



# A big “class” on jet physics

*(Learning rich jet representation via  $o(100)$  classes to accelerate the LHC resonance search programme)*

Congqiao Li (李聪乔), *Peking University*

***This talk is mainly based on***

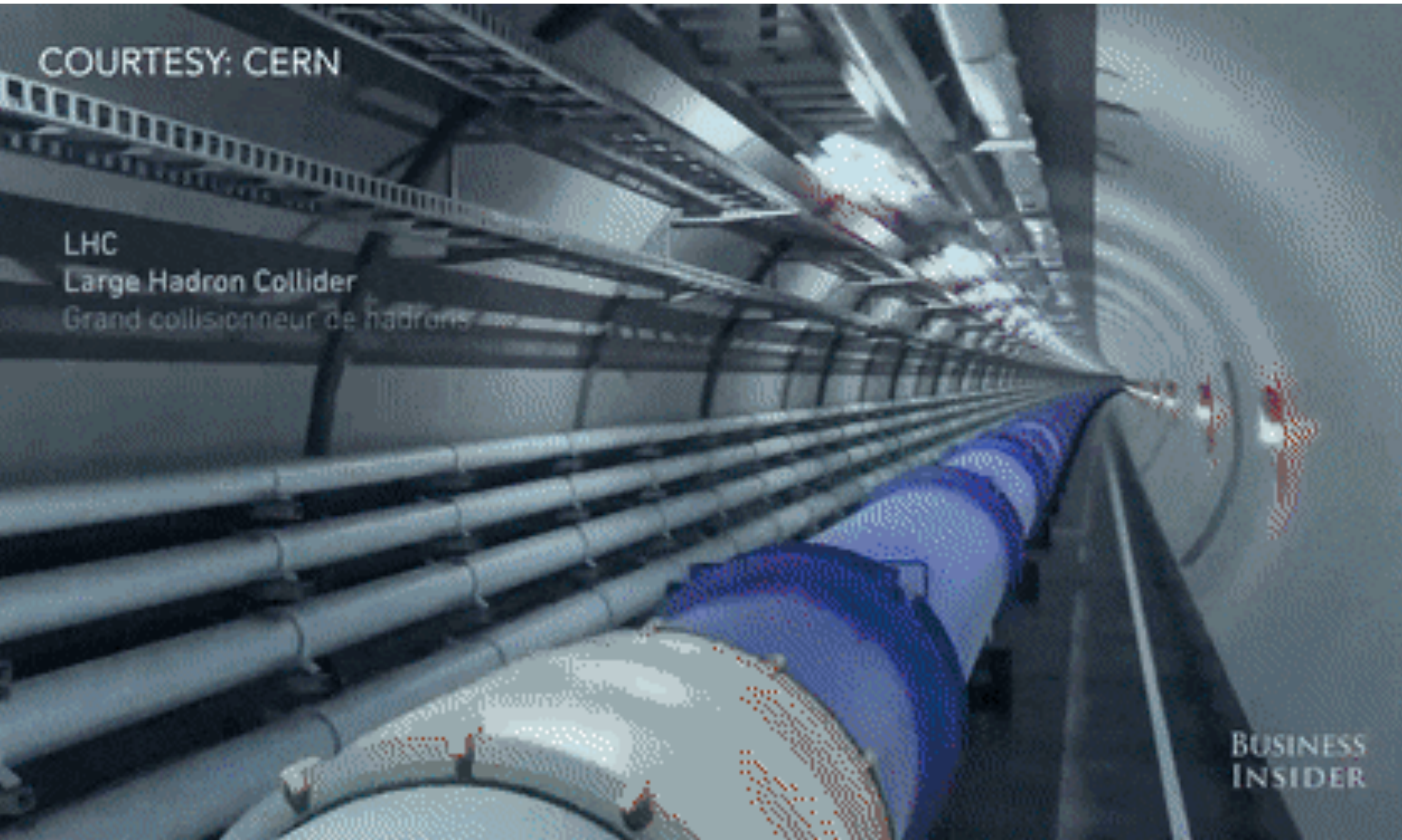
- 🪐 Accelerating LHC resonant search via *Sophon*: [arXiv:2405.12972](https://arxiv.org/abs/2405.12972)  
[\[Github\]](#) [\[Dataset\]](#) [\[Model\]](#) [\[Google Colab\]](#)
- 🪐 Development of *Global Particle Transformer (GloParT)* 3 within CMS

Larger than Larger: Large AI Models at the Frontiers of  
Experimental High-Energy Physics, Beijing  
7 January, 2025

***Let's begin with stories ...***



# The Large Hadron Collider



COURTESY: CERN

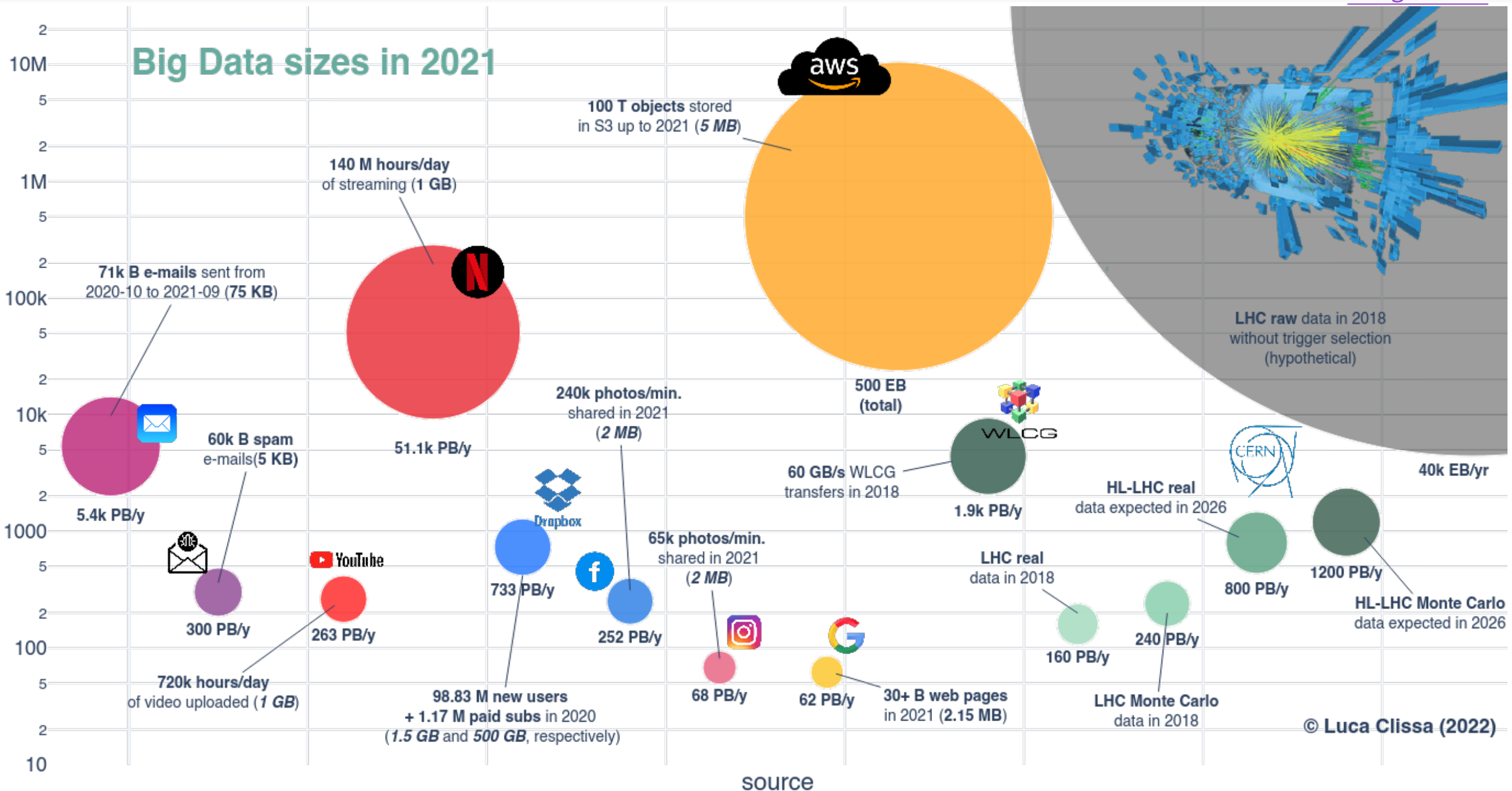
LHC  
Large Hadron Collider  
Grand collisionneur de hadrons

BUSINESS  
INSIDER



# In the age of big data...

Image credit



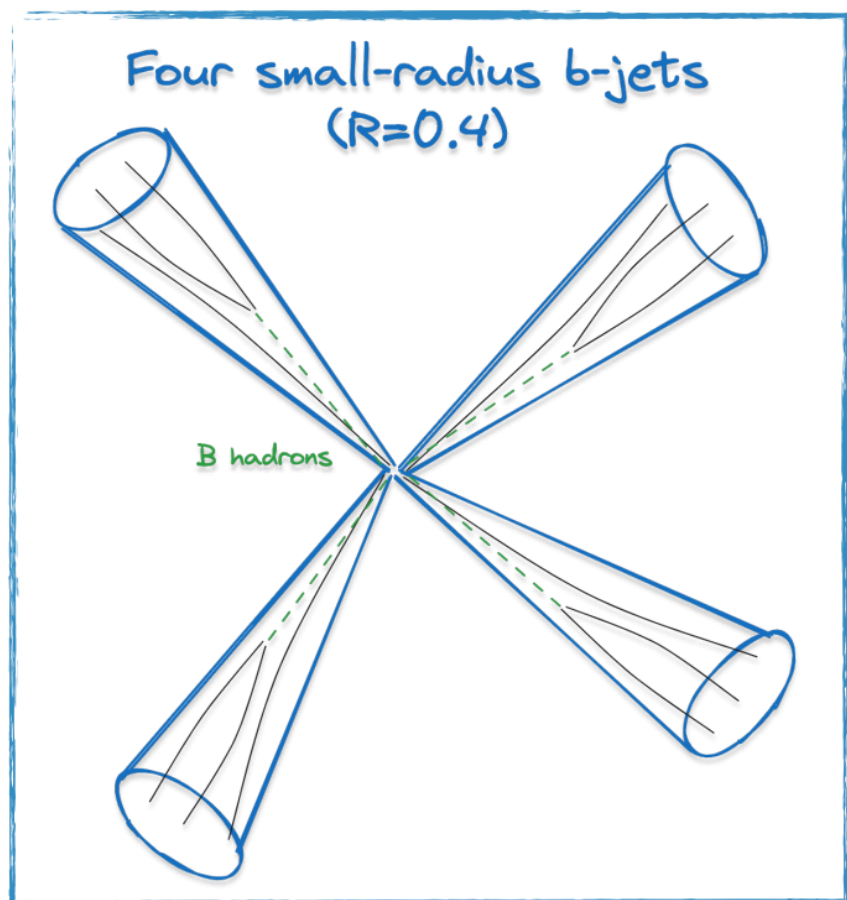
- Data analysis is a crucial subject in the particle physics of this age
- Novel engineering solutions are appearing!
- The advent of deep learning / AI is expected to play a transformative role



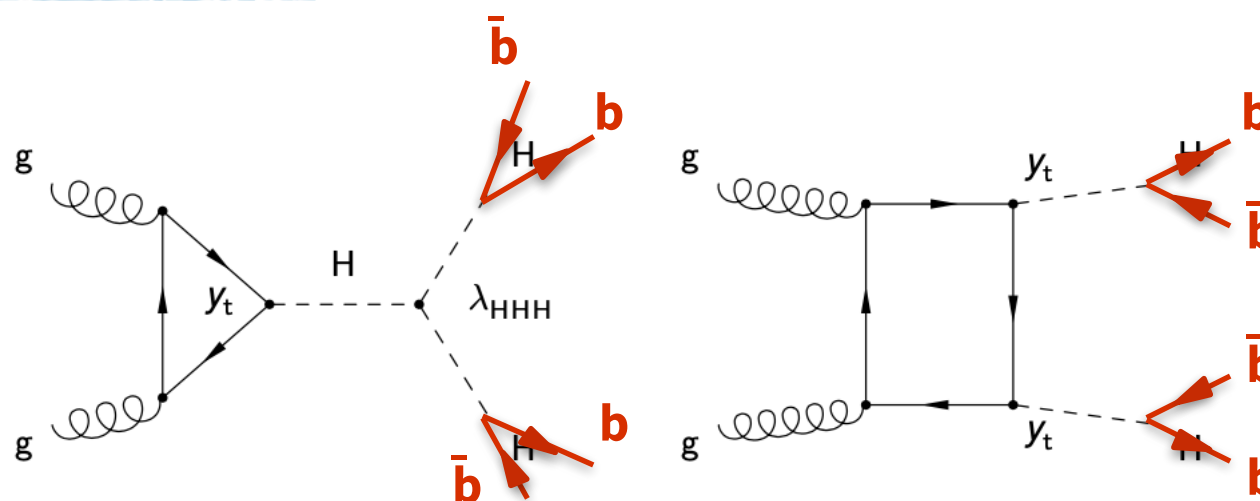
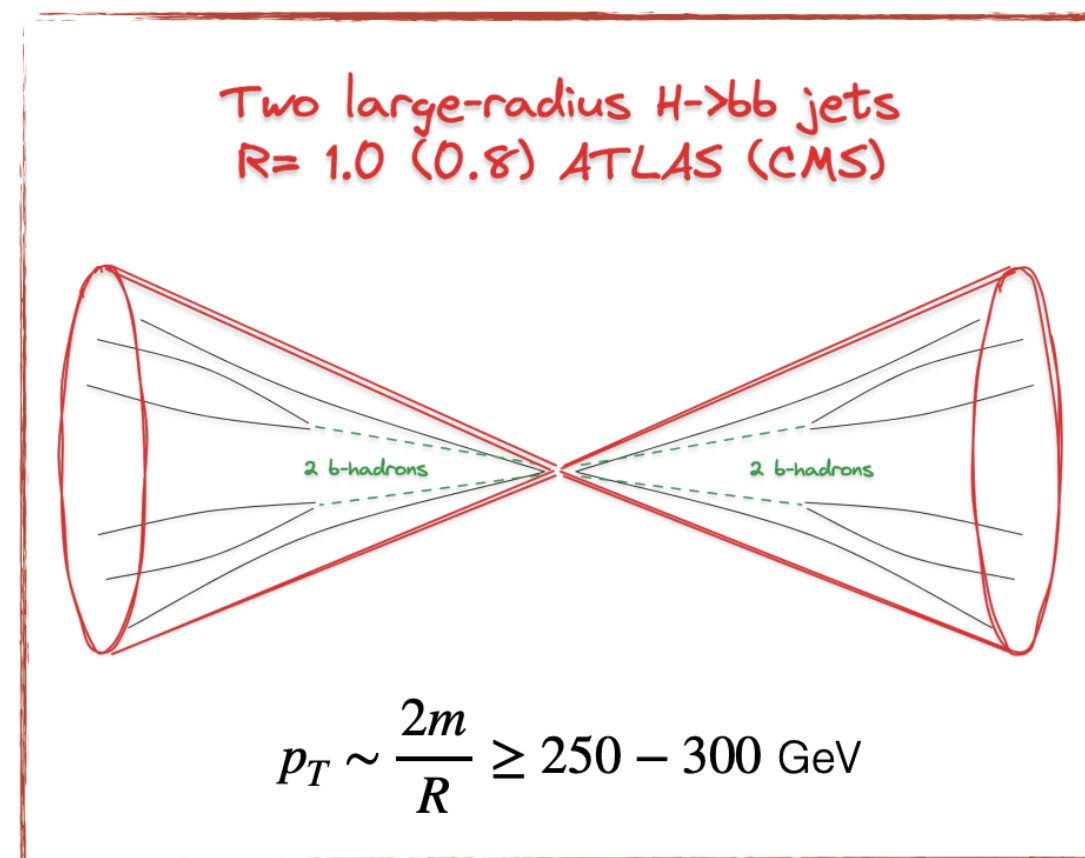
# Quick example: how to search for $HH \rightarrow 4b$

cartoons credit to [link](#)

Resolved regime



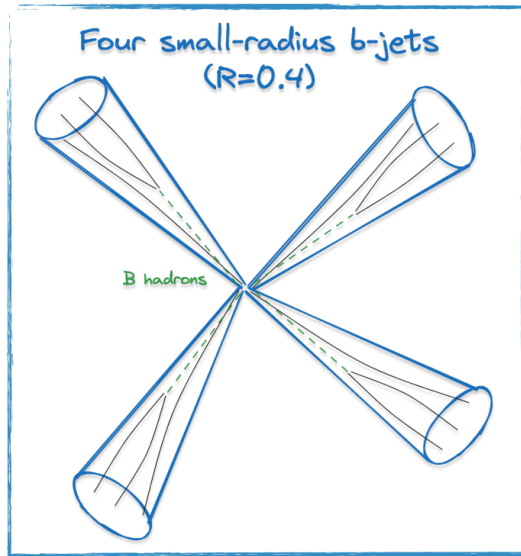
Boosted regime



To search for  $HH \rightarrow 4b$  signal at the LHC?

# Boosted regime as a booster?

## Resolved regime

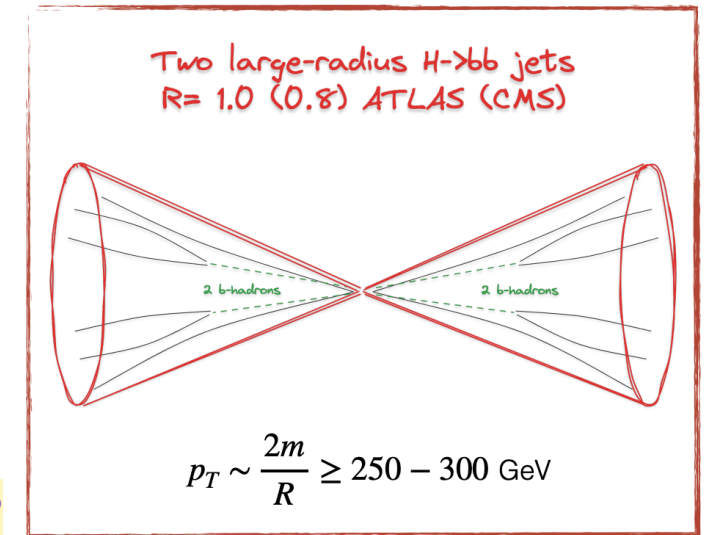


[PRL 129 \(2022\) 081802](#)

[results website](#)

- Boosted regime can only focus on high-pT behaviour
  - can only detect a tiny amount of signals ( $m_{HH} \gtrsim 600$  GeV)
  - but their sensitivity is close!

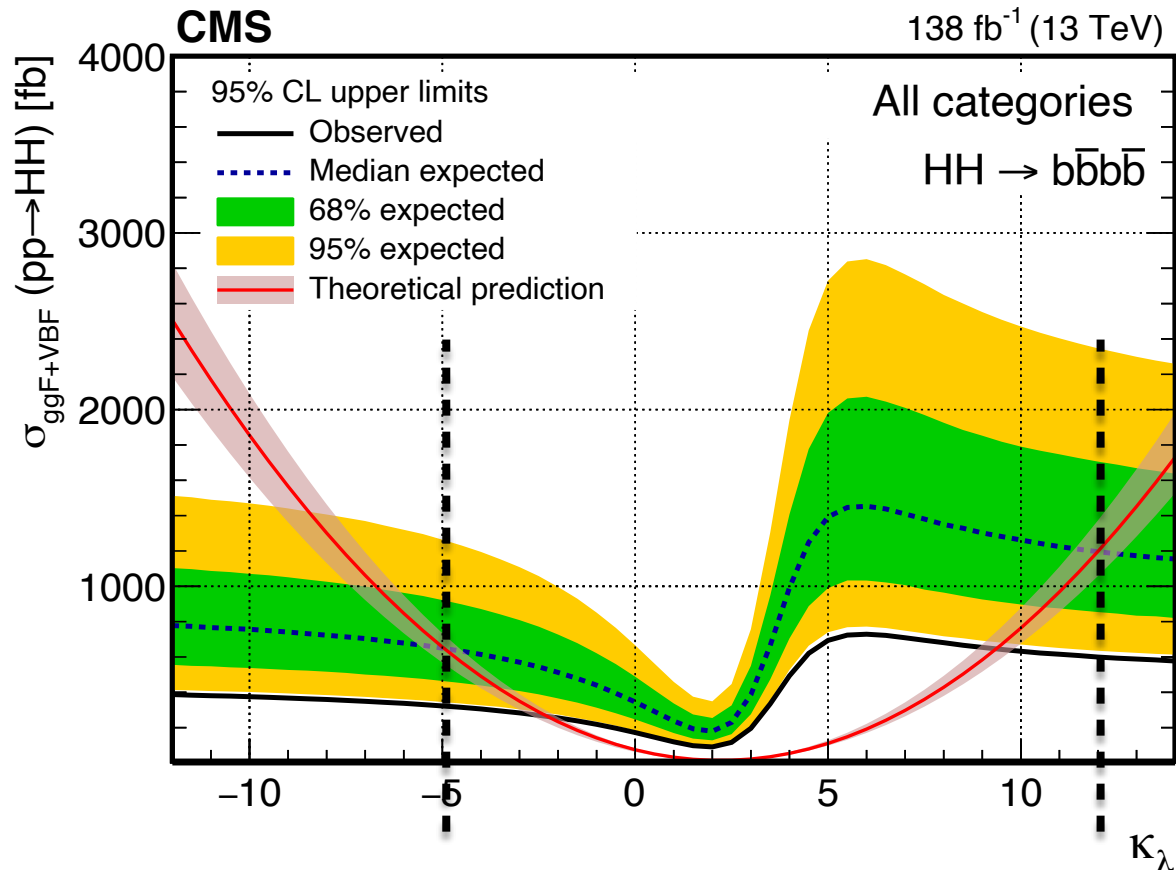
## Boosted regime



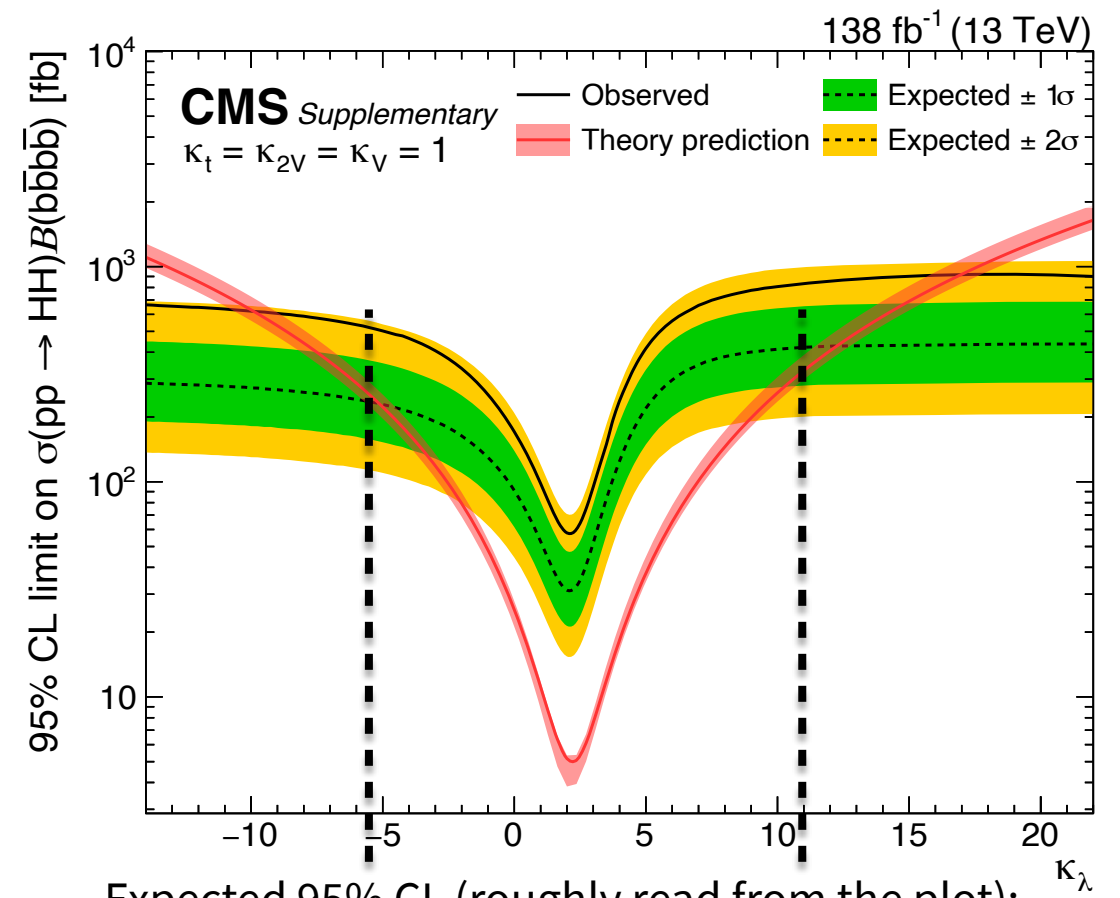
[PRL 131 \(2023\) 041803](#)

[results website](#)

$$p_T \sim \frac{2m}{R} \geq 250 - 300 \text{ GeV}$$



Exp (obs) 95% CL:  $-5.0 (-2.3) < \kappa_\lambda < 12.0 (9.4)$

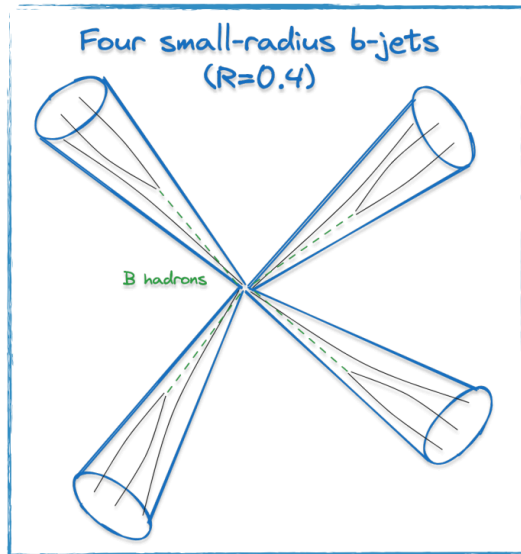


Expected 95% CL (roughly read from the plot):  
 $-5.0 < \kappa_\lambda < 12.5$



# Boosted regime as a booster?

## Resolved regime

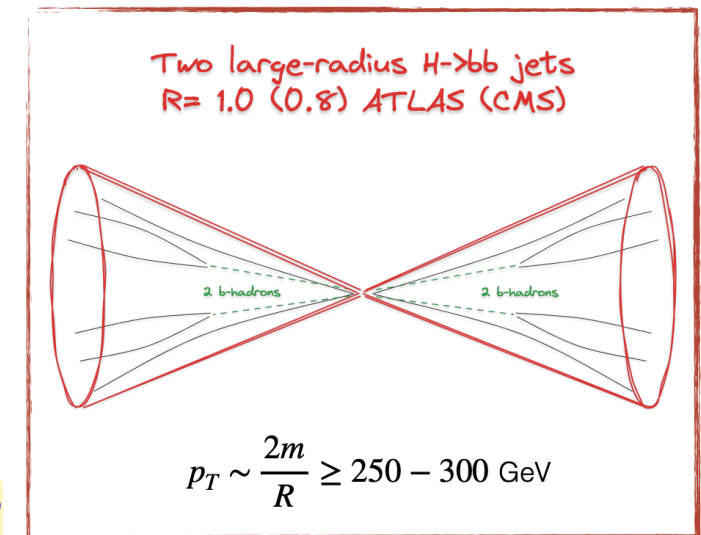


- Boosted regime can only focus on high-pT behaviour
  - can only detect a tiny amount of signals ( $m_{HH} \gtrsim 600$  GeV)
- but their sensitivity is close!

