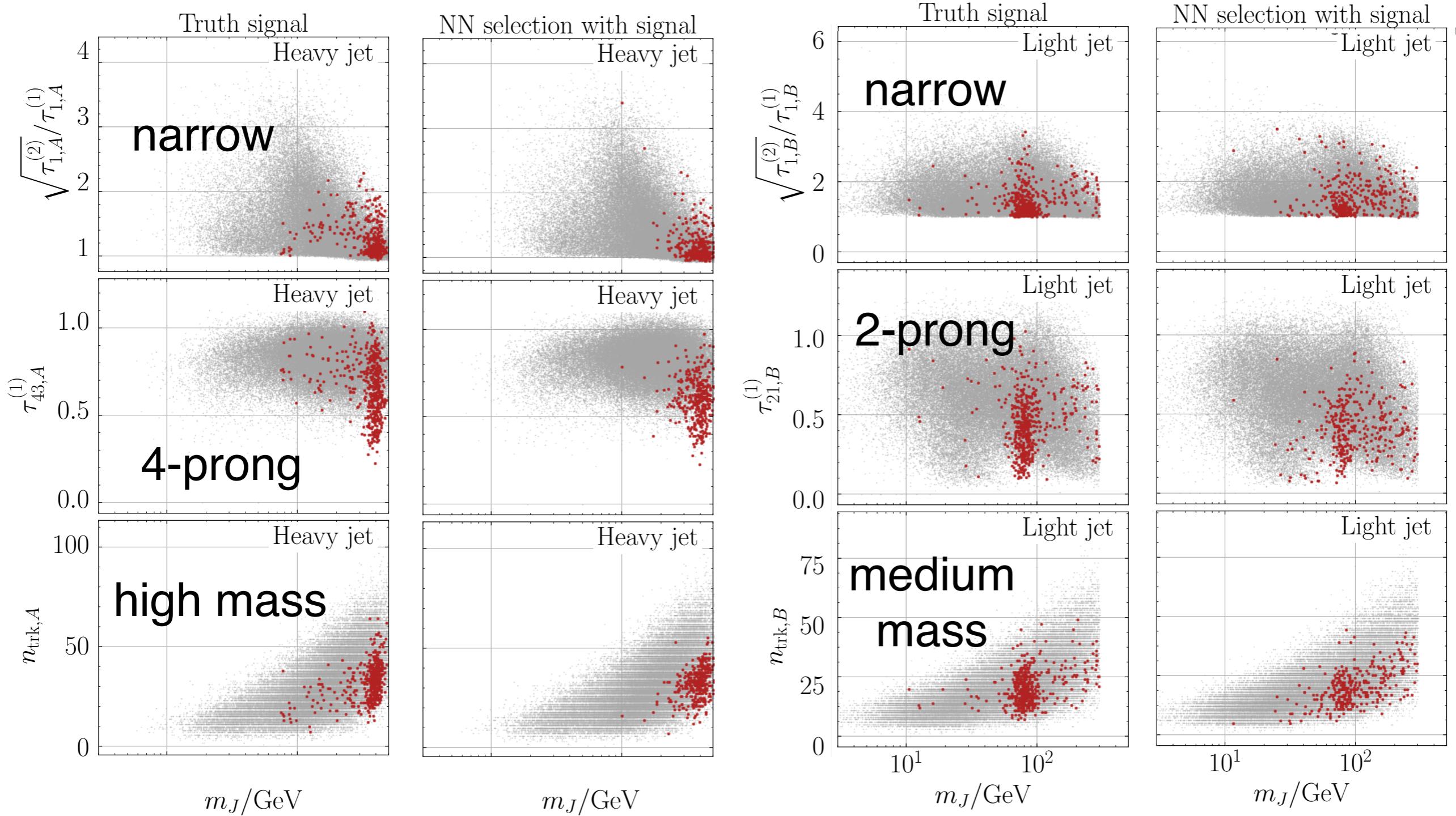
---- no cut on NN

— most 10% signal-region-like

— most 1% signal-region-like