[PRL 129 \(2022\) 081802](#)

[results website](#)

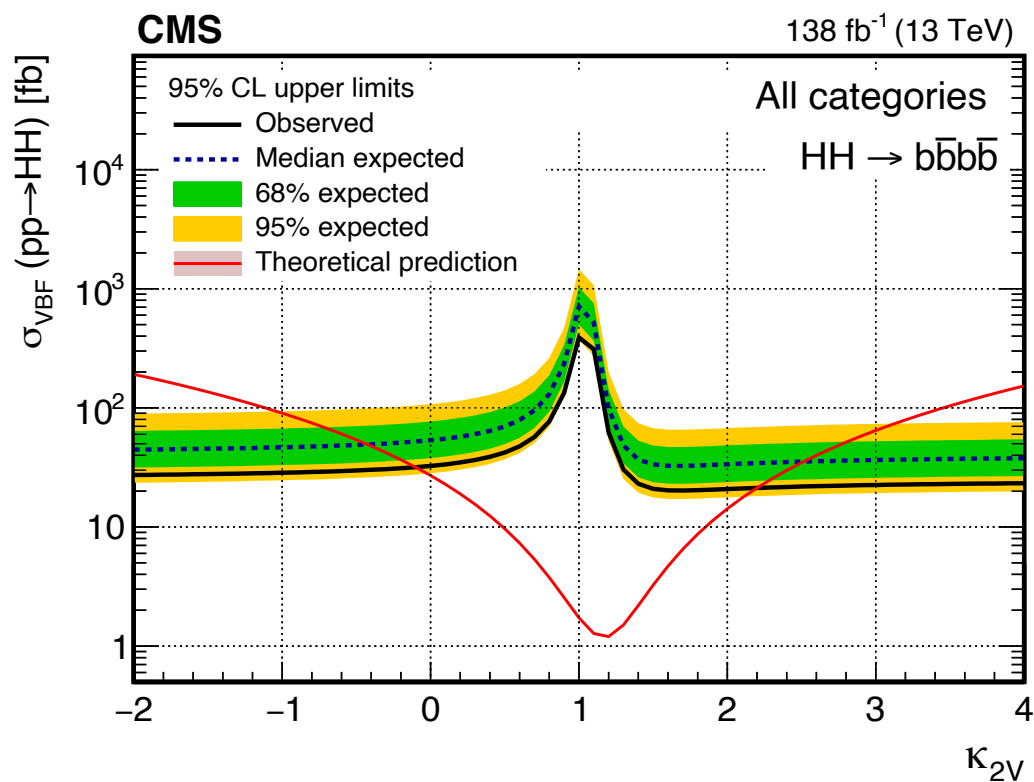
## Boosted regime



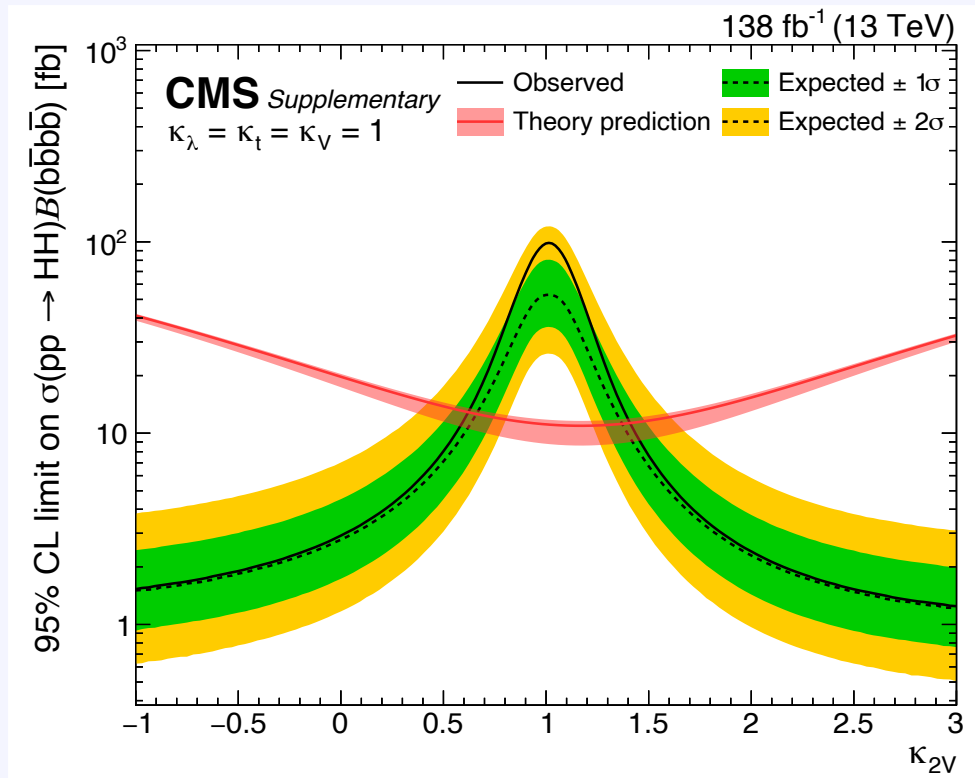
$$p_T \sim \frac{2m}{R} \geq 250 - 300 \text{ GeV}$$

[PRL 131 \(2023\) 041803](#)

[results website](#)



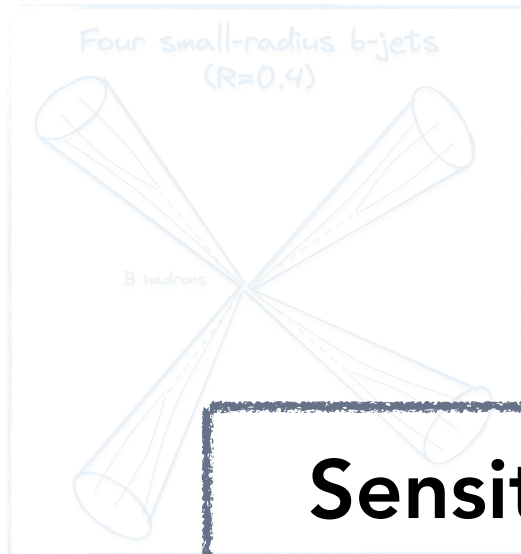
Exp (obs) 95% CL:  $-0.1 (-0.4) < \kappa_{2V} < 2.2 (2.5)$



Exp (obs) 95% CL:  $0.66 (0.64) < \kappa_{2V} < 1.37 (1.41)$

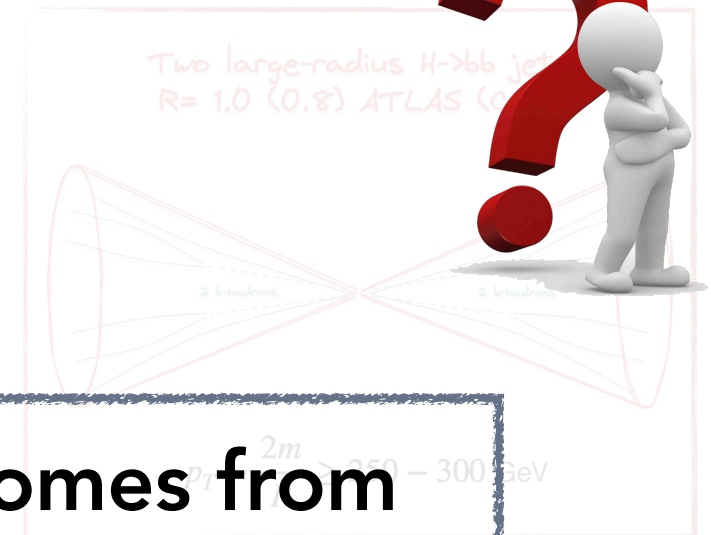
# Boosted regime as a booster?

Resolved regime



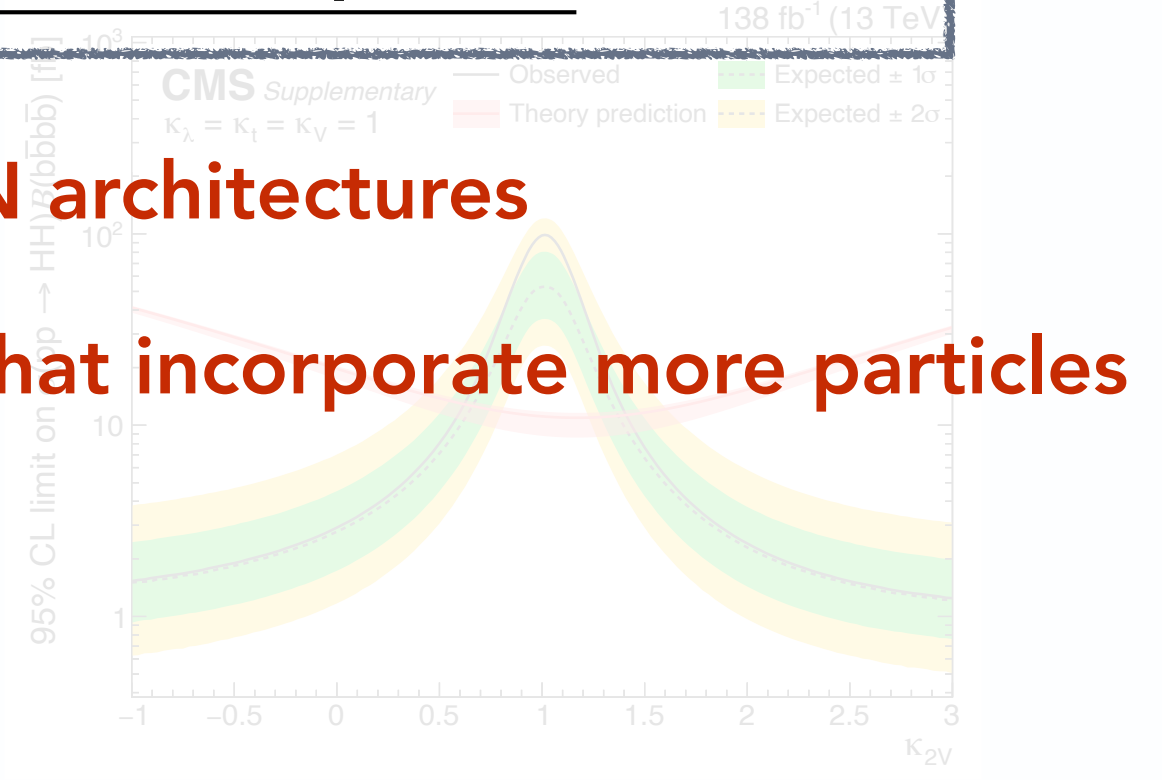
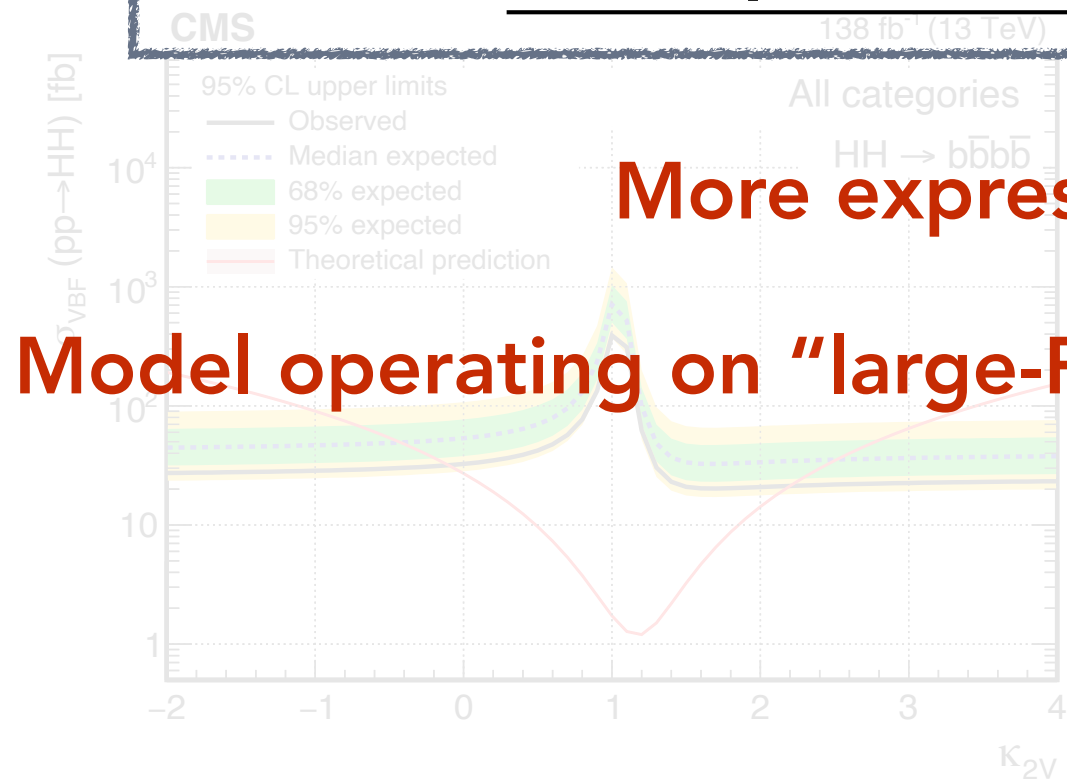
- Boosted regime can only focus on high-pT behaviour
  - can only detect a tiny amount of signals ( $m_{HH} \gtrsim 600$  GeV)
  - but their sensitivity is close!

Boosted regime



And even better results on  $\kappa_{2V}$  (where boosted regime is posed to give better sensitivity)

**Sensitivity in boosted topology largely comes from the superior BKG suppression power**



More expressive NN architectures

Model operating on "large-R jets" that incorporate more particles

Exp (obs) 95% CL:  $-0.1 (-0.4) < \kappa_{2V} < 2.2 (2.5)$

Exp (obs) 95% CL:  $0.66 (0.64) < \kappa_{2V} < 1.37 (1.41)$



# Outline

## → I. Backgrounds

- ❖ overview of boosted-jet taggers at the LHC
- ❖ what's next?

## → II. Introducing Sophon

- ❖ “large model for large-scale classification”; how are we led there?
- ❖ Sophon details & performance benchmark

## → III. Implications for LHC resonance search

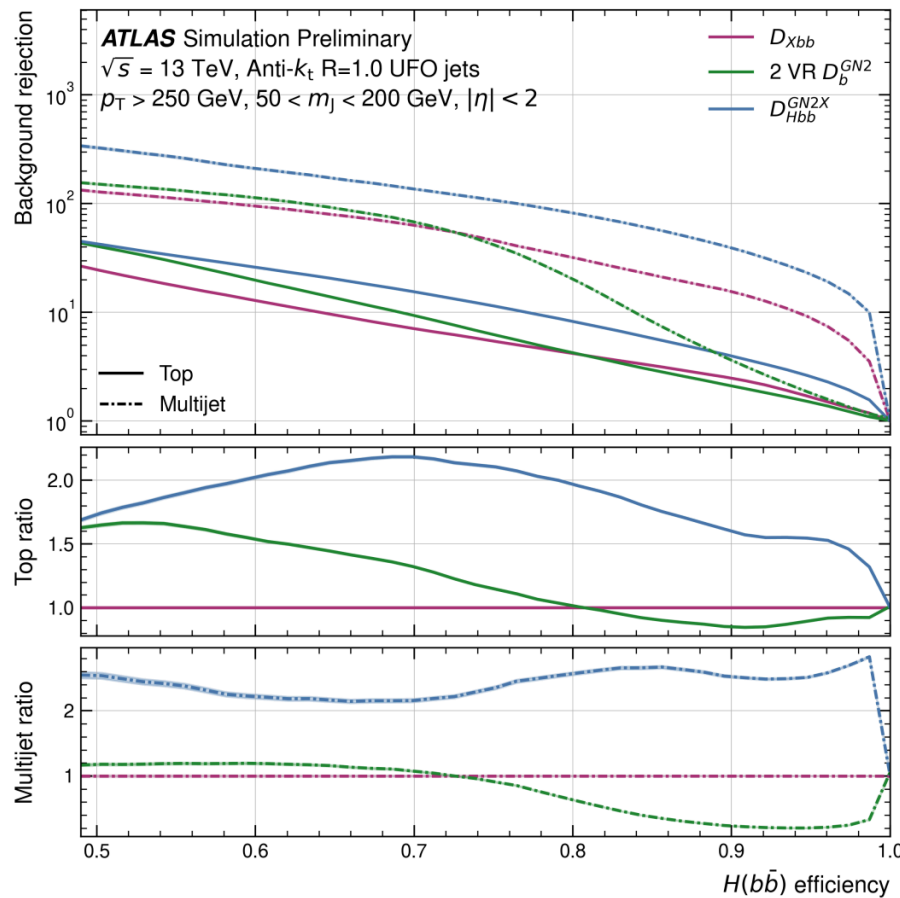
- ❖ model-specific and model-agnostic approaches
- ❖ *Global Particle Transformer*, the next-generation model in CMS
- ❖ More opportunities by Sophon/GloParT

## → IV. Open discussion

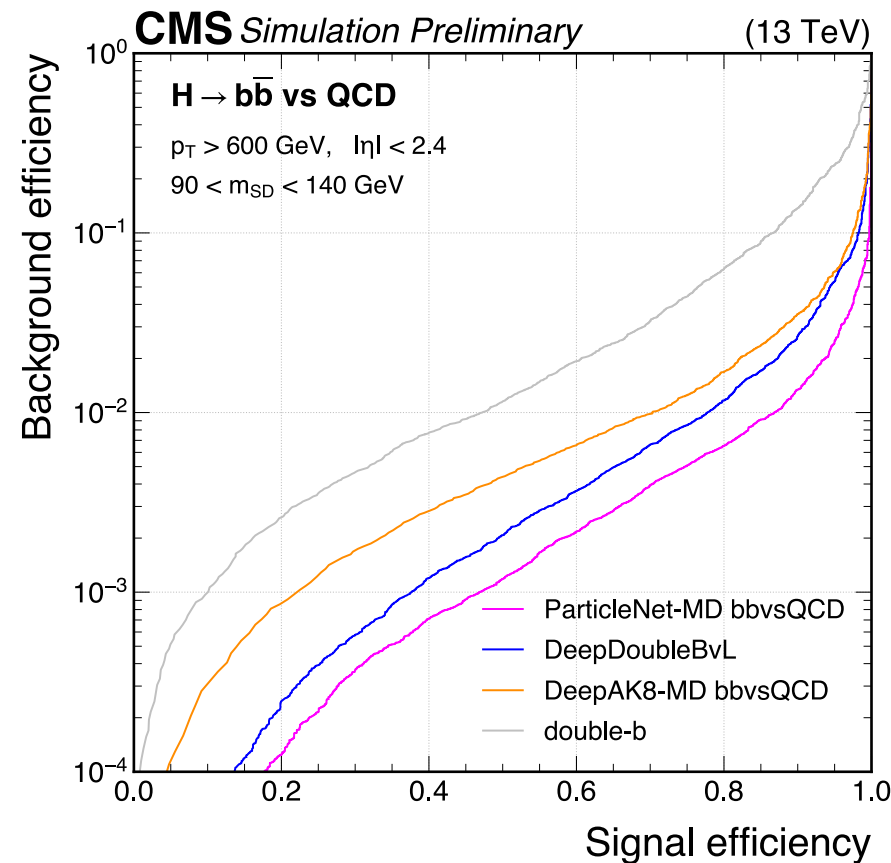
- ❖ datasets, training targets and scaling capabilities

# Inspiring progress on $H \rightarrow b\bar{b}/c\bar{c}$ tagging

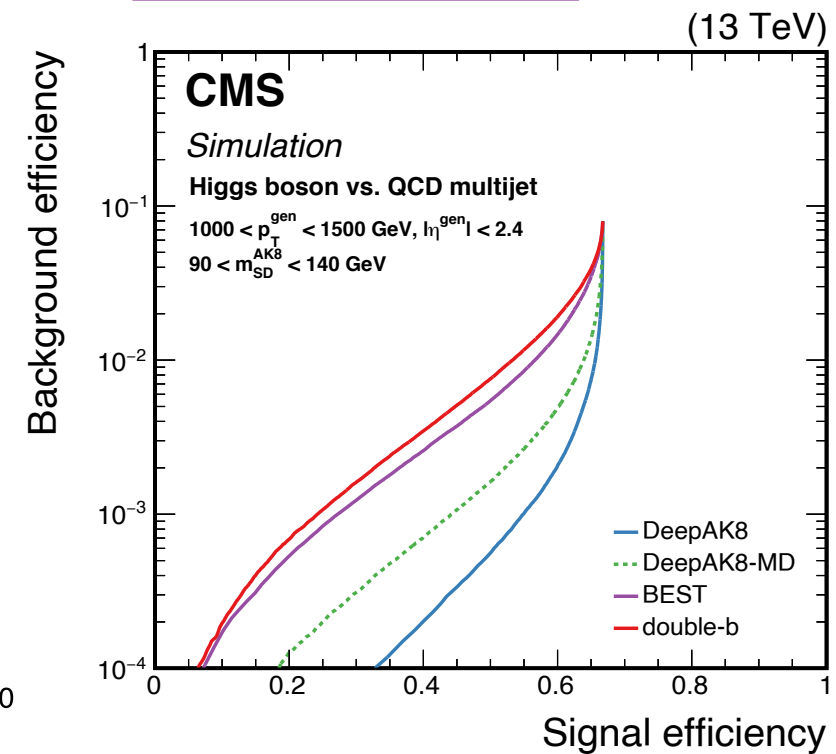
ATL-PHYS-PUB-2023-021



CMS-PAS-BTV-22-001



JINST 15 (2020) P06005



Transformer-based GN2X tagger:  
**~x3 QCD and x2 top background rejection**

DeepAK8 → ParticleNet:  
**x5 QCD background rejection**

Comparing with early approaches  
 Another **~x5 improvement achieved**



# Current boosted taggers in ATLAS/CMS

## CMS: DeepAK8 and ParticleNet algorithms

[JINST 15 \(2020\) P06005](#)

Output

Category	Label
Higgs	H (bb)
	H (cc)
	H (VV* → qq qq)
Top	top (bcq)
	top (bqq)
	top (bc)
	top (bq)
W	W (cq)
	W (qq)
Z	Z (bb)
	Z (cc)
	Z (qq)
QCD	QCD (bb)
	QCD (cc)
	QCD (b)
	QCD (c)
	QCD (others)

## DeepAK8-MD, ParticleNet-MD, and DeepDoubleX algorithms

[CMS-PAS-BTV-22-001](#)

- Focus on variable-mass resonance decays
- X → bb, cc, qq and QCD (5 subclasses)

## GN2X tagger

[ATL-PHYS-PUB-2023-021](#)

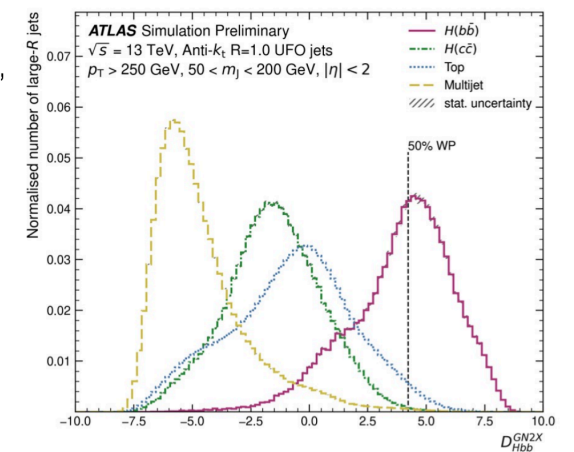
- including flat-mass H → bb, cc and t → bqq samples, with QCD

## GN2X Outputs

- GN2X adds a H → cc output class in addition to the H → bb, top and QCD classes from the previous tagger
- A discriminant score is built using a weighted log likelihood ratio similar to what's used for small-R tagging
- GN2X also includes the same auxiliary vertexing and track origin classification tasks present in GN1/GN2

$$D_{Hbb}^{GN2X} = \ln \left( \frac{P_{Hbb}}{f_{Hcc} \cdot P_{Hcc} + f_{top} \cdot P_{top} + (1 - f_{Hcc} - f_{top}) \cdot P_{QCD}} \right)$$

[Jackson's slides](#)



# How can we accelerate the pace?

## CMS: DeepAK8 and ParticleNet algorithms

[JINST 15 \(2020\) P06005](#)

Output

Category	Label
Higgs	H (bb)
	H (cc)
	H (VV* → qqqq)
Top	top (bcq)
	top (bqq)
	top (bc)
	top (bq)
W	W (cq)
	W (qq)
Z	Z (bb)
	Z (cc)

## DeepAK8-MD, ParticleNet-MD, and DeepDoubleX algorithms

[CMS-PAS-BTV-22-001](#)

- Focus on variable-mass resonance decays
- X → bb, cc, qq and QCD (5 subclasses)

## GN2X tagger

[ATL-PHYS-PUB-2023-021](#)

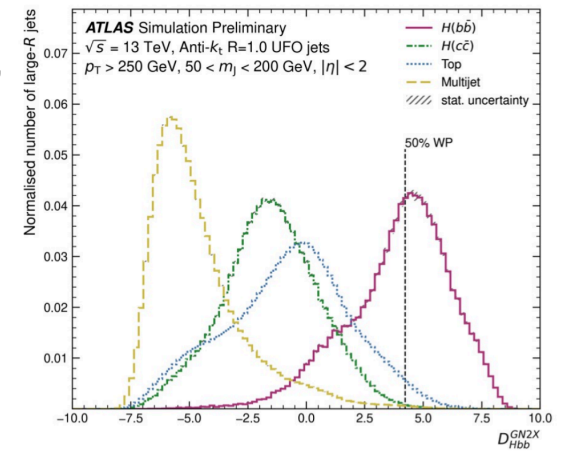
- including flat-mass H → bb, cc and t → bqq samples, with QCD

## GN2X Outputs

- GN2X adds a H → cc output class in addition to the H → bb, top and QCD classes from the previous tagger
- A discriminant score is built using a weighted log likelihood ratio similar to what's used for small-R tagging
- GN2X also includes the same auxiliary vertexing and track origin classification tasks present in GN1/GN2

$$\ln \left( \frac{P_{Hbb}}{f_{Hcc} \cdot P_{Hcc} + f_{top} \cdot P_{top} + (1 - f_{Hcc} - f_{top}) \cdot P_{QCD}} \right)$$

[Jackson's slides](#)



## All final states!

Major types	Index range	Label names
Resonant jets: X → 2 prong	0–14	bb, cc, ss, qq, bc, cs, bq, cq, sq, gg, ee, μμ, τ <sub>h</sub> τ <sub>e</sub> , τ <sub>h</sub> τ <sub>μ</sub> , τ <sub>h</sub> τ <sub>h</sub>
Resonant jets: X → 3 or 4 prong	15–160	bbbb, bbcc, bbss, bbqq, bbgg, bbee, bbμμ, bbτ <sub>h</sub> τ <sub>e</sub> , bbτ <sub>h</sub> τ <sub>μ</sub> , bbτ <sub>h</sub> τ <sub>h</sub> , bbb, bbc, bbs, bbq, bbq, bbe, bbμ, cccc, ccss, ccqq, ccgg, ccee, ccμμ, ccτ <sub>h</sub> τ <sub>e</sub> , ccτ <sub>h</sub> τ <sub>μ</sub> , ccτ <sub>h</sub> τ <sub>h</sub> , ccb, ccc, ccs, ccq, ccg, cce, ccμ, ssss, ssqq, ssgg, ssee, ssμμ, ssτ <sub>h</sub> τ <sub>e</sub> , ssτ <sub>h</sub> τ <sub>μ</sub> , ssτ <sub>h</sub> τ <sub>h</sub> , ssb, ssc, sss, ssq, ssg, sse, ssμ, qqqq, qqgg, qqee, qqμμ, qqτ <sub>h</sub> τ <sub>e</sub> , qqτ <sub>h</sub> τ <sub>μ</sub> , qqτ <sub>h</sub> τ <sub>h</sub> , qqb, qqc, qqe, qqμ, gggg, ggee, ggμμ, ggτ <sub>h</sub> τ <sub>e</sub> , ggτ <sub>h</sub> τ <sub>μ</sub> , ggτ <sub>h</sub> τ <sub>h</sub> , ggb, ggc, ggs, ggq, ggg, gge, ggμ, bee, cee, see, qee, gee, bμμ, cμμ, sμμ, qμμ, gμμ, bτ <sub>h</sub> τ <sub>e</sub> , cτ <sub>h</sub> τ <sub>e</sub> , sτ <sub>h</sub> τ <sub>e</sub> , qτ <sub>h</sub> τ <sub>e</sub> , gτ <sub>h</sub> τ <sub>e</sub> , bτ <sub>h</sub> τ <sub>μ</sub> , cτ <sub>h</sub> τ <sub>μ</sub> , sτ <sub>h</sub> τ <sub>μ</sub> , qτ <sub>h</sub> τ <sub>μ</sub> , gτ <sub>h</sub> τ <sub>μ</sub> , bτ <sub>h</sub> τ <sub>h</sub> , cτ <sub>h</sub> τ <sub>h</sub> , sτ <sub>h</sub> τ <sub>h</sub> , qτ <sub>h</sub> τ <sub>h</sub> , gτ <sub>h</sub> τ <sub>h</sub> , qqqb, qqqc, qqqs, bbcq, ccbs, ccbq, ccsq, sscq, qqbc, qqbs, qqcs, bcsq, bcs, bcq, bsq, csq, bcev, csev, bqev, cqev, sqev, qqev, bcμν, csμν, bqμν, cqμν, sqμν, qqμν, bcτ <sub>e</sub> ν, csτ <sub>e</sub> ν, bqτ <sub>e</sub> ν, cqτ <sub>e</sub> ν, sqτ <sub>e</sub> ν, qqτ <sub>e</sub> ν, bcτ <sub>μ</sub> ν, csτ <sub>μ</sub> ν, bqτ <sub>μ</sub> ν, cqτ <sub>μ</sub> ν, sqτ <sub>μ</sub> ν, qqτ <sub>μ</sub> ν, bcτ <sub>h</sub> ν, csτ <sub>h</sub> ν, bqτ <sub>h</sub> ν, cqτ <sub>h</sub> ν, sqτ <sub>h</sub> ν, qqτ <sub>h</sub> ν
QCD jets	161–187	bbccss, bbccs, bbcc, bbcss, bbcs, bbc, bbss, bbs, bb, bccss, bccs, bcc, bcss, bcs, bc, bss, bs, b, ccss, ccs, cc, css, cs, c, ss, s, others



# How can we accelerate the pace?

## *Sophon* (智子): Signature-Oriented Pre-training for Heavy-resonant Observation



	Index range	Label names
$X \rightarrow 2$ prong	0–14	$bb, cc, ss, qq, bc, cs, bq, cq, sq, gg, ee, \mu\mu, \tau_h\tau_e, \tau_h\tau_\mu, \tau_h\tau_h$
Resonant jets: $X \rightarrow 3$ or 4 prong	15–160	$bbbb, bbcc, bbss, bbqq, bbgg, bbee, bb\mu\mu, bb\tau_h\tau_e, bb\tau_h\tau_\mu, bb\tau_h\tau_h, bbb, bbc, bbs, bbq, bbg, bbe, bb\mu, cccc, ccss, ccqq, ccgg, ccee, cc\mu\mu, cc\tau_h\tau_e, cc\tau_h\tau_\mu, cc\tau_h\tau_h, ccb, ccc, ccs, ccq, ccg, cce, cc\mu, ssss, ssqq, ssgg, ssee, ss\mu\mu, ss\tau_h\tau_e, ss\tau_h\tau_\mu, ss\tau_h\tau_h, ssb, ssc, sss, ssq, ssg, sse, ss\mu, qqqq, qqgg, qqee, qq\mu\mu, qq\tau_h\tau_e, qq\tau_h\tau_\mu, qq\tau_h\tau_h, qqb, qqc, qqs, qqg, qqe, qq\mu, gggg, ggee, gg\mu\mu, gg\tau_h\tau_e, gg\tau_h\tau_\mu, gg\tau_h\tau_h, ggb, ggc, ggs, ggq, ggg, gge, gg\mu, bee, cee, see, qee, gee, b\mu\mu, c\mu\mu, s\mu\mu, q\mu\mu, g\mu\mu, b\tau_h\tau_e, c\tau_h\tau_e, s\tau_h\tau_e, q\tau_h\tau_e, g\tau_h\tau_e, b\tau_h\tau_\mu, c\tau_h\tau_\mu, s\tau_h\tau_\mu, q\tau_h\tau_\mu, g\tau_h\tau_\mu, b\tau_h\tau_h, c\tau_h\tau_h, s\tau_h\tau_h, q\tau_h\tau_h, g\tau_h\tau_h, qqqb, qqqc, qqqs, bbcq, cchs, ccbq, ccsq, sscq, qqbc, qqbs, qqcs, bcsq, bcs, bcq, bsq, csq, bcev, csev, bqev, cqev, sqev, qqev, bc\mu\nu, cs\mu\nu, bq\mu\nu, cq\mu\nu, sq\mu\nu, qq\mu\nu, bc\tau_e\nu, cs\tau_e\nu, bq\tau_e\nu, cq\tau_e\nu, sq\tau_e\nu, qq\tau_e\nu, bc\tau_\mu\nu, cs\tau_\mu\nu, bq\tau_\mu\nu, cq\tau_\mu\nu, sq\tau_\mu\nu, qq\tau_\mu\nu, bc\tau_h\nu, cs\tau_h\nu, bq\tau_h\nu, cq\tau_h\nu, sq\tau_h\nu, qq\tau_h\nu$
QCD jets	161–187	$bbccss, bbccs, bbcc, bbcss, bbcs, bbc, bbss, bbs, bb, bccss, bccs, bcc, bcss, bcs, bc, bss, bs, b, ccss, ccs, cc, css, cs, c, ss, s, \text{others}$

### → Major concerns:

- ❖ Will the model achieve the best performance for each specific task?
- ❖ What can we use this model for, beyond its identification task supported by its final states?

# How can we accelerate the pace?

## *Sophon* (智子): Signature-Oriented Pre-training for Heavy-resonant Observation

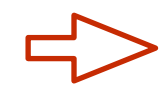


	Index range	Label names
$X \rightarrow 2$ prong	0–14	$bb, cc, ss, qq, bc, cs, bq, cq, sq, gg, ee, \mu\mu, \tau_h\tau_e, \tau_h\tau_\mu, \tau_h\tau_h$
Resonant jets: $X \rightarrow 3$ or 4 prong	15–160	$bbbb, bbcc, bbss, bbqq, bbgg, bbee, bb\mu\mu, bb\tau_h\tau_e, bb\tau_h\tau_\mu, bb\tau_h\tau_h, bbb, bbc, bbs, bbq, bbg, bbe, bb\mu, cccc, ccss, ccqq, ccgg, ccee, cc\mu\mu, cc\tau_h\tau_e, cc\tau_h\tau_\mu, cc\tau_h\tau_h, ccb, ccc, ccs, ccq, ccg, cce, cc\mu, ssss, ssqq, ssgg, ssee, ss\mu\mu, ss\tau_h\tau_e, ss\tau_h\tau_\mu, ss\tau_h\tau_h, ssb, ssc, sss, ssq, ssg, sse, ss\mu, qqqq, qqgg, qqee, qq\mu\mu, qq\tau_h\tau_e, qq\tau_h\tau_\mu, qq\tau_h\tau_h, qqb, qqc, qqs, qqg, qqe, qq\mu, gggg, ggee, gg\mu\mu, gg\tau_h\tau_e, gg\tau_h\tau_\mu, gg\tau_h\tau_h, ggb, ggc, ggs, ggq, ggg, gge, gg\mu, bee, cee, see, qee, gee, b\mu\mu, c\mu\mu, s\mu\mu, q\mu\mu, g\mu\mu, b\tau_h\tau_e, c\tau_h\tau_e, s\tau_h\tau_e, q\tau_h\tau_e, g\tau_h\tau_e, b\tau_h\tau_\mu, c\tau_h\tau_\mu, s\tau_h\tau_\mu, q\tau_h\tau_\mu, g\tau_h\tau_\mu, b\tau_h\tau_h, c\tau_h\tau_h, s\tau_h\tau_h, q\tau_h\tau_h, g\tau_h\tau_h, qqqb, qqqc, qqqs, bbcq, cchs, ccbq, ccsq, sscq, qqbc, qqbs, qqcs, bcsq, bcs, bcq, bsq, csq, bcev, csev, bqev, cqev, sqev, qqev, bc\mu\nu, cs\mu\nu, bq\mu\nu, cq\mu\nu, sq\mu\nu, qq\mu\nu, bc\tau_e\nu, cs\tau_e\nu, bq\tau_e\nu, cq\tau_e\nu, sq\tau_e\nu, qq\tau_e\nu, bc\tau_\mu\nu, cs\tau_\mu\nu, bq\tau_\mu\nu, cq\tau_\mu\nu, sq\tau_\mu\nu, qq\tau_\mu\nu, bc\tau_h\nu, cs\tau_h\nu, bq\tau_h\nu, cq\tau_h\nu, sq\tau_h\nu, qq\tau_h\nu$
QCD jets	161–187	$bbccss, bbccs, bbcc, bbcss, bbcs, bbc, bbss, bbs, bb, bccss, bccs, bcc, bcss, bcs, bc, bss, bs, b, ccss, ccs, cc, css, cs, c, ss, s, \text{others}$



### → Major concerns:

- ❖ Will the model achieve the best performance for each specific task?
- ❖ What can we use this model for, beyond its identification task supported by its final states?



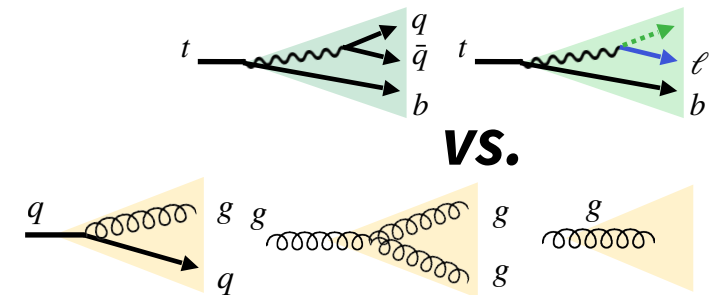
*Yes, it does. The Particle Transformer can support each task to reach its optimal performance*

*We can regard it as a true “**based model**” and fine-tune it for wider ranges of downstream tasks*

# Propose “Large model for large-scale classification”

## View from jet tagging

- Instead of training dedicated jet taggers, we consider **multi-class classification** with  $N(\text{class})$  reaches  $\mathcal{O}(100)$ 
  - ❖ statistical insights: an ideal multi-class classifier is a stack of ideal binary classifiers
- The model should be **large** → carry enough capacity
- The classes should be comprehensive → **tagging ability can be further generalized by fine-tuning**

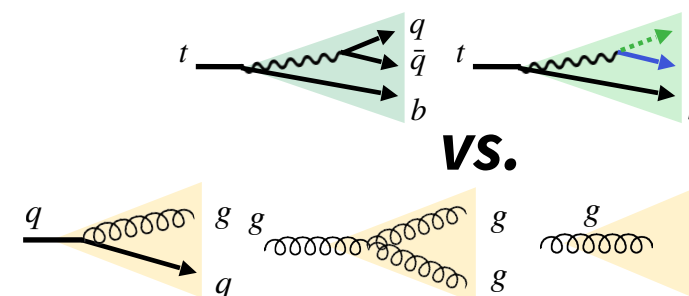




# Propose “Large model for large-scale classification”

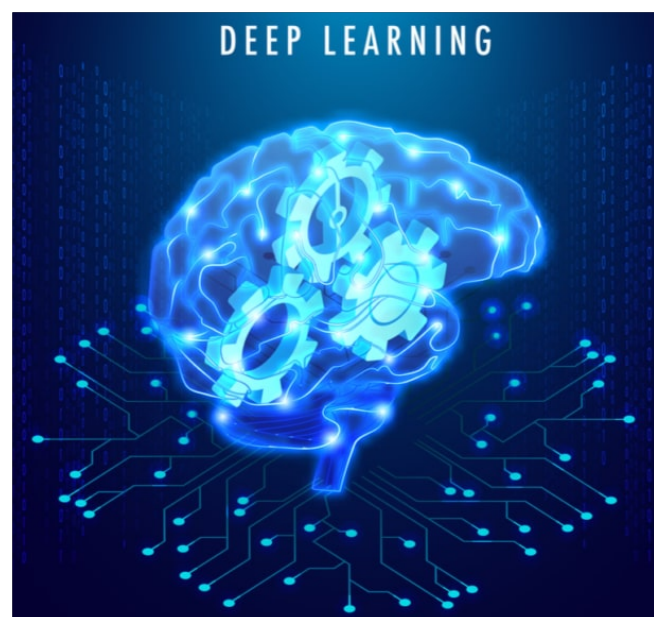
## View from jet tagging

- Instead of training dedicated jet taggers, we consider **multi-class classification** with  $N(\text{class})$  reaches  $\mathcal{O}(100)$ 
  - ❖ statistical insights: an ideal multi-class classifier is a stack of ideal binary classifiers
- The model should be **large** → carry enough capacity
- The classes should be comprehensive → **tagging ability can be further generalized by fine-tuning**



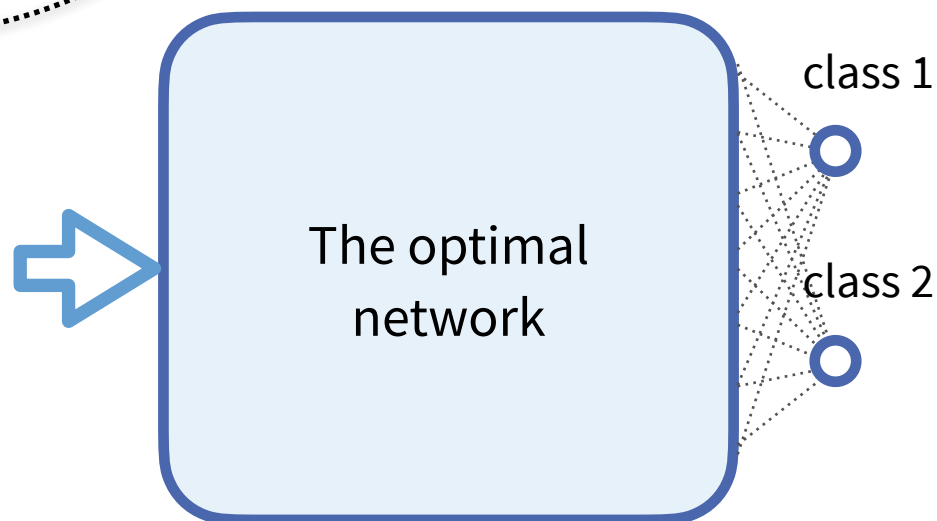
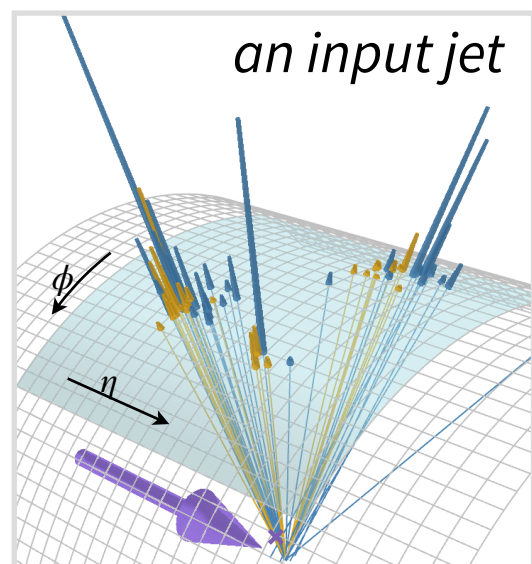
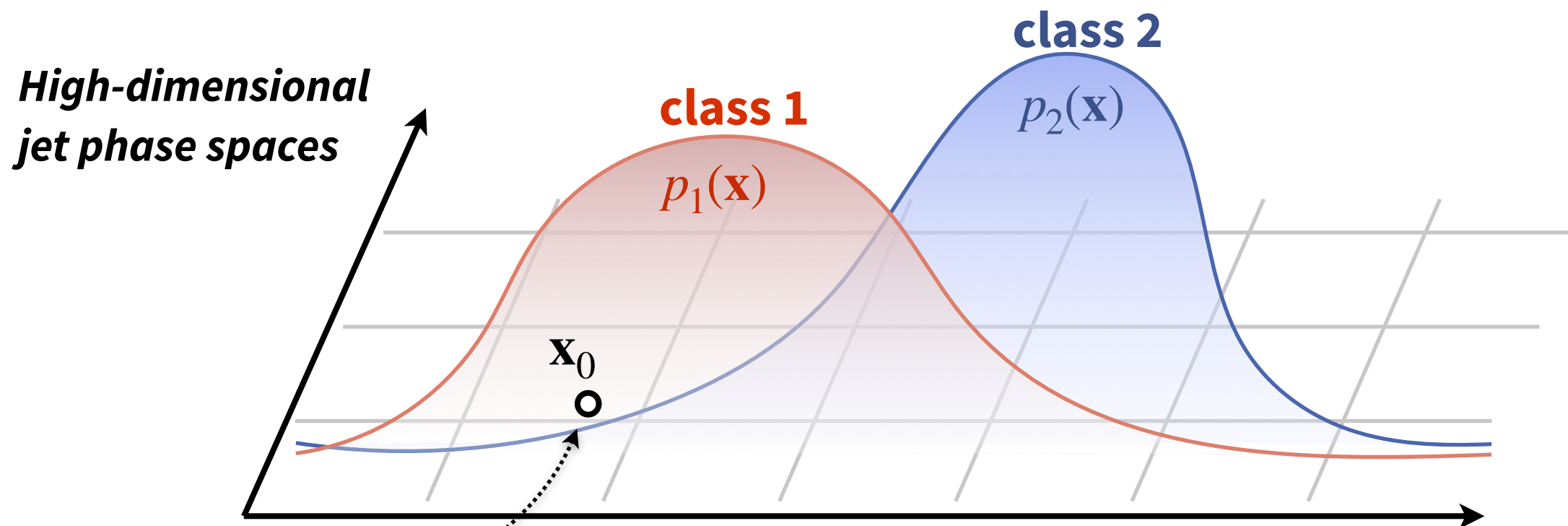
## View from a pre-training solution

- Based on a comprehensive jet dataset, we hope to pre-train **a base model** to facilitate **all** LHC analyses exploring the large- $R$  jet
- Set the training task: let the model learn to connect **“what a jet is like”** to **“which truth signature the jet reveals”** (= jet label in our case)
  - ❖ “jet labels” are simple signatures to explore
    - pre-training it as a classifier is just a starting point in this sense!



# Statistical essence of jet tagging problem

- **Question: where is the limit of jet tagging?**
- **Answer: the probability density ratio of two classes provides the optimal tagging**



- ❖ Ideal classifier network results in  
 $g_1 : g_2 : \dots = p_1(\mathbf{x}_0) : p_2(\mathbf{x}_0) : \dots$
- ❖ It is a direct estimation of  $p$
- ❖ The **network capacity** decides how close the estimation is

# Statistical property of multi-class classifier

→ Statistical theory shows that:

A **multi-class** classifier with minimum **cross-entropy loss** **estimates the probability ratios** on the input classes:

$$g_i(\mathbf{x}) = \frac{p(\text{class} = i | \mathbf{x})}{\sum_{j=1}^{N_{\text{out}}} p(\text{class} = j | \mathbf{x})}$$

hence it contains **all the information** the ideal  $N(N - 1)$  binary classifiers can do



# Statistical property of multi-class classifier

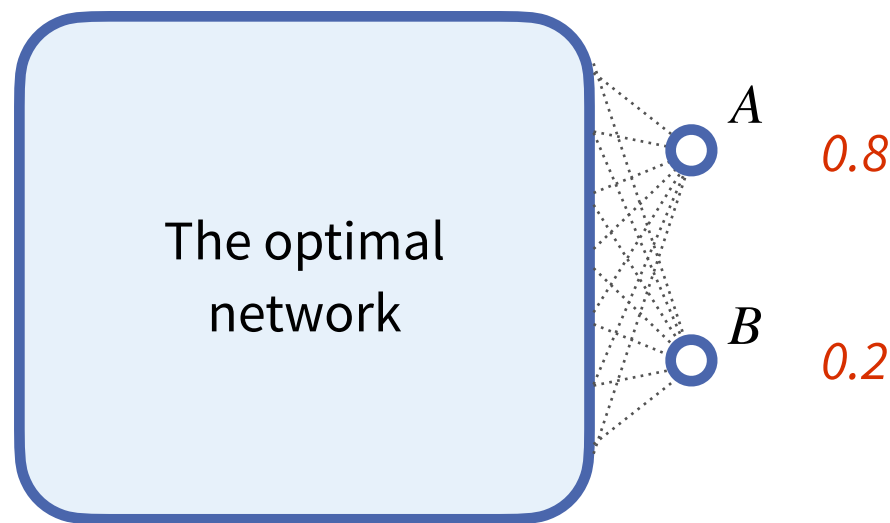
→ Statistical theory shows that:

A **multi-class** classifier with minimum **cross-entropy loss** **estimates the probability ratios** on the input classes:

$$g_i(\mathbf{x}) = \frac{p(\text{class} = i | \mathbf{x})}{\sum_{j=1}^{N_{\text{out}}} p(\text{class} = j | \mathbf{x})}$$

hence it contains **all the information** the ideal  $N(N - 1)$  binary classifiers can do

**Two properties:**



splitting class A

○  $A_1$  0.55  
○  $A_2$  0.25  
○  $B$  0.2

$p_A = p_{A_1} + p_{A_2}$   
remains the same

adding class C

○  $A$  0.6  
○  $B$  0.15  
○  $C$  0.25

$p_A/p_B$   
remains the same

# Statistical property of multi-class classifier

→ Statistical theory shows that:

A **multi-class** classifier with minimum **cross-entropy loss** **estimates the probability ratios** on the input classes:

$$g_i(\mathbf{x}) = \frac{p(\text{class} = i | \mathbf{x})}{\sum_{j=1}^{N_{\text{out}}} p(\text{class} = j | \mathbf{x})}$$

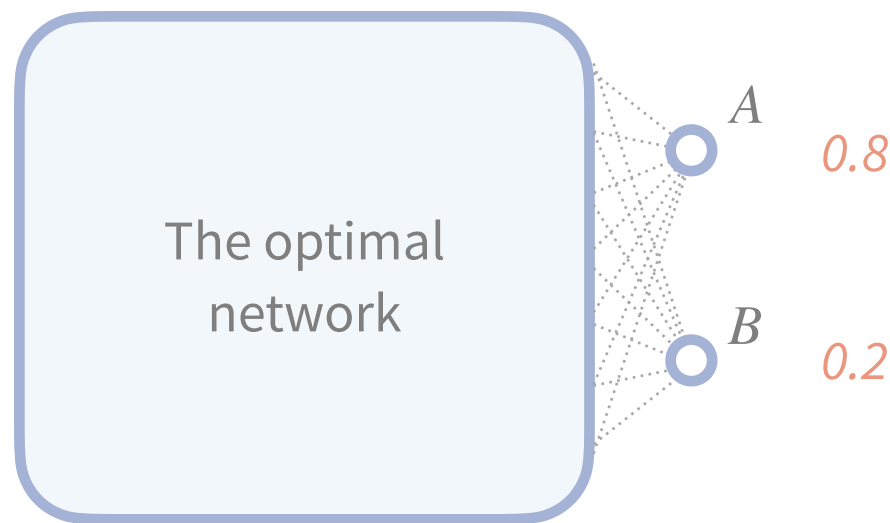
## The key question in this context

Does the model's capacity still enable us to reach the best achievable performance in existing tasks?