— most 0.2% signal-region-like

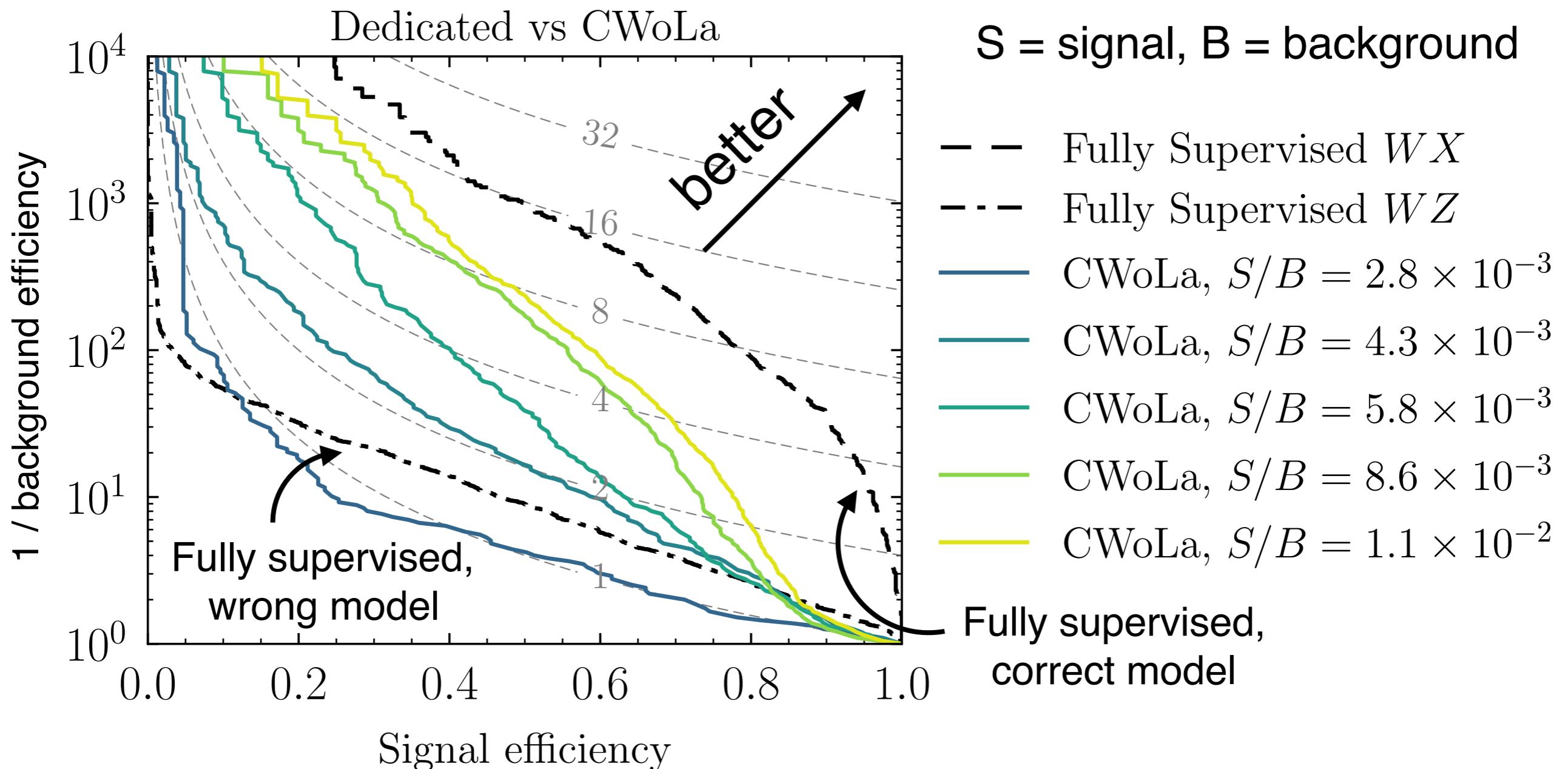
What is the network learning?



Learns to find the signal !

$X \rightarrow A(\rightarrow \text{qq}) + B(\rightarrow \text{qqqq})$

CWoLa hunting vs. Full Supervision

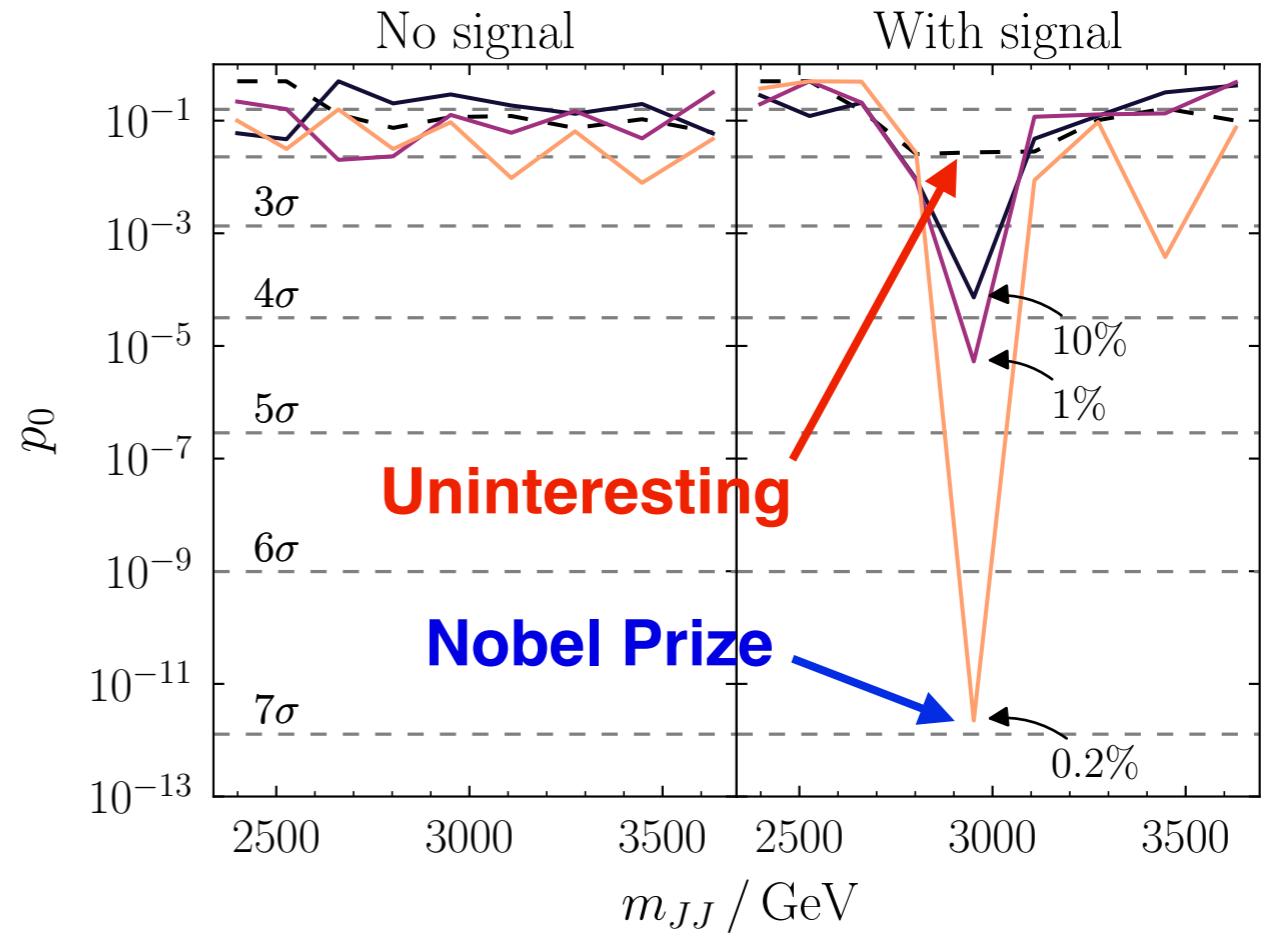


If you know what you are looking for, you should look for it. If you don't know, then CWoLa hunting may be able to catch it!

Outlook*

Jack Collins, Kiel Howe, **BPN**, PRL 121, 241803 (2018)

If we continue to find nothing (and even if we find something!), it is important to broaden our scope.



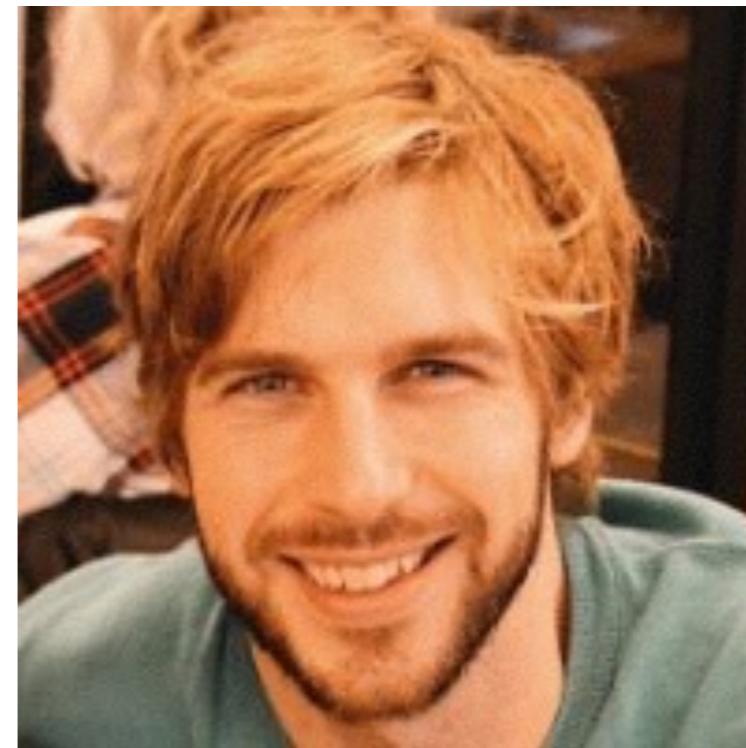
There may be BSM hiding, but we need hyper-variate vision to find it ... ML can help us learn something fundamental about nature !