**Our result will show: Yes.**

hence it contains **all the information** the ideal  $N(N - 1)$  binary classifiers can do

**Two properties:**



splitting class A

○  $A_1$  0.55  
○  $A_2$  0.25  
○ B 0.2

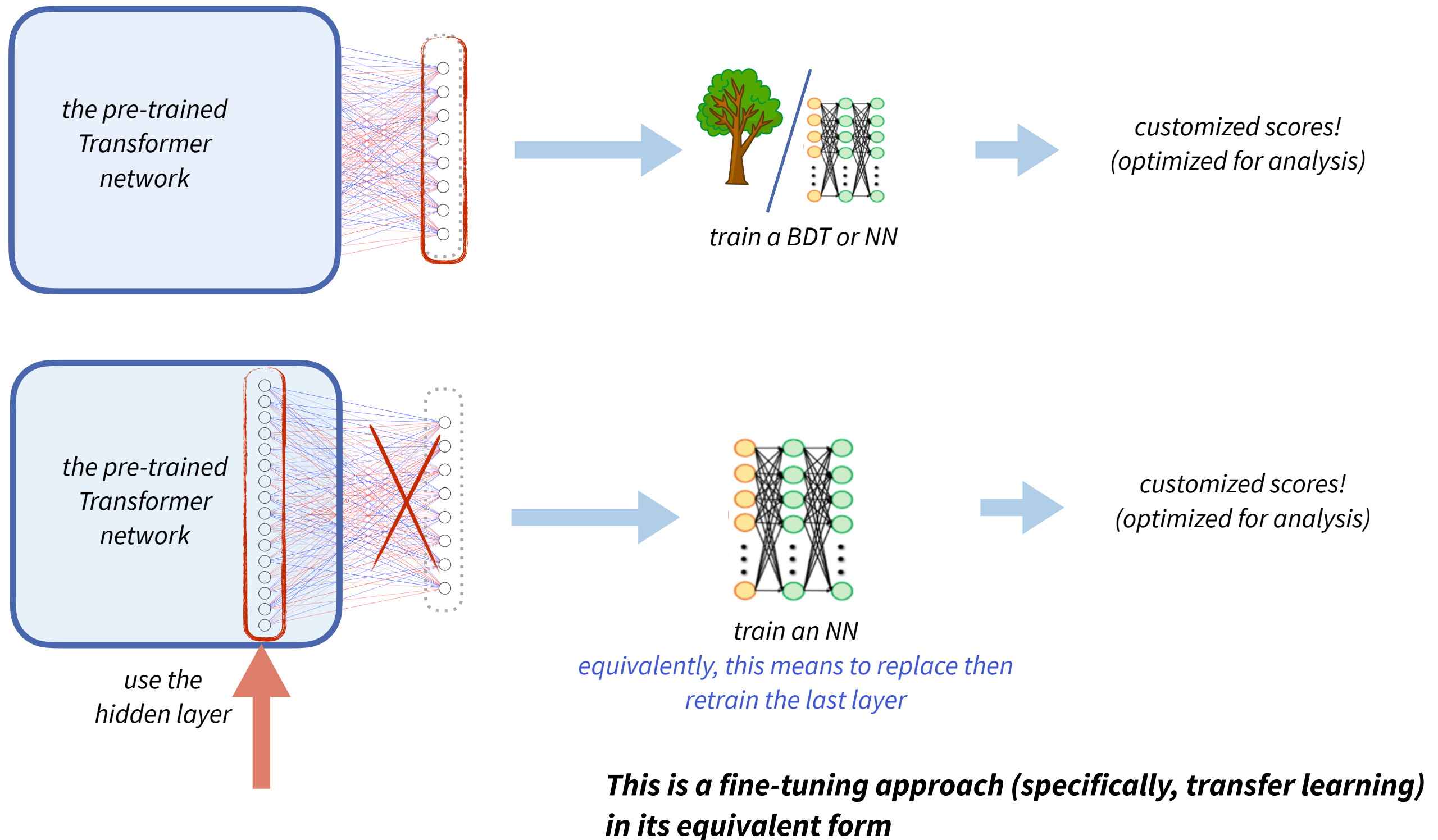
$p_A = p_{A_1} + p_{A_2}$   
remains the same

adding class C

○ A 0.6  
○ B 0.15  
○ C 0.25

$p_A/p_B$   
remains the same

# A glance into fine-tuning spirits





# Introducing Sophon (智子)

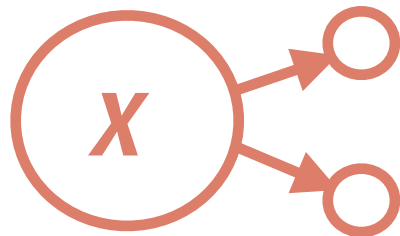
[arXiv:2405.12972](https://arxiv.org/abs/2405.12972)

<https://github.com/jet-universe/sophon>

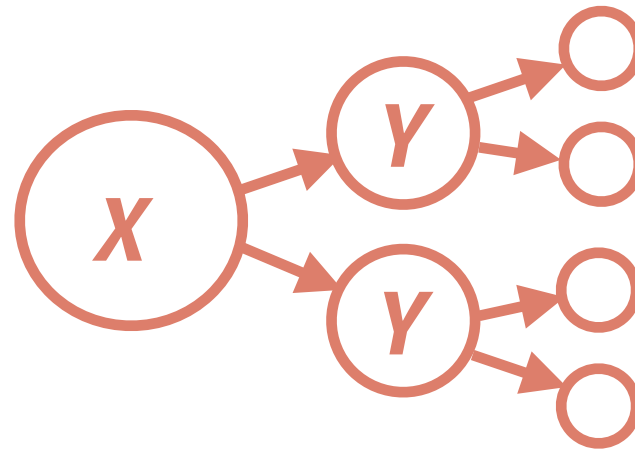


- Signature-Oriented Pre-training for Heavy-resonant Observation
- the model is based on Particle Transformer (ParT) architecture
- a pre-trained model on a newly developed comprehensive dataset: JetClass-II
  - *finely categorized labels:*

**Resonant jet:**  
 $X \rightarrow 2$  prong

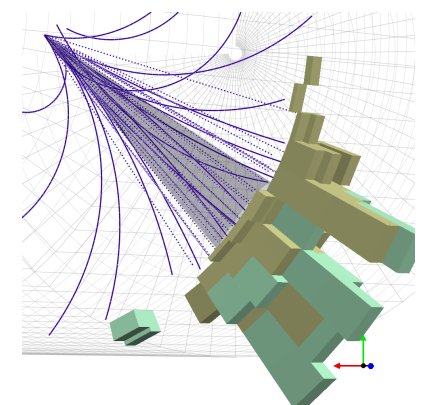


**Resonant jet:**  
 $X \rightarrow 3/4$  prong



$bb/cc/ss/qq/gg/ee/\mu\mu/\tau\tau$   
 $bc/bq/cs/cq$   
 $ev/\mu\nu/\nu\nu$

**QCD jets**



contributed  
final states:

$bb/cc/ss/qq/gg/ee/\mu\mu/\tau\tau$   
 $bc/bq/cs/cq$

all combination of  $Y$  decays,  
resulting to 4-prong or 3-prong

**Key property:** we do not focus on any specific  $X$  and  $Y$  masses  
Their masses are variables: ranges from 20-500 GeV

# Introducing Sophon (智子)

[arXiv:2405.12972](https://arxiv.org/abs/2405.12972)
<https://github.com/jet-universe/sophon>


- Signature-Oriented Pre-training for Heavy-resonant Observation
- the model is based on Particle Transformer (ParT) architecture
- a pre-trained model on a newly developed comprehensive dataset: **JetClass-II**

▸ *finely categorized labels:*

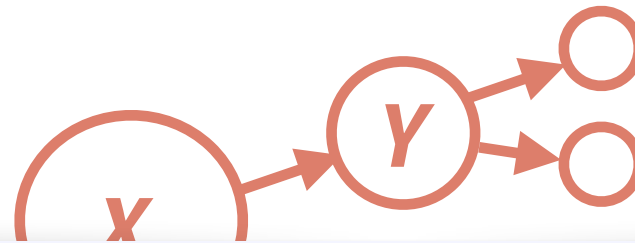
**Resonant jet:**

$X \rightarrow 2$  prong



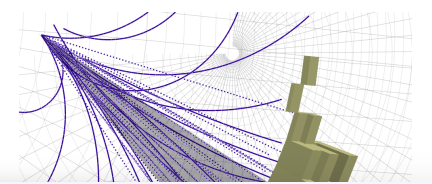
**Resonant jet:**

$X \rightarrow 3/4$  prong



*bb/cc/ss/qq/gg/ee/μμ/ττ*  
*bc/bq/cs/cq*  
*ev/μν/νν*

**QCD jets**



Major types	Index range	Label names	<b>All final states!</b>
Resonant jets: $X \rightarrow 2$ prong	0–14	$bb, cc, ss, qq, bc, cs, bq, cq, sq, gg, ee, \mu\mu, \tau_h\tau_e, \tau_h\tau_\mu, \tau_h\tau_h$	
Resonant jets: $X \rightarrow 3$ or 4 prong	15–160	$bbbb, bbcc, bbss, bbqq, bbgg, bbee, bb\mu\mu, bb\tau_h\tau_e, bb\tau_h\tau_\mu, bb\tau_h\tau_h, bbb, bbc, bbs, bbq, bbg, bbe, bb\mu, cccc, ccss, ccqq, ccgg, ccee, cc\mu\mu, cc\tau_h\tau_e, cc\tau_h\tau_\mu, cc\tau_h\tau_h, ccb, ccc, ccs, ccq, ccg, cce, cc\mu, ssss, ssqq, ssgg, ssee, ss\mu\mu, ss\tau_h\tau_e, ss\tau_h\tau_\mu, ss\tau_h\tau_h, ssb, ssc, sss, ssq, ssg, sse, ss\mu, qqqq, qqgg, qqee, qq\mu\mu, qq\tau_h\tau_e, qq\tau_h\tau_\mu, qq\tau_h\tau_h, qqb, qqc, qqe, qq\mu, gggg, ggee, gg\mu\mu, gg\tau_h\tau_e, gg\tau_h\tau_\mu, gg\tau_h\tau_h, ggb, ggc, ggs, ggq, ggg, gge, gg\mu, bee, cee, see, qee, gee, b\mu\mu, c\mu\mu, s\mu\mu, q\mu\mu, g\mu\mu, b\tau_h\tau_e, c\tau_h\tau_e, s\tau_h\tau_e, q\tau_h\tau_e, g\tau_h\tau_e, b\tau_h\tau_\mu, c\tau_h\tau_\mu, s\tau_h\tau_\mu, q\tau_h\tau_\mu, g\tau_h\tau_\mu, b\tau_h\tau_h, c\tau_h\tau_h, s\tau_h\tau_h, q\tau_h\tau_h, g\tau_h\tau_h, qqqb, qqqc, qqqs, bbcq, ccbs, ccbq, ccsq, sscq, qqbc, qqbs, qqcs, bcsq, bcs, bcq, bsq, csq, bcev, csev, bqev, cqev, sqev, qqev, bc\mu\nu, cs\mu\nu, bq\mu\nu, cq\mu\nu, sq\mu\nu, qq\mu\nu, bc\tau_e\nu, cs\tau_e\nu, bq\tau_e\nu, cq\tau_e\nu, sq\tau_e\nu, qq\tau_e\nu, bc\tau_\mu\nu, cs\tau_\mu\nu, bq\tau_\mu\nu, cq\tau_\mu\nu, sq\tau_\mu\nu, qq\tau_\mu\nu, bc\tau_h\nu, cs\tau_h\nu, bq\tau_h\nu, cq\tau_h\nu, sq\tau_h\nu, qq\tau_h\nu$	
QCD jets	161–187	$bbccss, bbccs, bbcc, bbcss, bbcs, bbc, bbss, bbs, bb, bccss, bccs, bcc, bcss, bcs, bc, bss, bs, b, ccss, ccs, cc, css, cs, c, ss, s, \text{others}$	

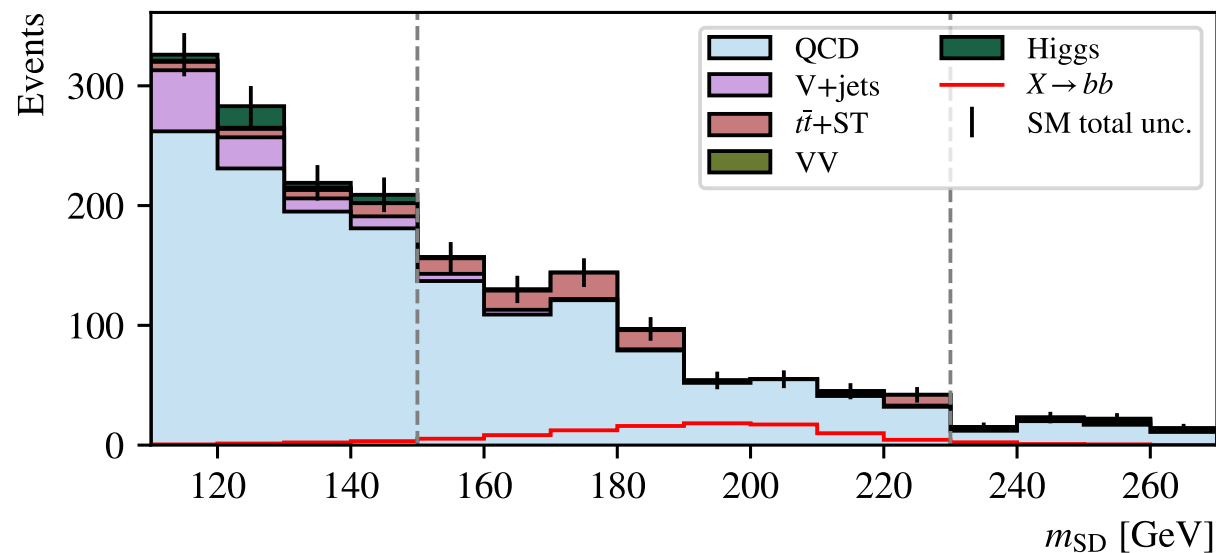
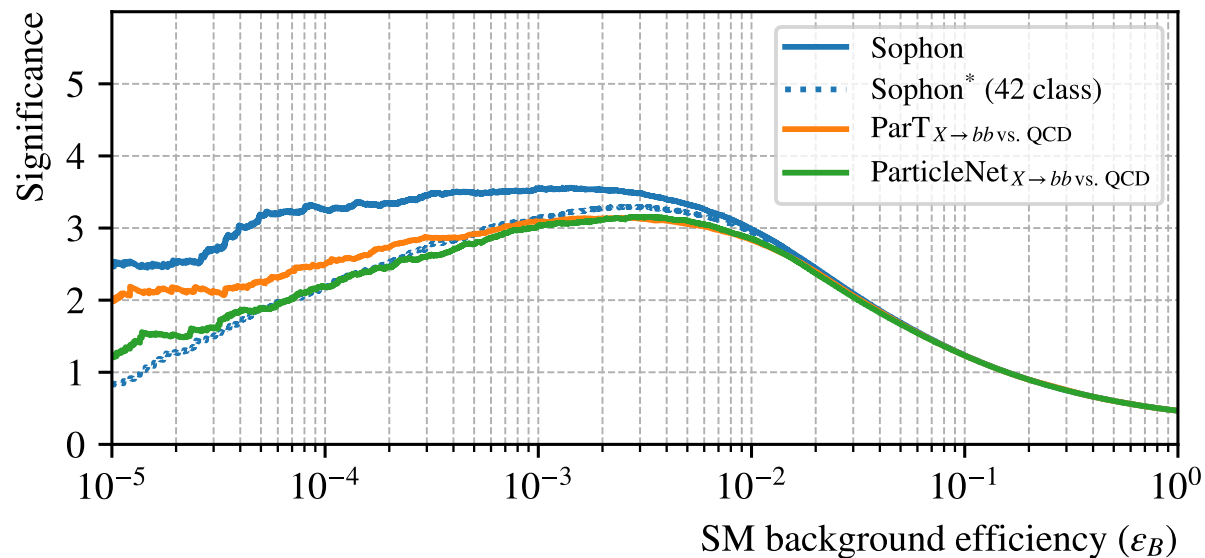
# Sophon: performance benchmark

[arXiv:2405.12972](https://arxiv.org/abs/2405.12972)

Search significance:

$$Z = \sqrt{2((s+b)\log(1+s/b) - s)}$$

## Direct tagging ability



- Apply tagger selection
- Check discrimination power of  $X(200 \text{ GeV}) \rightarrow \mathbf{bb}$  signal vs. all backgrounds

- **Sophon** (training on 188 classes) has best performance

$$\text{discr}(X \rightarrow bb \text{ vs. QCD}) = \frac{g_{X \rightarrow bb}}{g_{X \rightarrow bb} + \sum_{l=1}^{27} g_{\text{QCD}_l}}$$

- Performance gain does come from large-scale classification (compared to **Sophon\*** (42 classes))
- **ParT** and **ParticleNet** for binary classification: they represent the best performance we can reach in experiment now



# Sophon: performance benchmark

[arXiv:2405.12972](https://arxiv.org/abs/2405.12972)

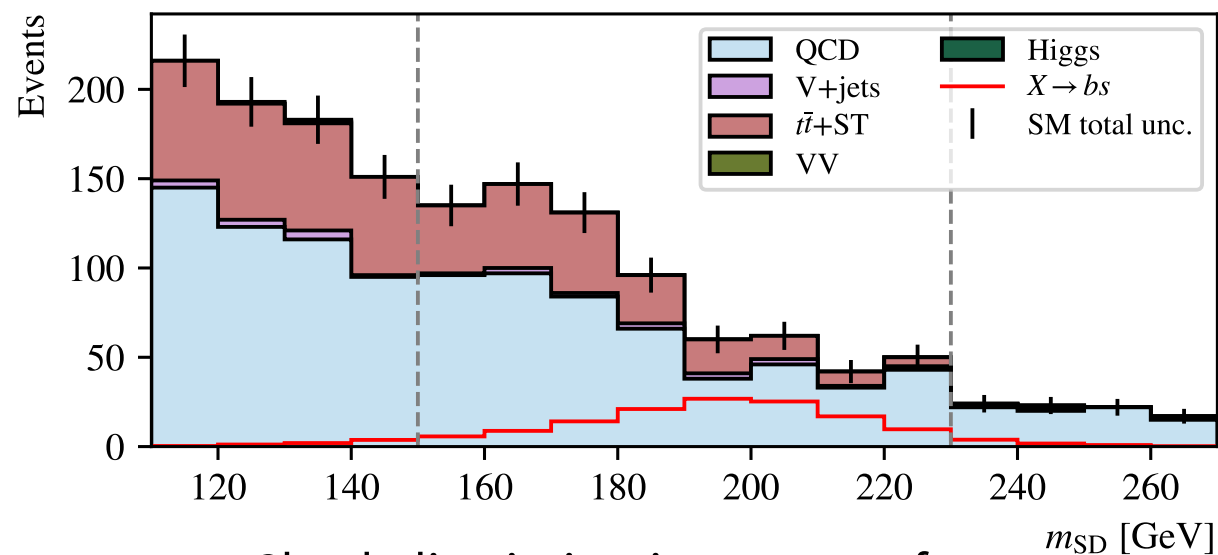
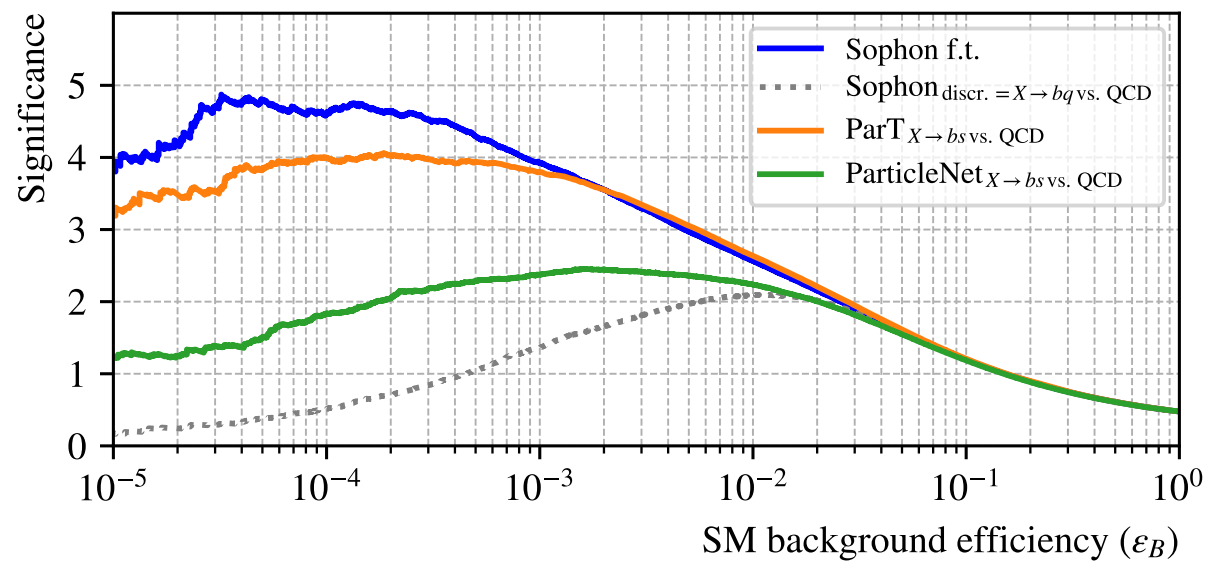
Search significance:

$$Z = \sqrt{2((s+b)\log(1+s/b) - s)}$$



## Transfer learning ability

(adapt the tagger to a new classification task)



- Check discrimination power of  $X(200 \text{ GeV}) \rightarrow \mathbf{bs}$  signal vs. all backgrounds

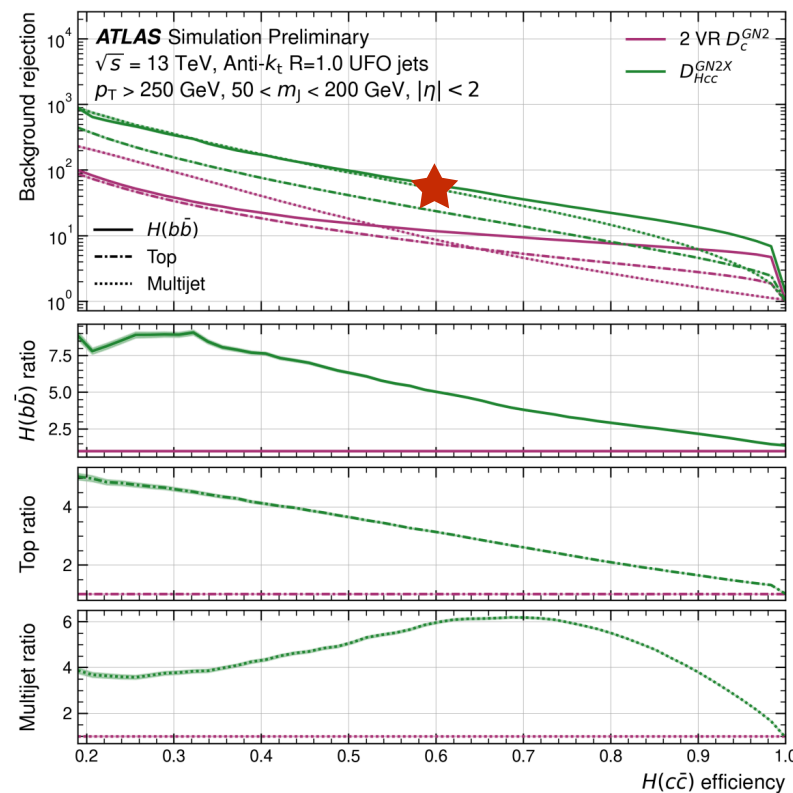
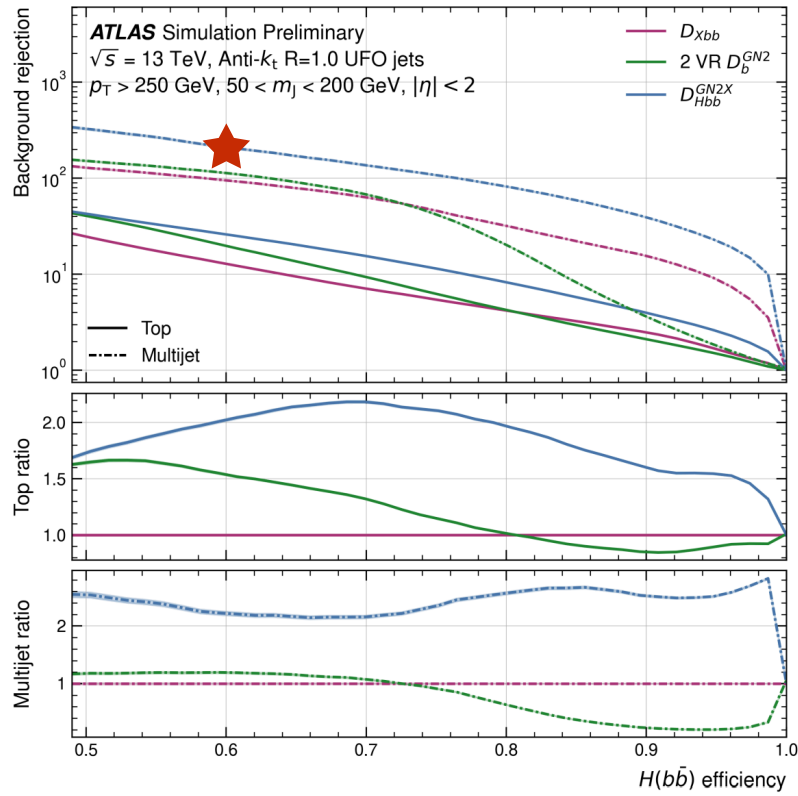
- **Sophon** (training on 188 classes) reaches the best performance **after fine-tuned (via transfer learning)**
- **ParT** and **ParticleNet** for binary  $X \rightarrow bs$  vs QCD classification: they reveal the best performance we can reach in the experiment now