*Image from [this article](#). This Koala is actually being freed - I do not condone violence against these animals!

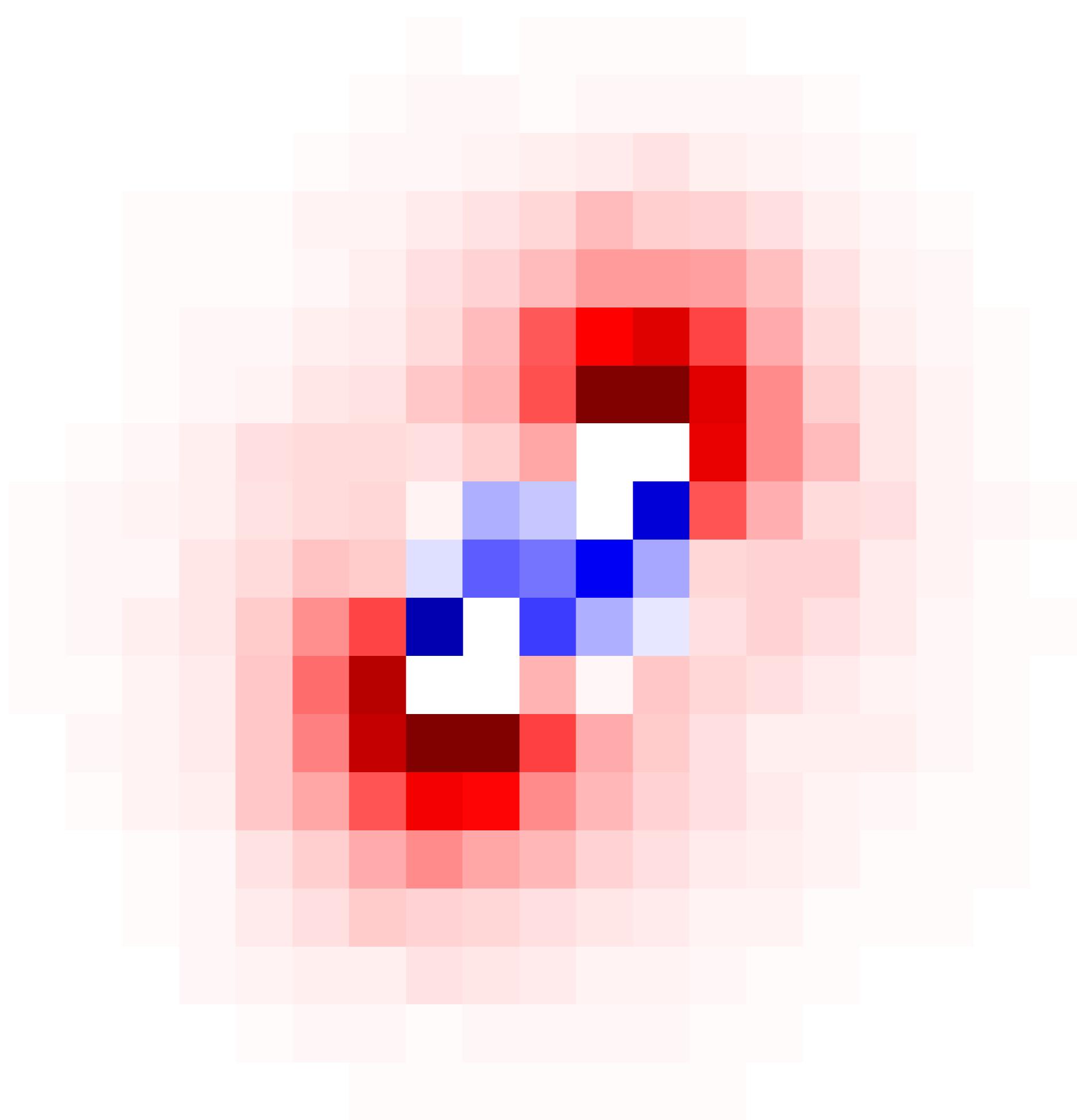
**...many thanks to a fun and insightful
collaboration with these guys!**



Jack Collins



Kiel Howe



Fin.

Autoencoders

Marco Farina, Yuichiro Nakai, David Shih,
<https://arxiv.org/abs/1808.08992>

Theo Heimel, Gregor Kasieczka, Tilman Plehn,
Jennifer M Thompson, <https://arxiv.org/abs/1808.08979>

Olmo Cerri, Thong Q. Nguyen, Maurizio Pierini, Maria Spiropulu,
Jean-Roch Vlimant, <https://arxiv.org/abs/1811.10276>

Statistics