# Sophon: close to real experimental performance?

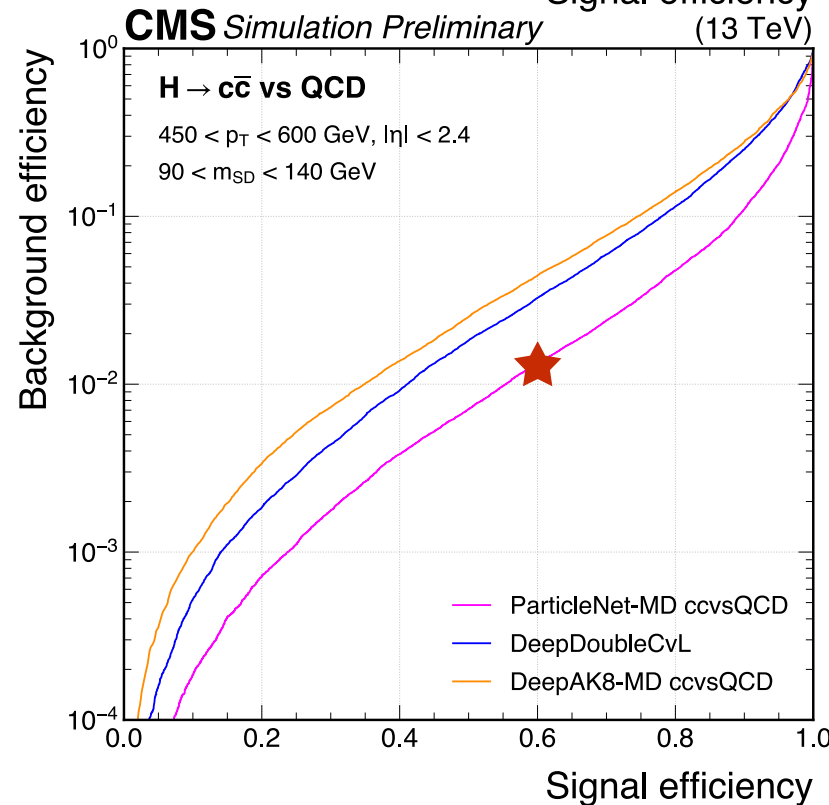
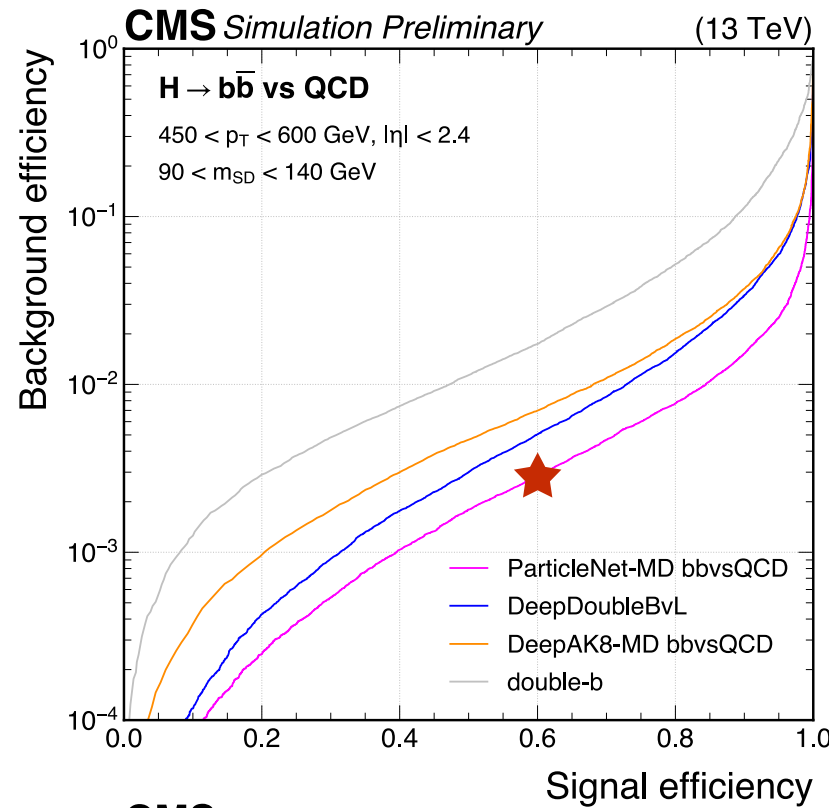
★ marks

QCD BKG rej at  
signal eff. = 60%

## ATLAS results [ATL-PHYS-PUB-2023-021](#)

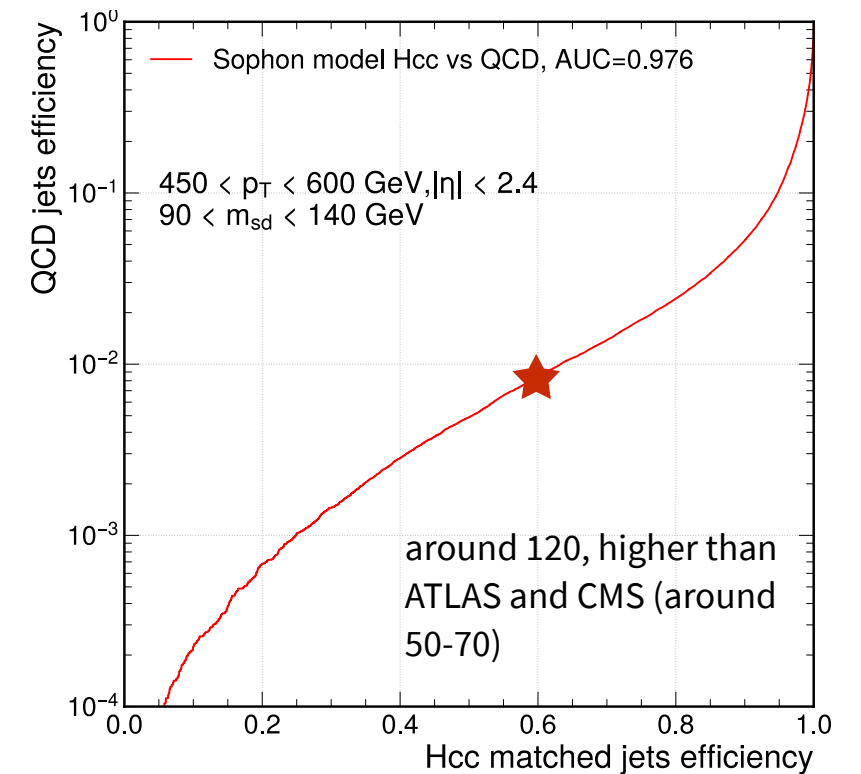
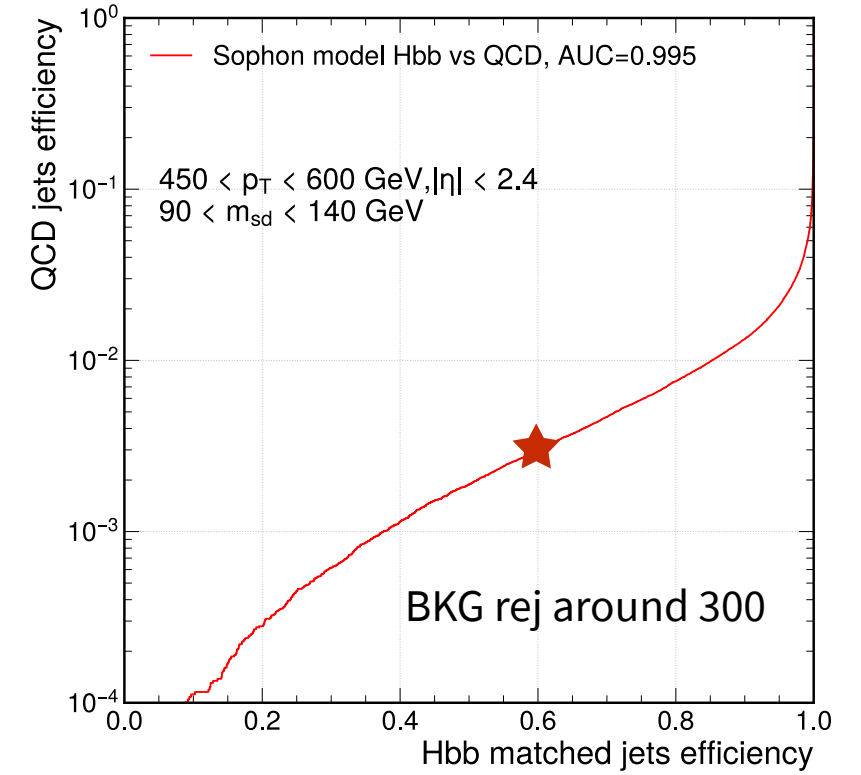


## CMS results [CMS-PAS-BTV-22-001](#)



## Sophon results

(performance on Delphes)

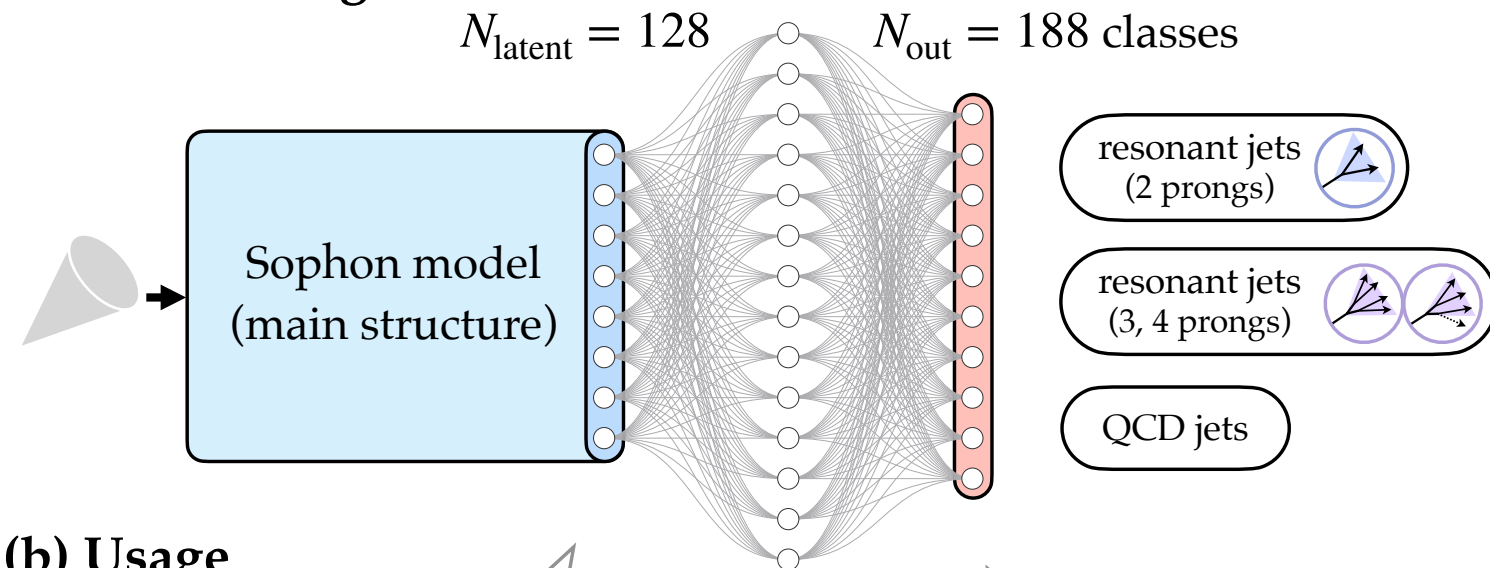


# Implications for LHC resonance search

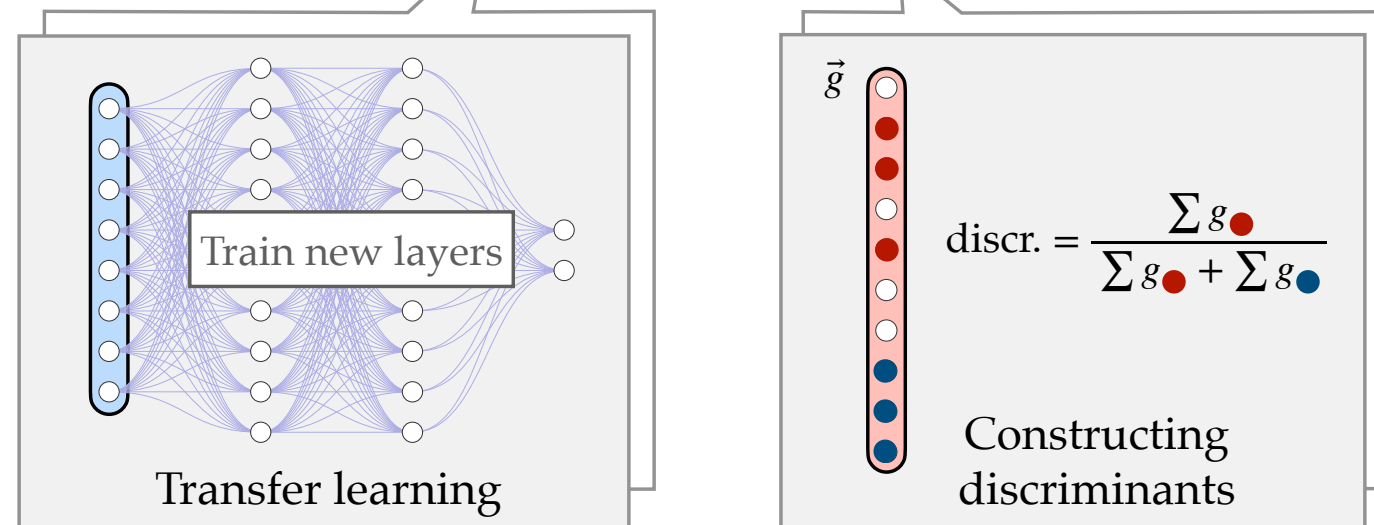
---

# Using Sophon

## (a) Pre-training



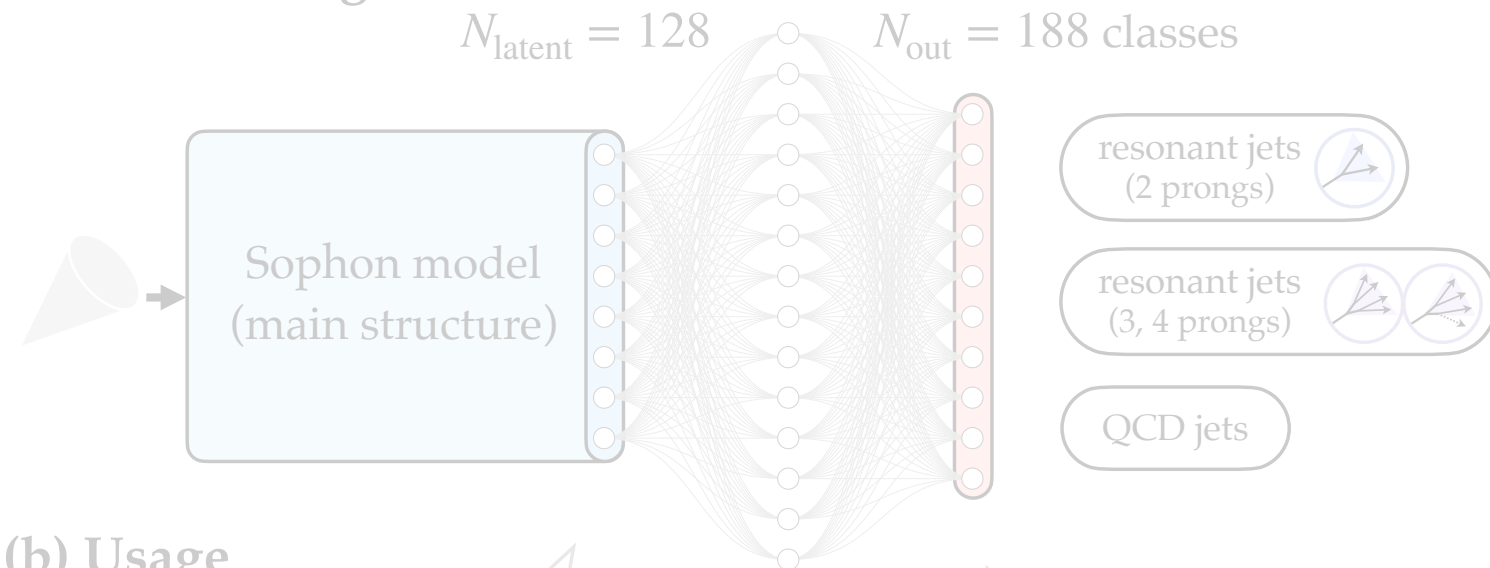
## (b) Usage



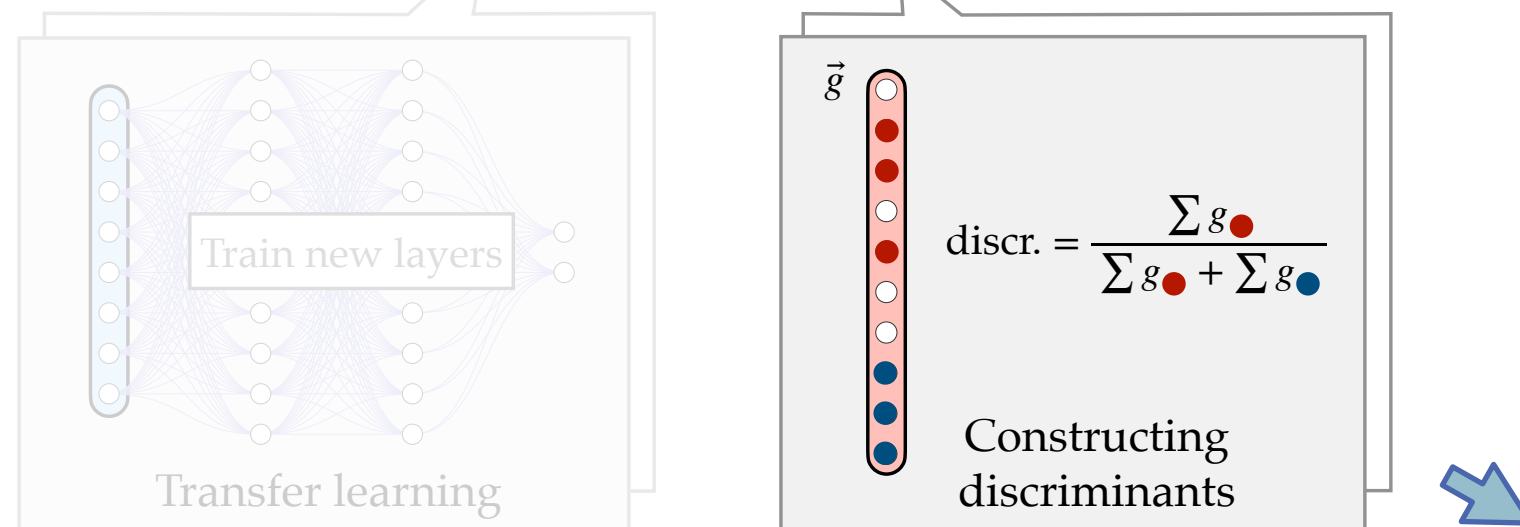


# Using Sophon

(a) Pre-training



(b) Usage



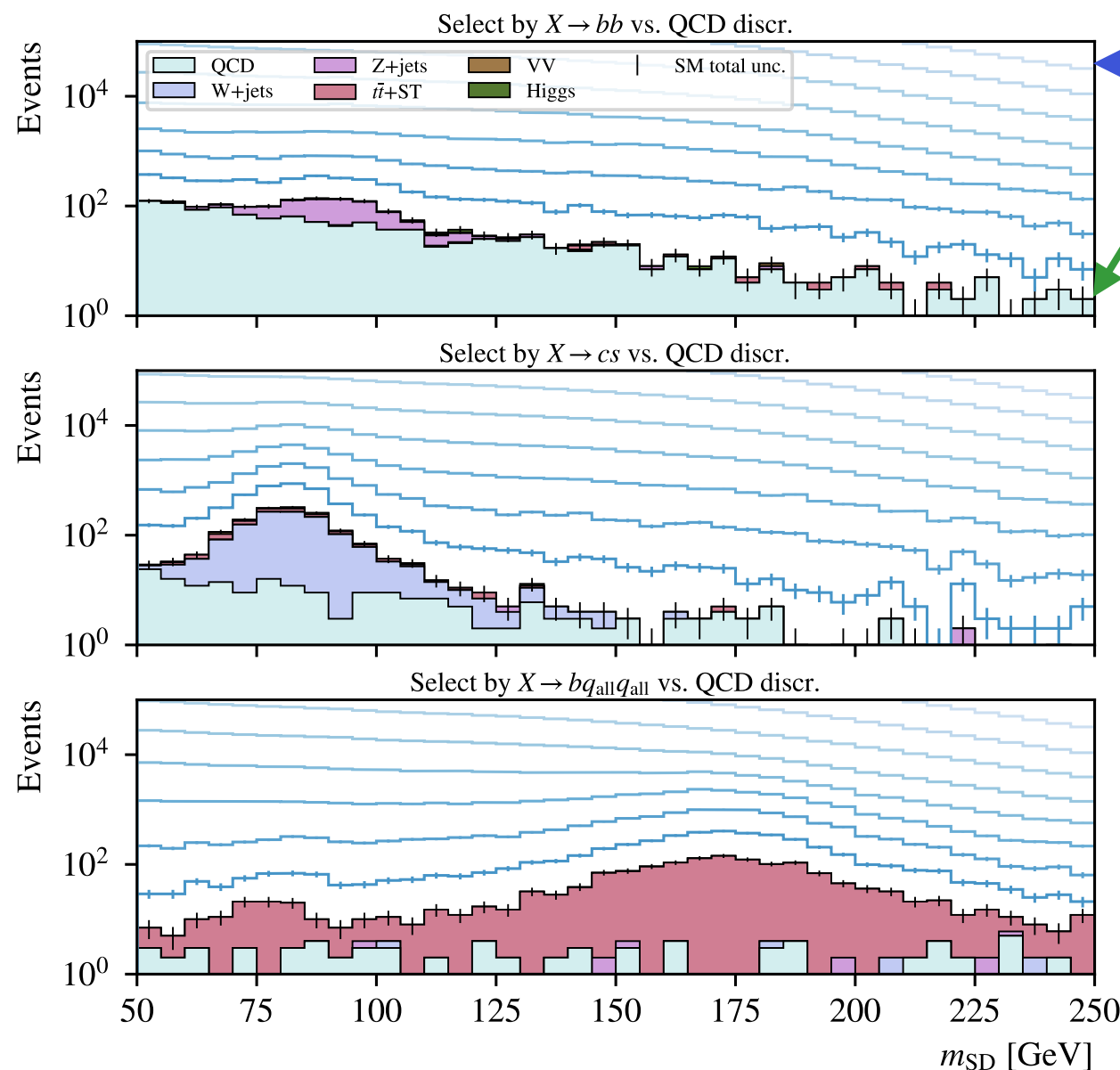
**Use it out of the box!**

Construct a dedicated discr.  
→ perform a bump hunt

# Can we rediscover the SM particles?

→ Simulate 40fb<sup>-1</sup> LHC collision events,  $\sqrt{s} = 13$  TeV, nPU=50

- ❖ focus on the large- $R$  jet trigger (triggered with  $\Sigma p_T$  threshold and trimmed mass)
- ❖ abundant QCD backgrounds
- ❖ **rediscover Z/W/t particles** simply from the large- $R$  jet’s **mass spectrum**



Without selection  
Select at eff. = 1e-4

○ Select by Sophon’s different discriminants

$$\text{discr} = \frac{g_A}{g_A + \sum_{l=1}^{27} g_{\text{QCD}_l}} \begin{cases} \textcircled{1}: A = \{bb\} \\ \textcircled{2}: A = \{cs\} \\ \textcircled{3}: A = \{ccb, ssb, qqb, bcs, bcq, bsq\} \end{cases}$$

# More heavy resonances

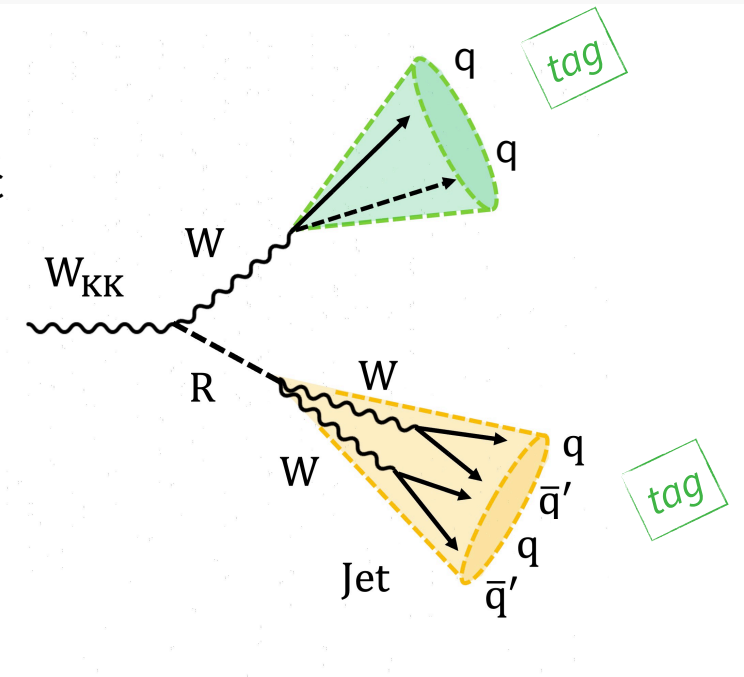
→ Consider triboson signal:

$W'$  ( $m_{W'} = 3 \text{ TeV}$ )  $\rightarrow W\phi$  ( $m_\phi = 400 \text{ GeV}$ )  $\rightarrow WWW$  (fully hadronic decays)

→ Optimize an event-level discr. from tagger discr.

$$\text{discr} = \sum_{\text{jet}=1,2} \frac{g_{A,\text{jet}}}{g_{A,\text{jet}} + \sum_{l=1}^{27} g_{\text{QCD},l,\text{jet}}} \quad (\text{sum for jets 1, 2})$$

$$A = \begin{cases} 0.3 \times \{cs, qq\} \\ + 0.1 \times \{ccss, qqcs, qqqq\} \\ + 0.6 \times \{ccs, ccq, ssc, ssq, qqc, qqqs, qqqs\} \end{cases}$$

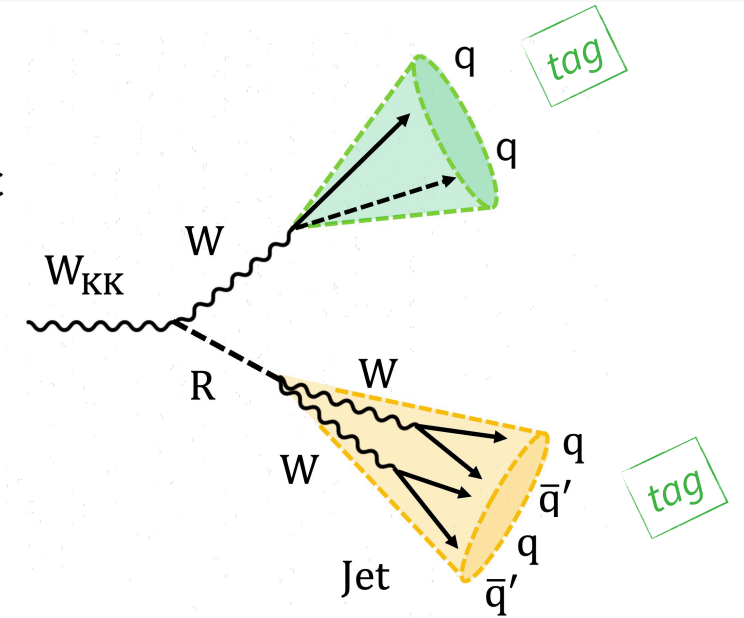


# More heavy resonances

→ Consider triboson signal:

$$W' (m_{W'} = 3 \text{ TeV}) \rightarrow W\phi (m_\phi = 400 \text{ GeV}) \rightarrow WWW \text{ (fully hadronic decays)}$$

→ Optimize an event-level discr. from tagger discr.



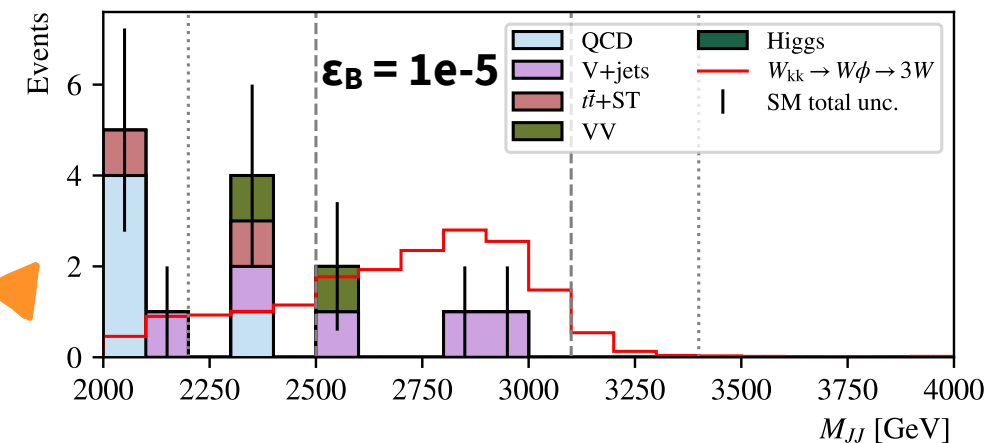
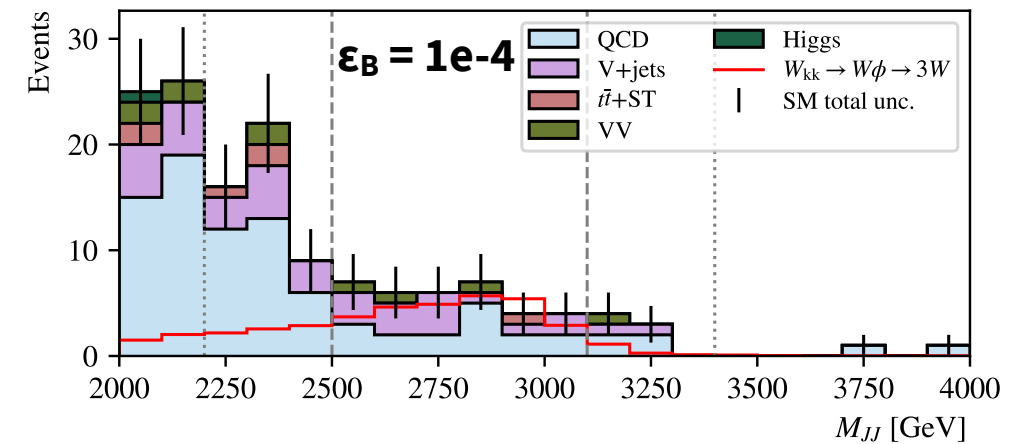
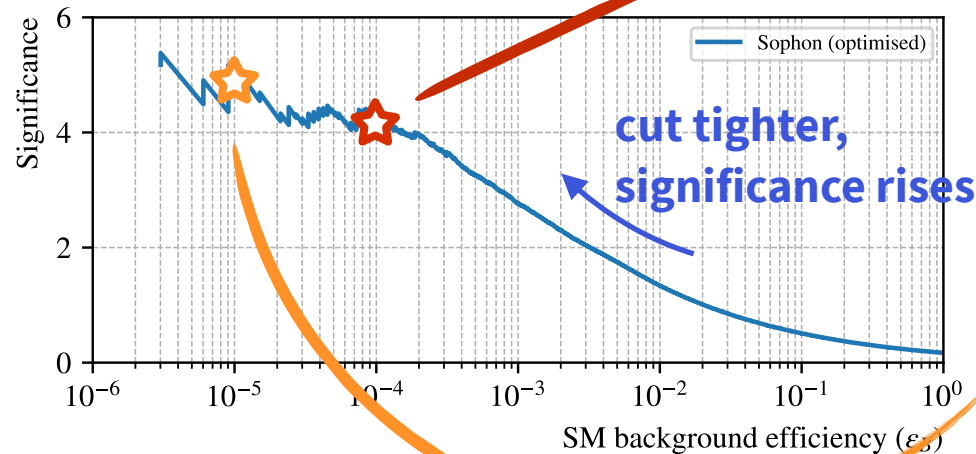
$$\text{discr} = \sum_{\text{jet}=1,2} \frac{g_{A,\text{jet}}}{g_{A,\text{jet}} + \sum_{l=1}^{27} g_{\text{QCD},l,\text{jet}}} \quad (\text{sum for jets 1, 2})$$

$$A = \begin{cases} 0.3 \times \{cs, qq\} \\ + 0.1 \times \{ccss, qqcs, qqqq\} \\ + 0.6 \times \{ccs, ccq, ssc, ssq, qqc, qqs, qqq\} \end{cases}$$

Search significance

$$Z = \sqrt{2((s+b)\log(1+s/b) - s)}$$

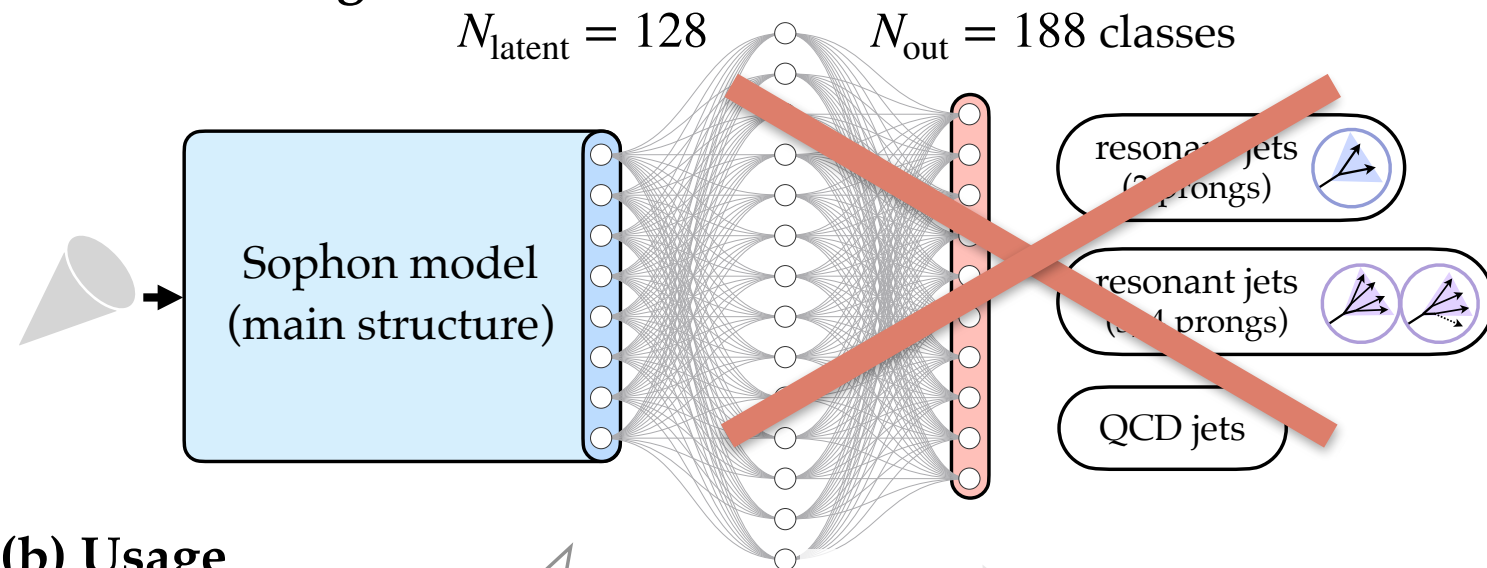
in dijet inv. mass window  
2500–3100 GeV



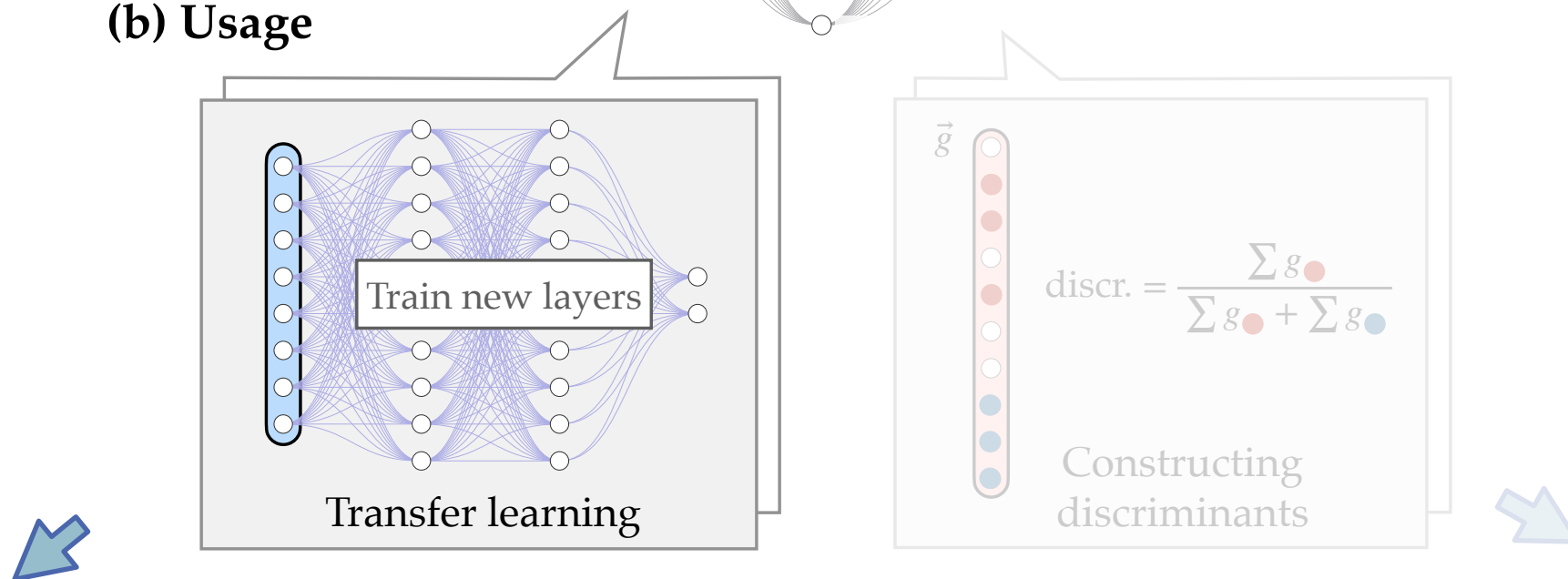


# Sophon's transfer learning

(a) Pre-training



(b) Usage



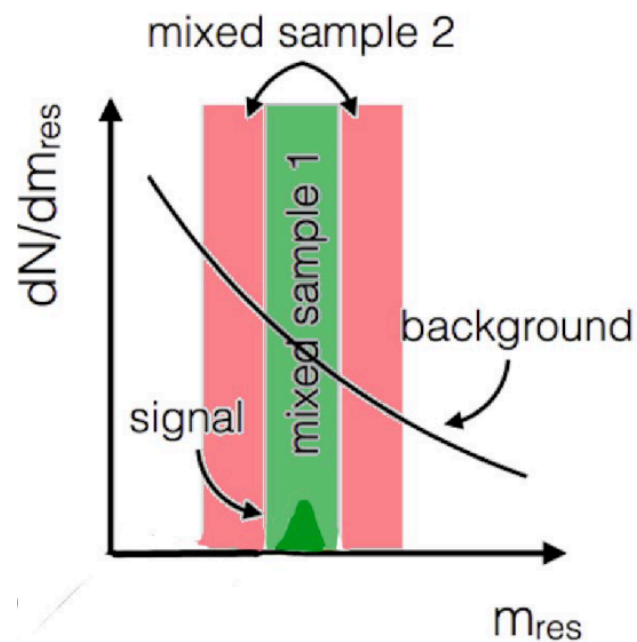
Use it out of the box!

- Transfer to uncovered tagging scenarios...
- facilitate anomaly detection (weakly-supervised, autoencoder)...
- *more potential to unlock!*

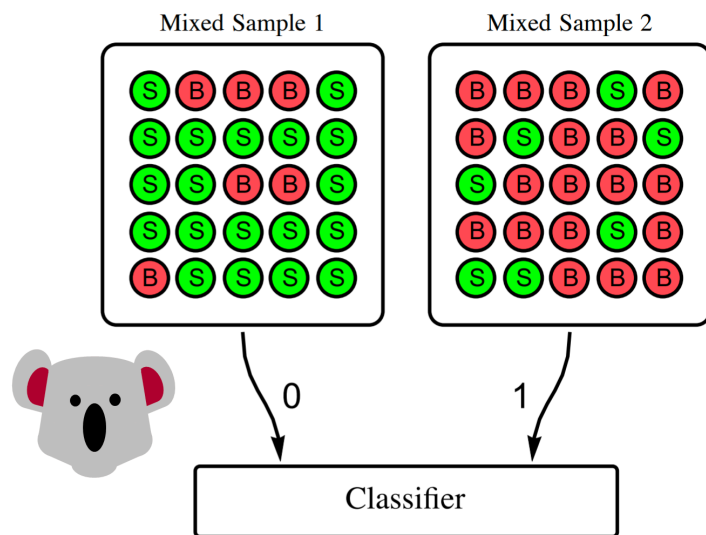
Construct a dedicated discr.  
→ perform a bump hunt

# Background: anomaly detection in weakly-supervised approach

*JHEP 10 (2017) 174*



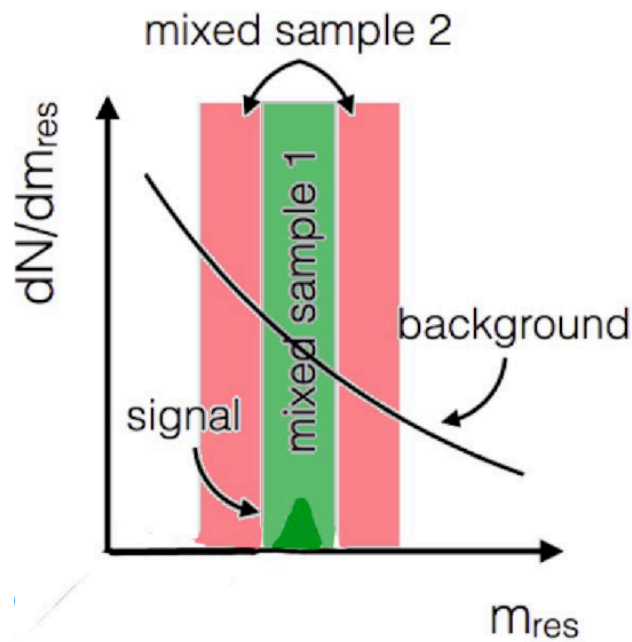
- Recall the early work: CWoLa (classification without labels) Hunting
- ❖ allow to detect anomalies purely from data
  - ❖ train a classifier for mass window vs mass sideband (mixed sample 1 vs 2)
  - ❖ many improved approaches in recent years → very active field



Equivalent effect for training **S** vs **B**

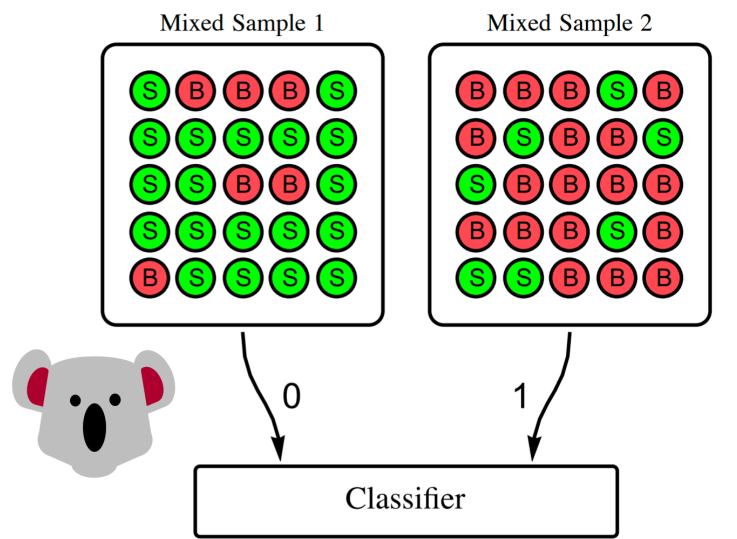
# Background: anomaly detection in weakly-supervised approach

[JHEP 10 \(2017\) 174](#)

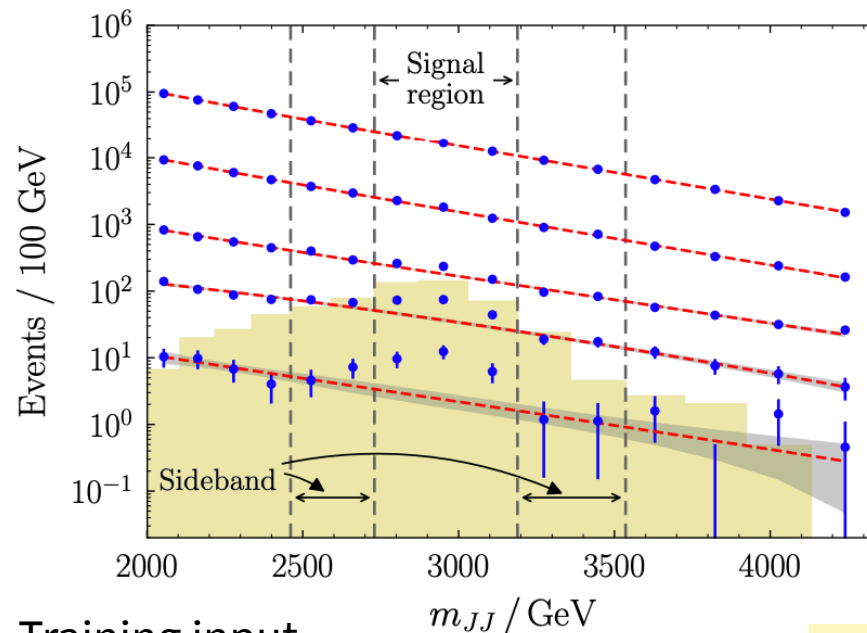


→ Recall the early work: CWoLa (classification without labels) Hunting

- ❖ allow to detect anomalies purely from data
- ❖ train a classifier for mass window vs mass sideband (mixed sample 1 vs 2)
- ❖ many improved approaches in recent years → very active field



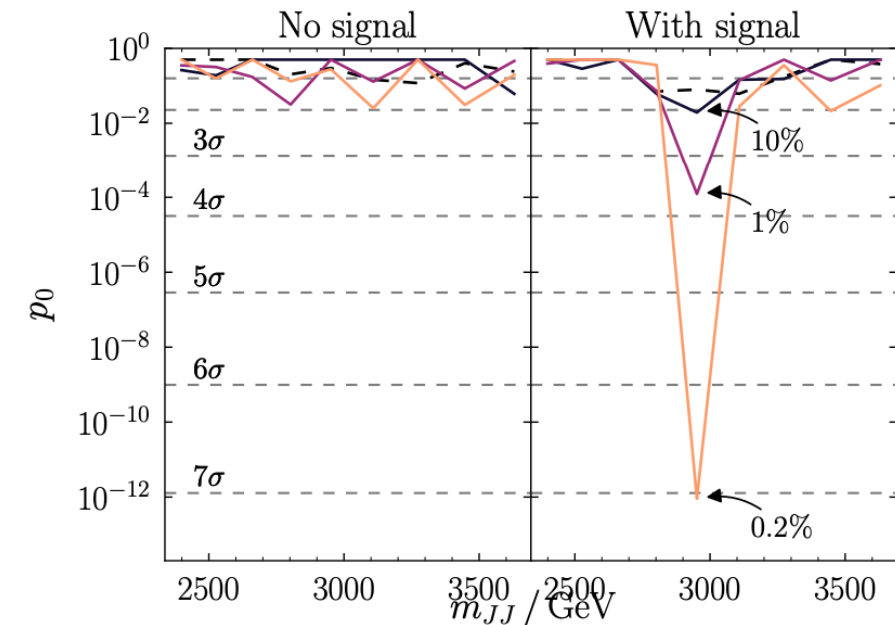
Equivalent effect for training **S** vs **B**



Training input

$$m_J, \sqrt{\tau_1^{(2)} / \tau_1^{(1)}}, \tau_{21}, \tau_{32}, \tau_{43}, n_{\text{trk}},$$

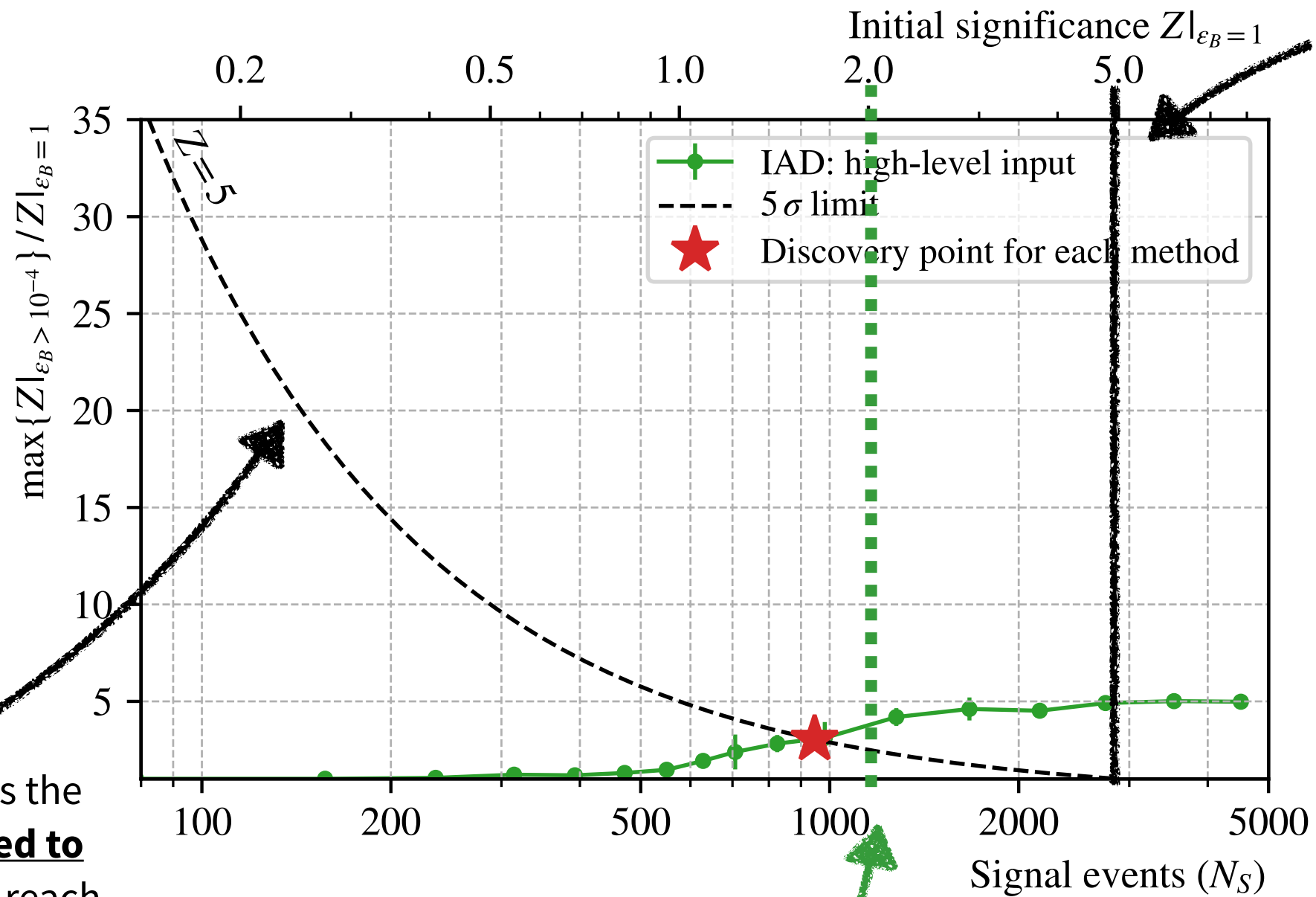
can discover  $W' \rightarrow W\phi \rightarrow WWW$  signals  
see  $2\sigma \rightarrow 7\sigma$  improvement



[PRL, 121 \(2018\) 24, 241803](#)

[PRD, 99 \(2019\) 1, 014038](#)

# Dijet search capabilities



“If signal events reach this point, **with initial Z=5**, then we have already discovered the signal without needing to make a cut”

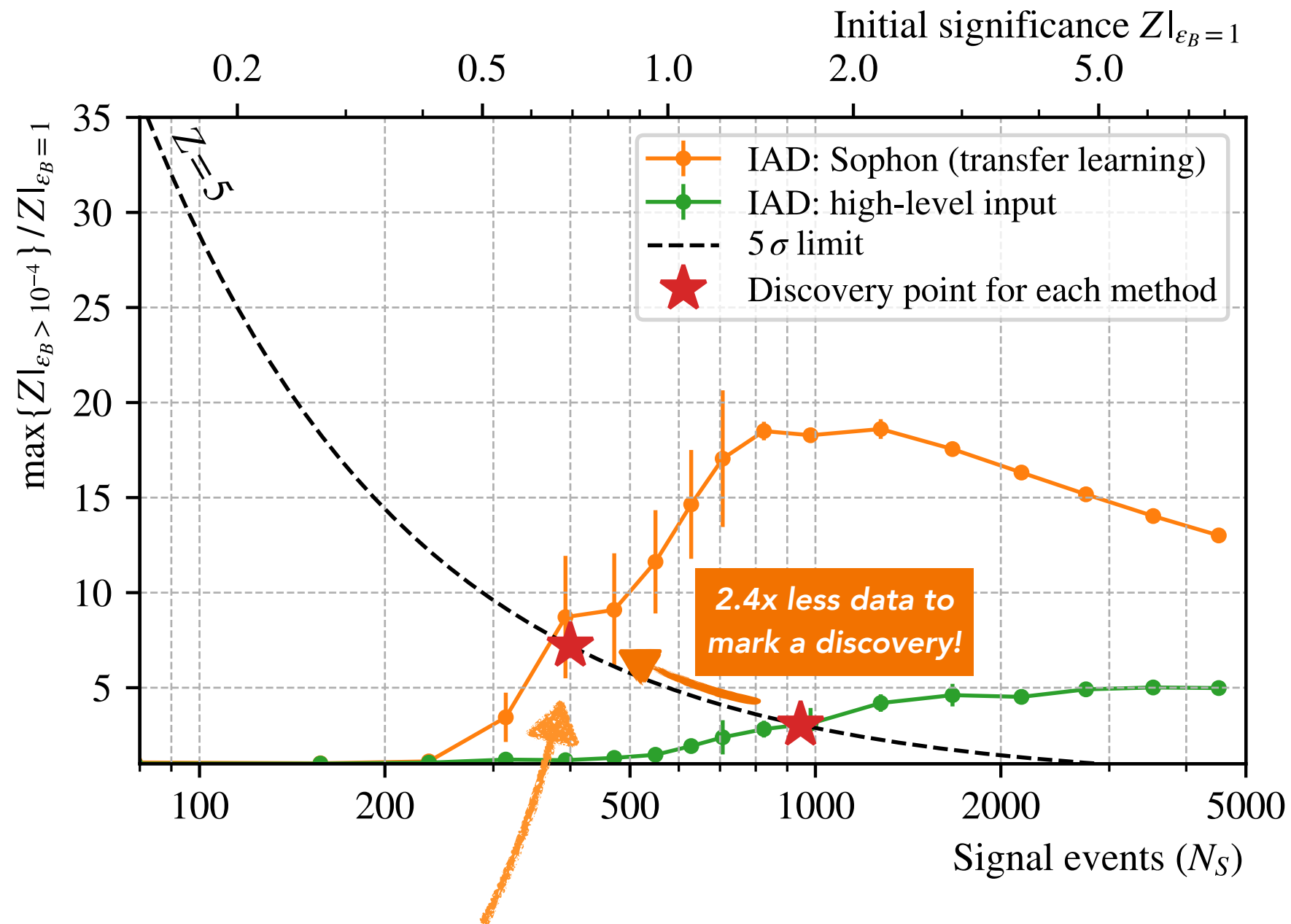
“How much does the **significance need to be increased** to reach the  $5\sigma$  discovery”

**a similar  $2\sigma \rightarrow 7\sigma$  is reached** with conventional AD approach; roughly reproduce the result in

[PRL, 121 \(2018\) 24, 241803](#)



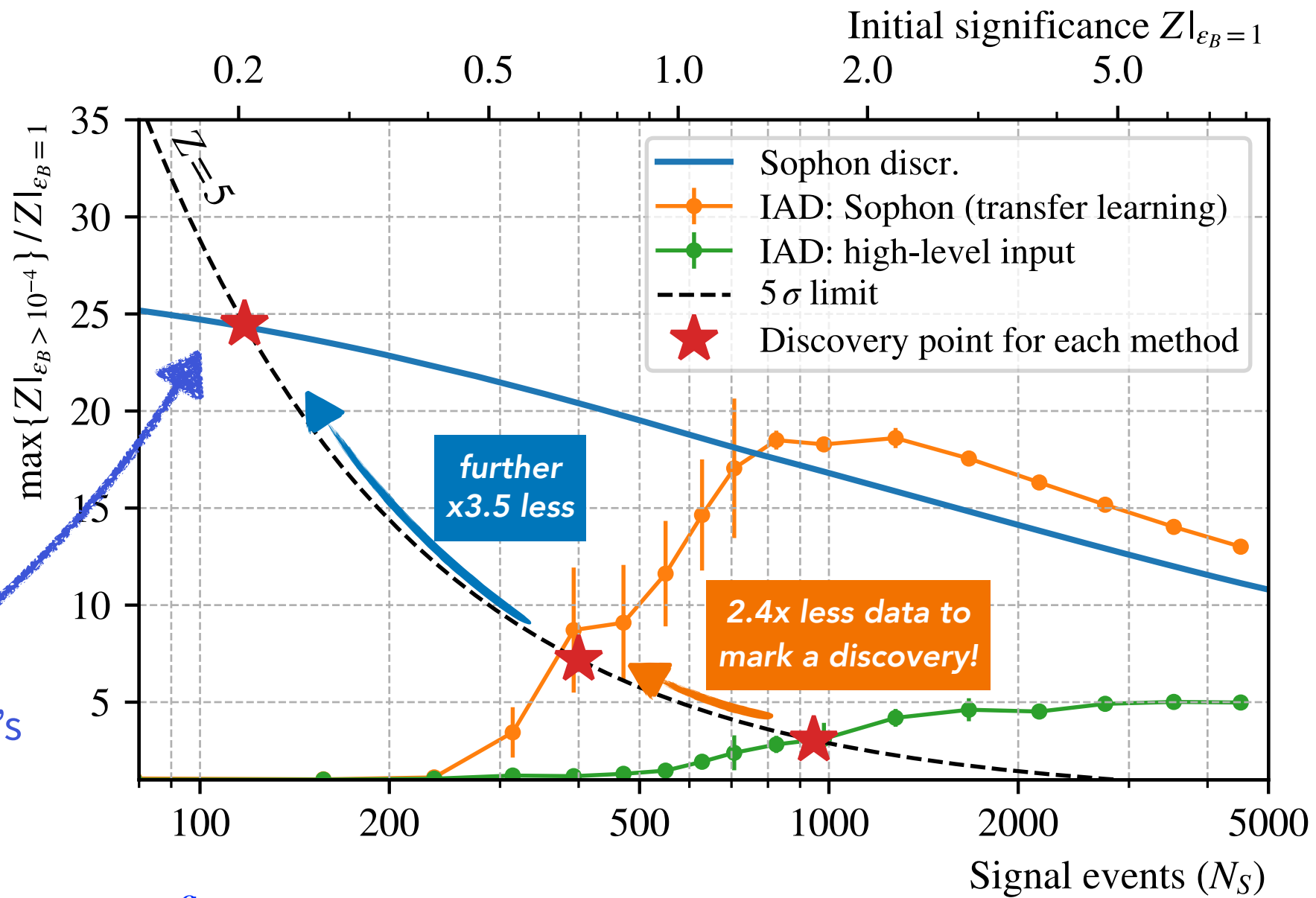
# Dijet search capabilities



Combining Sophon’s transfer learning (using Sophon’s “knowledge”) with AD marks a success

- More sensitive at low signal injection (**even starting at  $\sim 0.6\sigma$** )
- Much improved S vs B distinguishability than using high-level input

# Dijet search capabilities



using Sophon's constructed discriminant

$$\text{discr} = \sum_{\text{jet}=1,2} \frac{g_{A,\text{jet}}}{g_{A,\text{jet}} + \sum_{l=1}^{27} g_{\text{QCD}_l,\text{jet}}}$$

$$A = \begin{cases} 0.3 \times \{cs, qq\} \\ + 0.1 \times \{ccss, qqcs, qqqq\} \\ + 0.6 \times \{ccs, ccq, ssc, ssq, qqc, qqs, qqq\} \end{cases}$$

# CMS’s path to develop Global Particle Transformer

## Philosophy to develop **Global Particle Transformer (GloParT)** in CMS

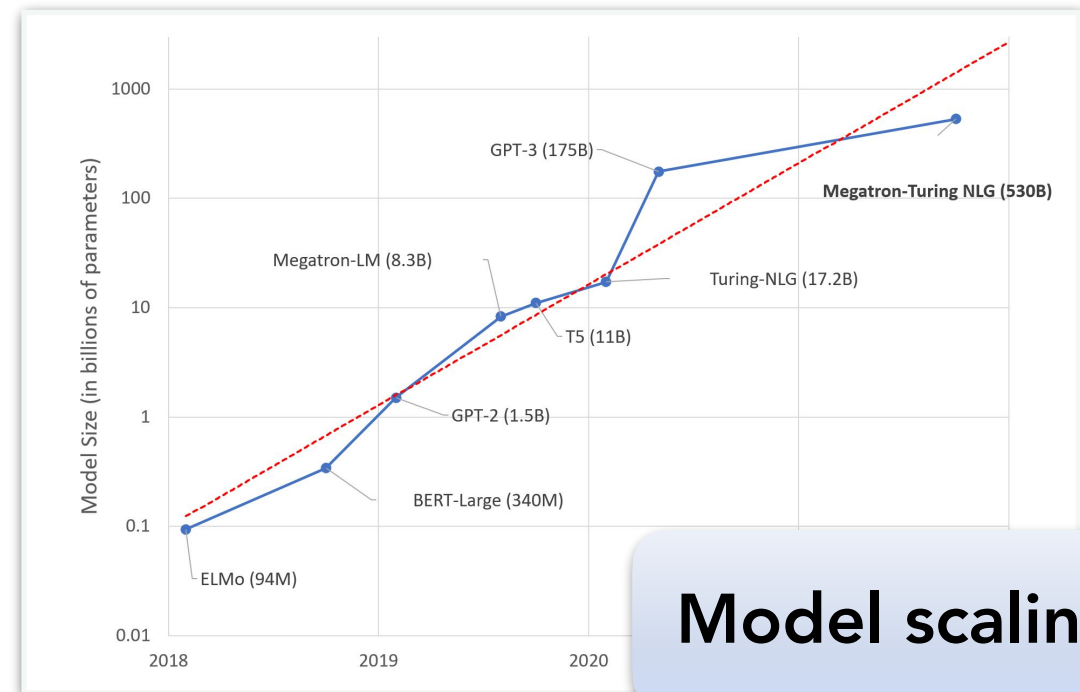
**Good probability density estimators**

- What is  $p$ ? - the “differential cross section” of a process  $A$  on very high-dim space
- discriminating process  $A$  vs.  $B$ : estimate  $p_A(\mathbf{x})/p_B(\mathbf{x})$  as best as we can
- need a model to **cover a variety of processes  $A, B, C, D, \dots$**

$A \rightarrow BC$	$B = \text{SM}$									$B = \text{BSM}$	
	$e$	$\mu$	$\tau$	$q/g$	$b$	$t$	$\gamma$	$Z/W$	$H$		
$C = \text{SM}$	$e$	$Z'$	$\tilde{R}$	$\tilde{R}$	$LQ$	$LQ$	$LQ$	$L^*$	$L^*$	$L^*$	Many
	$\mu$		$Z'$	$\tilde{R}$	$LQ$	$LQ$	$LQ$	$L^*$	$L^*$	$L^*$	
	$\tau$			$Z'$	$LQ$	$LQ$	$LQ$	$L^*$	$L^*$	$L^*$	
	$q/g$				$Z'$	$W'$	$T'$	$Q^*$	$Q^*$	$Q'$	
	$b$					$Z'$	$W'$	$Q^*$	$Q^*$	$B'$	
	$t$						$Z'$	$Q^*$	$T'$	$T'$	
	$\gamma$							$H$	$H$	$Z_{KK}$	
	$Z/W$							$H$	$H$	$H^\pm/A$	
	$H$									$H$	
	$C = \text{BSM}$	Consider just the di-object search for resonant $A \rightarrow B C$									

J.Kim *et al.* JHEP  
04 (2020) 30  
[1907.06659](https://arxiv.org/abs/1907.06659)

**Generalization ability**



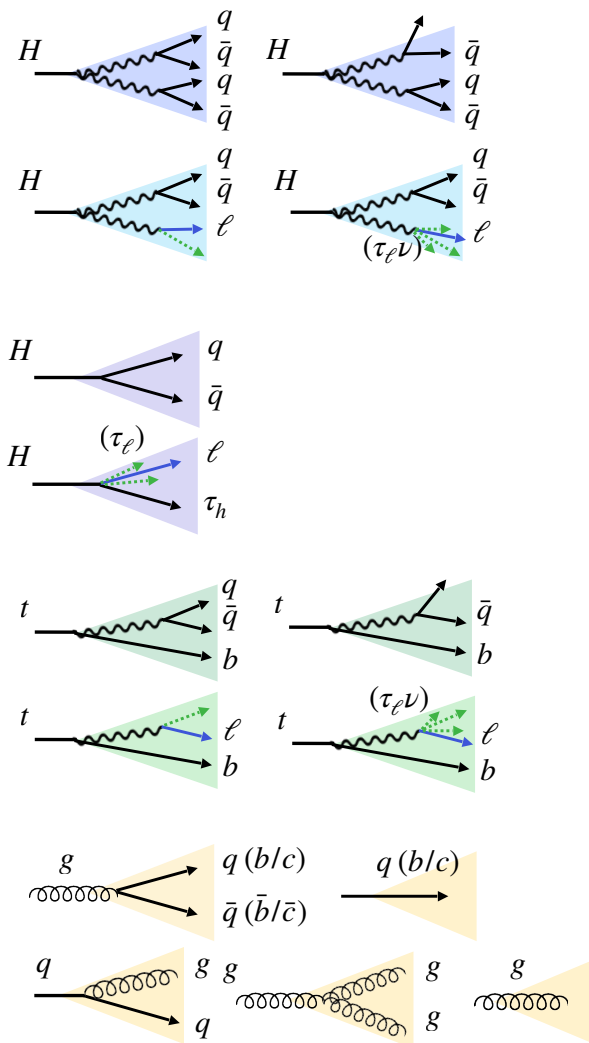
**Model scaling up**

- one upstream pre-training, broad downstream applicability

# CMS’s path to develop Global Particle Transformer

[CMS-PAS-HIG-23-012](#)

Process	Final state/ prongness	heavy flavour	# of classes
H→VV (full-hadronic)	qqqq	0c/1c/2c	3
	qqq		3
H→WW (semi-leptonic)	eνqq	0c/1c	2
	μνqq		2
	τ <sub>e</sub> νqq		2
	τ <sub>μ</sub> νqq		2
	τ <sub>h</sub> νqq		2
H→qq		bb	1
		cc	1
		ss	1
		qq (q=u/d)	1
H→ττ	τ <sub>e</sub> τ <sub>h</sub>		1
	τ <sub>μ</sub> τ <sub>h</sub>		1
	τ <sub>h</sub> τ <sub>h</sub>		1
t→bW (hadronic)	bqq	1b + 0c/1c	2
	bq		2
t→bW (leptonic)	bēν	1b	1
	bμν		1
	bτ <sub>e</sub> ν		1
	bτ <sub>μ</sub> ν		1
	bτ <sub>h</sub> ν		1
QCD		b	1
		bb	1
		c	1
		cc	1
		others (light)	1



The early version (**GloParT stage-1**) has been released with the HH→bbWW search

❖ **GloParT stage-1**: a fatjet tagger for **37-category classification**

- bbWW is the second work in the series of “boosted HH search”
- tagging **boosted H→WW→4q signature** for the first time
- set a tight constraint to  $\kappa_{2V}$



# CMS's path to develop Global Particle Transformer

CMS-PAS-HIG-23-012

The early version (GloParT stage-1) has been released with the  $HH \rightarrow bbWW$  search

Process	Final state/prongness	heavy flavour	# of classes
$H \rightarrow WW$ (full-hadronic)	qqqq		3
	evqq		2
$H \rightarrow WW$ (semi-leptonic)	$\mu\nu qq$		2
	$\tau\nu qq$		2
$H \rightarrow qq$		cc	1
		ss	1
		qq (q=u/d)	1
$H \rightarrow \tau\tau$	$T_e T_h$		1
	$T_\mu T_h$		1
	$T_h T_h$		1
$t \rightarrow bW$ (hadronic)	bqq	1b + 0c/1c	2
	bq		2
$t \rightarrow bW$ (leptonic)	b $\nu$		1
	b $T_e \nu$	1b	1
	b $T_\mu \nu$		1
	b $T_h \nu$		1

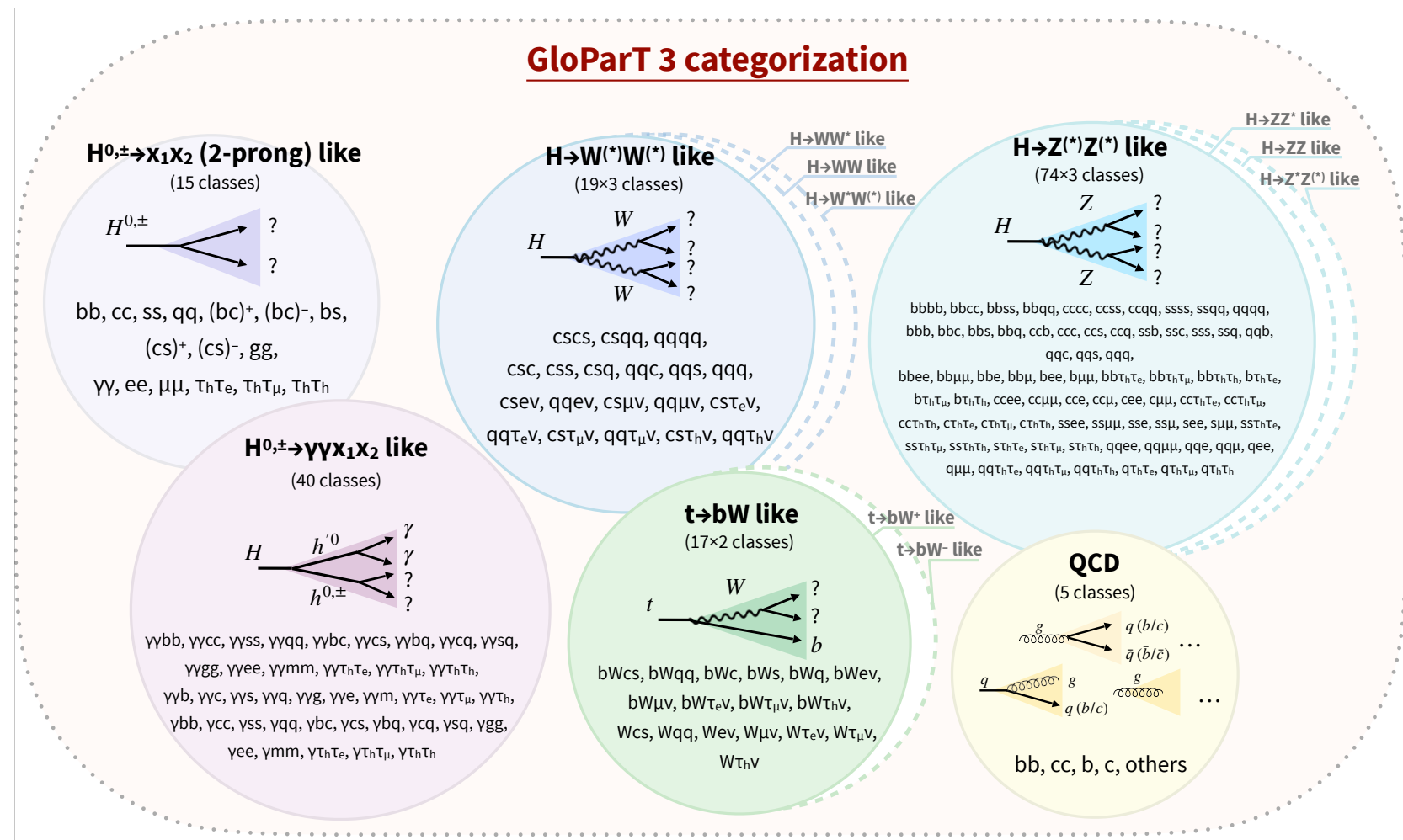
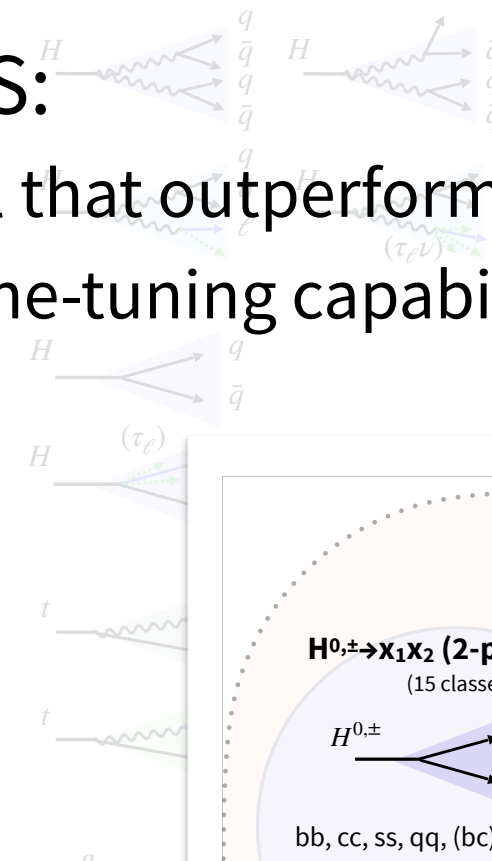
## GloParT-3 in CMS:

- ✓ A universal model that outperforms all existing taggers across existing tasks
- ✓ Manifest strong fine-tuning capability in various downstream tasks

## A more powerful model is now in place

- the 3<sup>rd</sup> version of the Global Particle Transformer (**GloParT-3**) now integrated in CMS

- 750 nodes in total



# Discussion: Implications to ATLAS/CMS experiments?

- “Sophon/GloParT methodology” releases a lot of new opportunities for future LHC experiments
  - ❖ it creates a **“global large- $R$  jet tagger”** → bring benefits of the advanced NN to ~all hadronic final-state searches
  - ❖ **also viewed as a pre-trained jet model**: a base model tailored for a broad range of LHC analyses
- How to use the experimental version of the Sophon model?
  - ❖ used in conventional analyses: except for some well-calibrated nodes, the major challenge will be the calibration of peculiar signals (not easy to find proxies)
  - ❖ **analyses that only use data (simulation-free)**: develop discriminants dedicated to different signals → cut tight on the data events → peak finding on some mass observable (single jet / di-jet / jet+lepton...)
    - **could be helpful in broadly searching for BSM resonance!**
  - ❖ anomaly detection: weakly-supervised approaches / further improvements?  
(see backup for recent ATLAS/CMS results)

# Discussion: JetClass-II and Sophon

[arXiv:2405.12972](https://arxiv.org/abs/2405.12972)

- Developed the **JetClass-II** dataset and the **Sophon** model
- **JetClass-II** [[Hugging Face dataset](#)] covers more comprehensive phase spaces and can be a good playground to develop future foundation models
  - ❖ can be used to train models for various jet-related tasks, e.g. jet classification, regression, generation or reconstruction...
  - ❖ its extensive phase space coverage and high statistics enable model developers to **focus on specific regions of interest**, or **work with the entire dataset**
  - ❖ generation details can be found in this [repository](#)
- The **Sophon model** [[Hugging Face](#)] can be helpful in delivering future LHC pheno research
  - ❖ optimizing sensitivity for dedicated searches/anomaly detection/novel paradigms
  - ❖ performing studies on the pheno dataset/model can inspire how we do real experiments at the LHC



# Discussion: Close-up & Future Fantasies...

## → Future foundation/base model at LHC:

- ❖ which dataset will it be trained on (data/simulation)?
- ❖ which training targets (generation/classification/embedding prediction)?
- ❖ maybe, most importantly, which goal do we want to achieve?

## → Some specific (maybe preliminary) points to discuss

- ❖ SSL or signature-oriented pre-training?
  - In HEP, we do not lack data labels; simulation will also reach expected accuracy in future
  - should we do SSL, or supervised pre-training exploring all GEN labels, or both?
- ❖ foundation model requires a “foundation dataset”; which will be the foundation dataset in HEP?
  - should cover a vast phase-space (beyond SM simulation/real data): produced by philosophies like JetClass-II?
  - should be in the most general form, e.g. independent of experiments.. - GEN-level data be the most suitable playground?
- ❖ The limit of scaling capabilities?
  - are we reaching the limit for classification tasks? Will fine-tuning from a large pre-trained model provide better results?

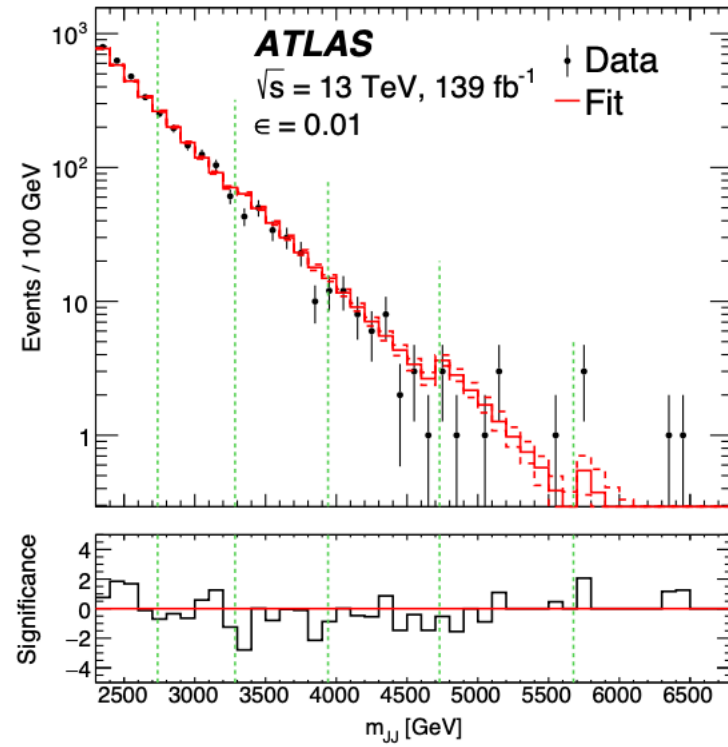
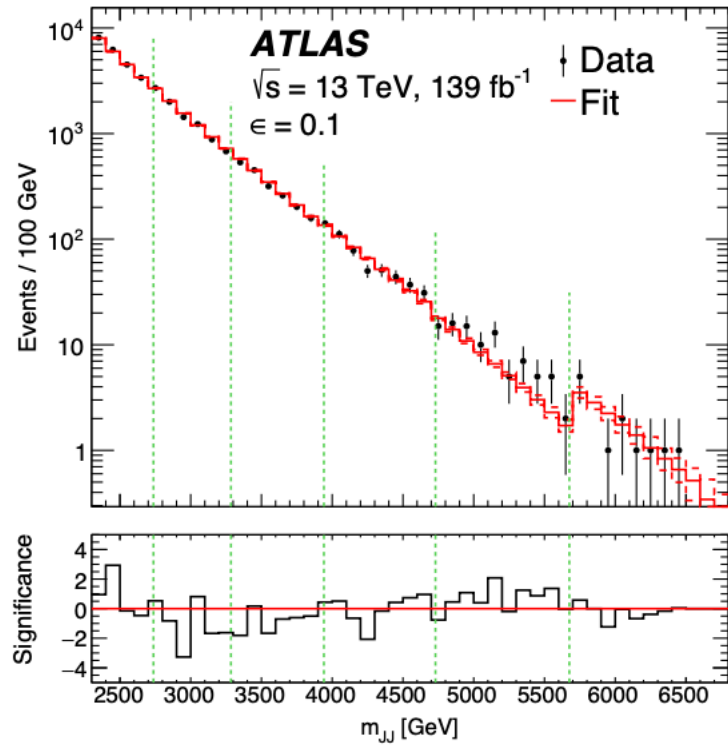


# Backup

---

# Recent ATLAS/CMS anomaly detection results

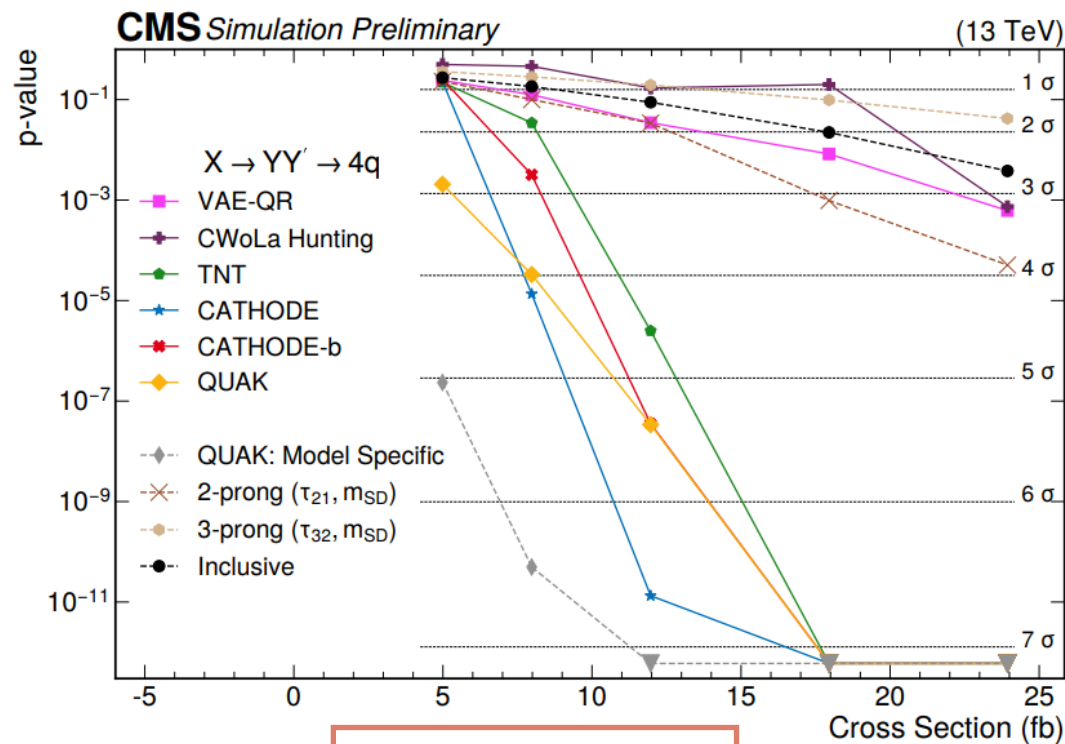
## Weakly supervised approach



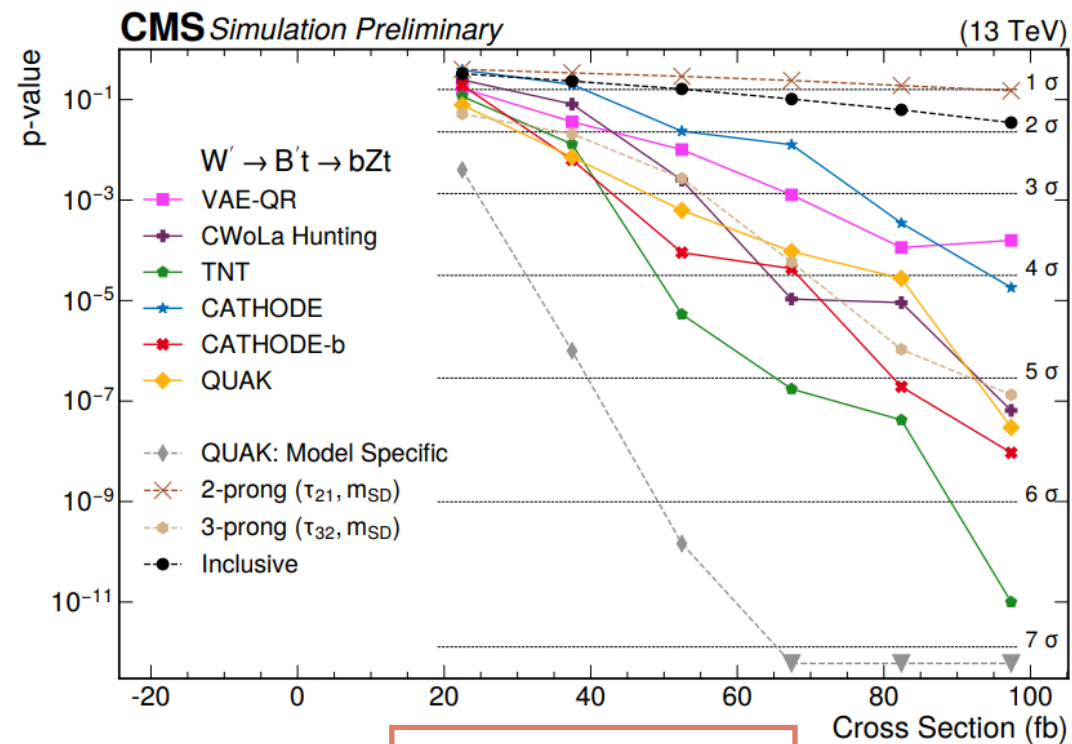
### Application in ATLAS

[PRL 125, 131801 \(2020\)](#)

Select on the trained weakly supervised classifier in each signal regions, and search for the peak



2-prong + 2-prong



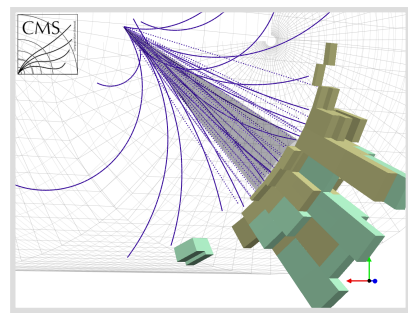
3-prong + 3-prong

[CMS-PAS-EXO-22-026](#)

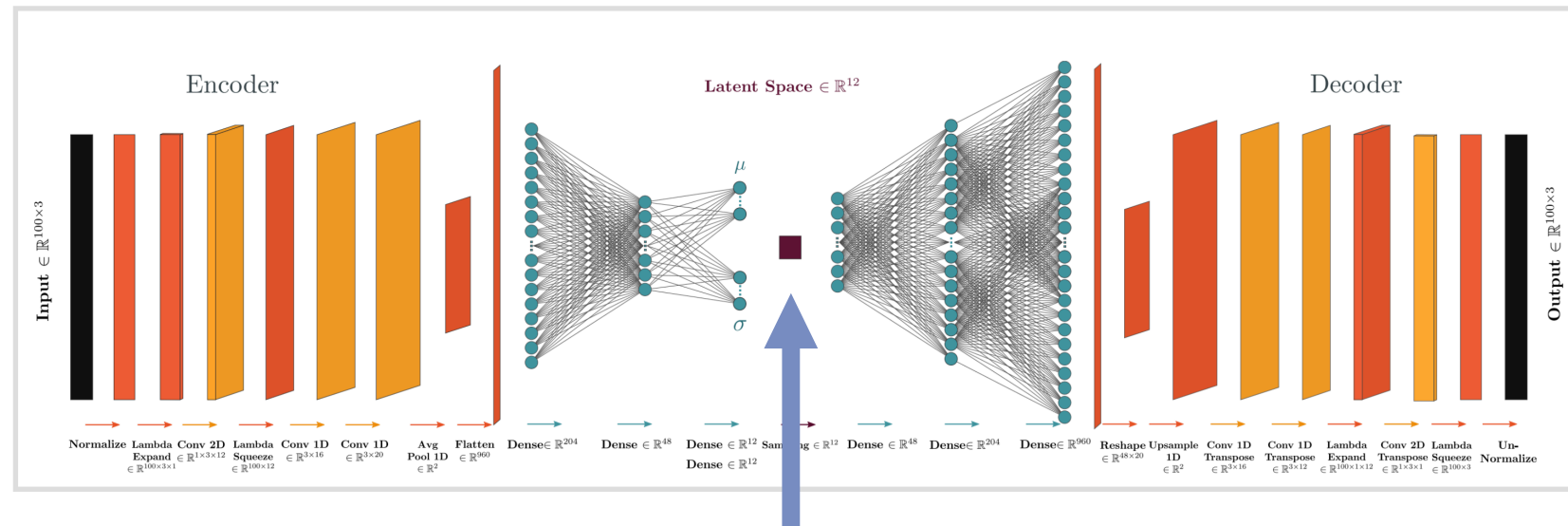
In a recent work, CMS systematically test all model-agnostic approaches to search for resonance

# Recent ATLAS/CMS anomaly detection results

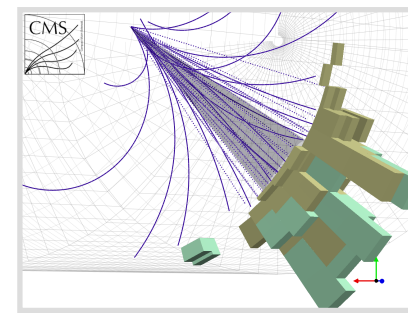
## Autoencoder approach



input jet



a compressed jet representation

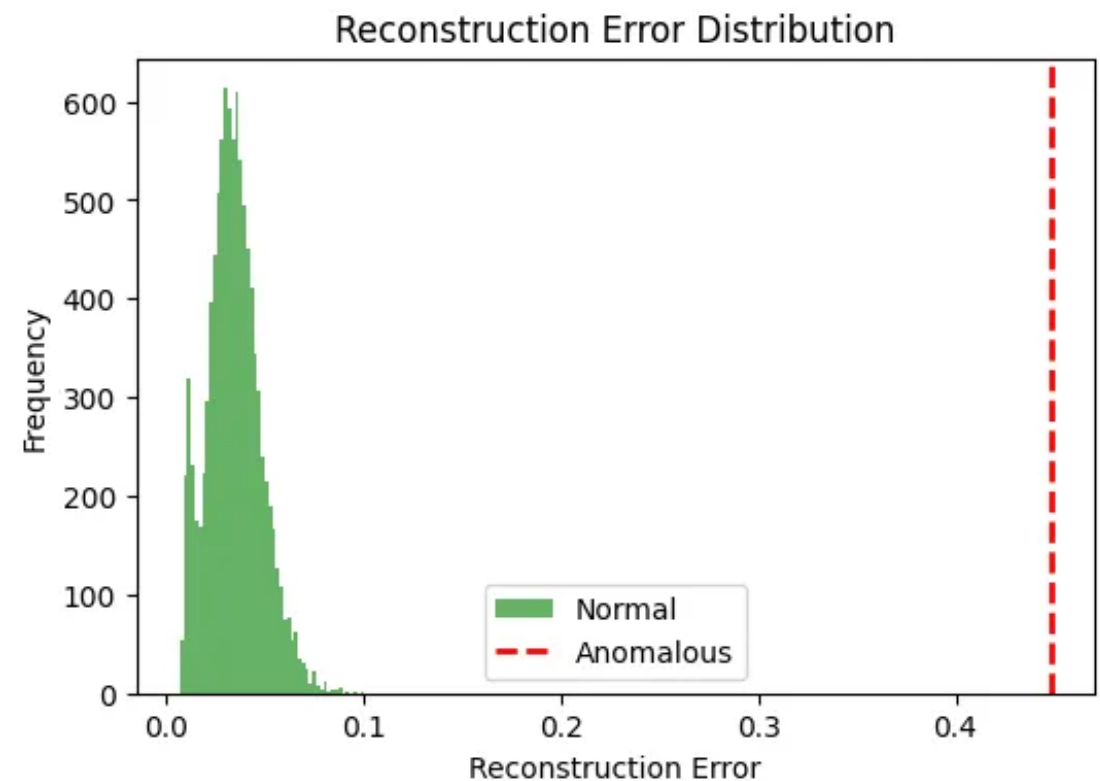


reconstructed jet

→ A view on (variational) autoencoder for anomaly detection

- ❖ Training on SM background jet → **anomalous jet will produce outlier latent scores** → make selection on the score

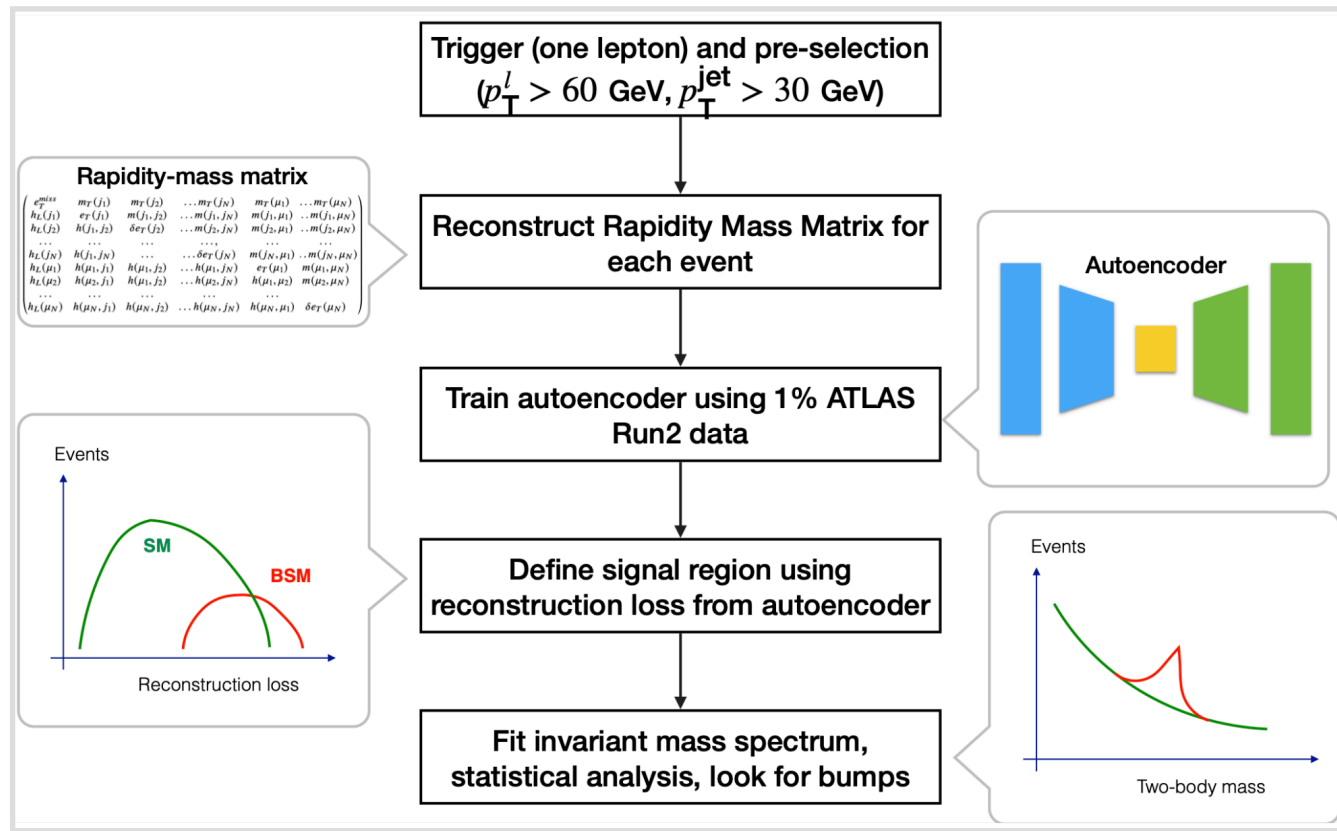
→ Use autoencoder for anomaly detection has industry basis



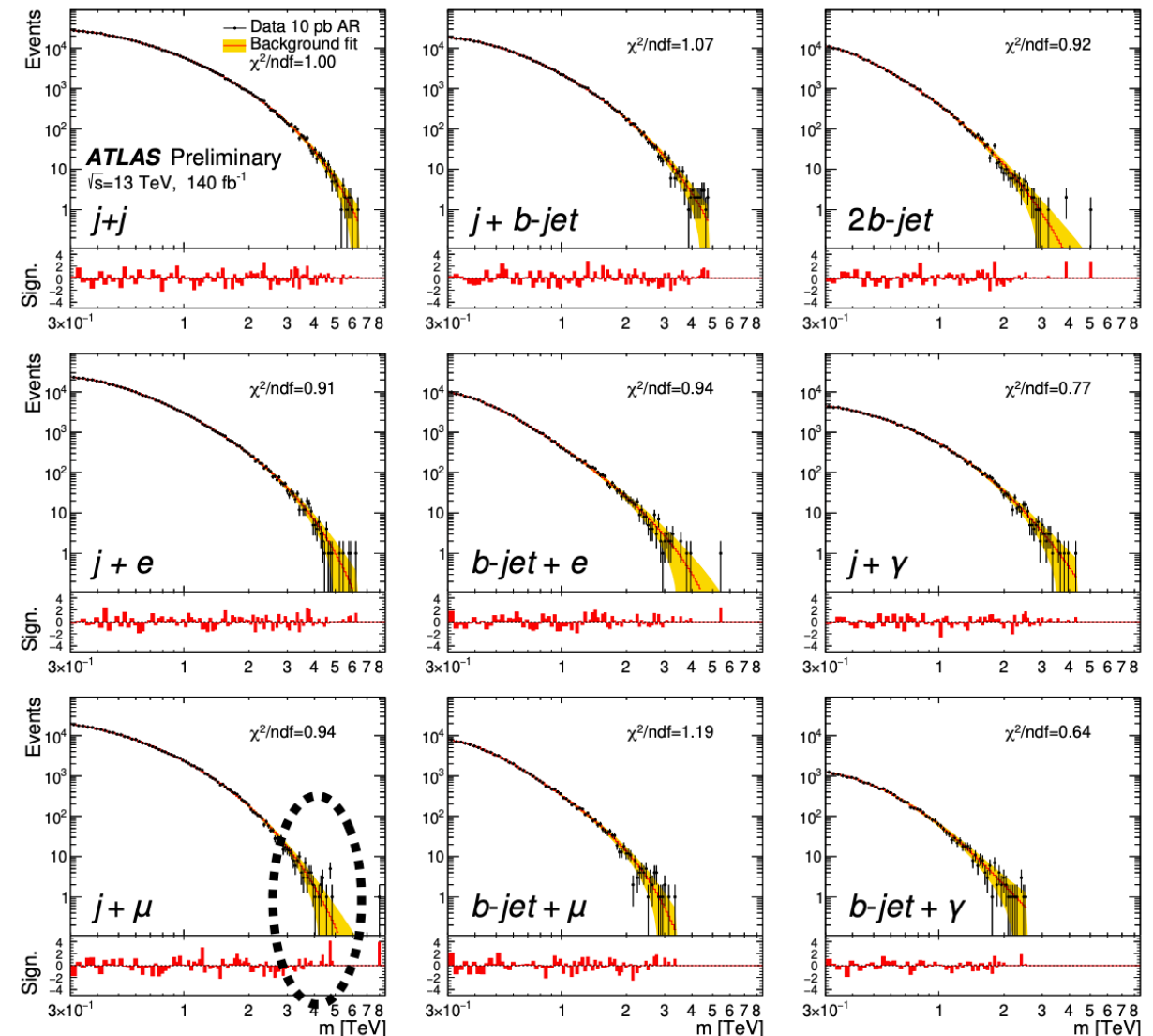


# Recent ATLAS/CMS anomaly detection results

## Autoencoder approach



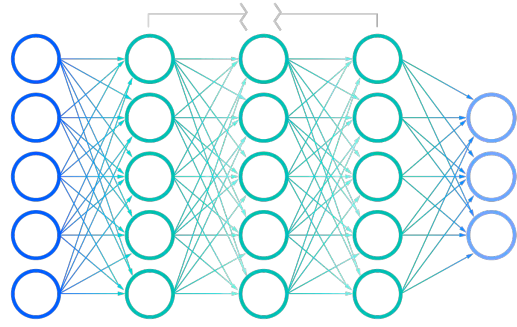
PRL 132 (2024) 081801



- ATLAS applies full-event-level anomaly detection
- Train “autoencoder” and select on the score
- Search in 9 invariant masses including di-jet, di-b-jet, with three anomaly regions

# Evolution of jet NNs

*feed-forward NN (high-level inputs)* ...▶... *1D/2D CNN, RNN (low-level inputs)* ...▶... *graph NN, Transformers (low-level inputs)* ...▶... ??



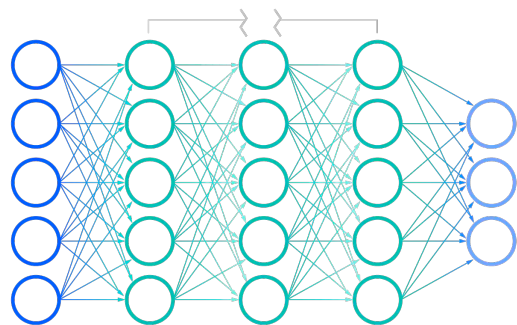
## Shallow networks

- ◆ Using high-level features directly as input to a shallow network



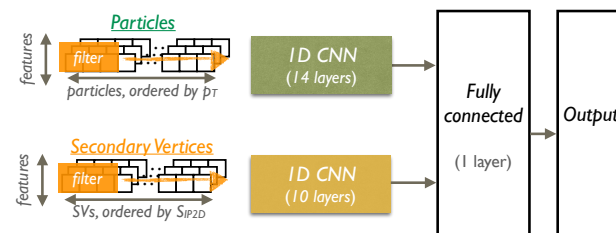
# Evolution of jet NNs

feed-forward NN (high-level inputs) ... 1D/2D CNN, RNN (low-level inputs) ... graph NN, Transformers (low-level inputs) ... ??



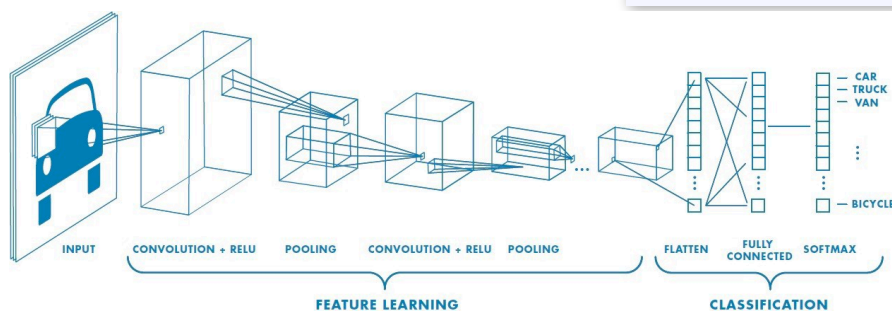
## Shallow networks

- ◆ Using high-level features directly as input to a shallow network

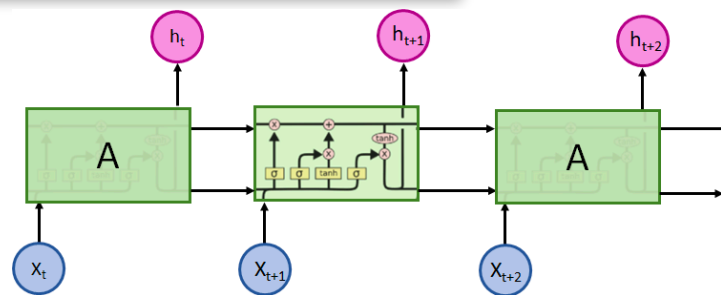


## Deep NN with low-level inputs

- ◆ Using particle-level features
- ◆ Input data structure determines the type of networks
  - jet as a *image* (fixed-grid data structure)
  - jet as a *sequence* → 1D CNN or RNN



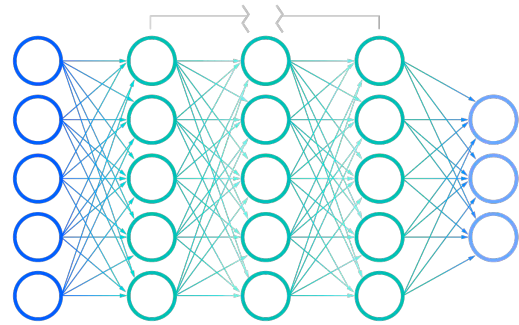
Typical CNN



Typical RNN

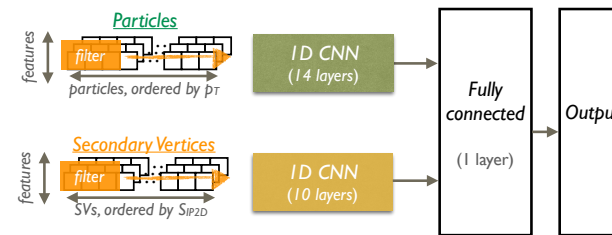
# Evolution of jet NNs

feed-forward NN (high-level inputs) ... 1D/2D CNN, RNN (low-level inputs) ... graph NN, Transformers (low-level inputs) ... ??



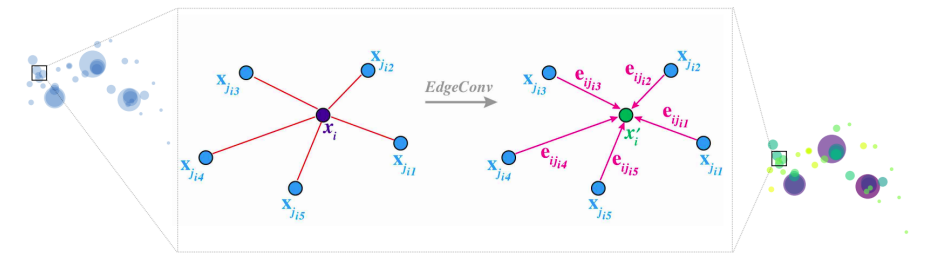
Shallow networks

- Using high-level features directly as input to a shallow network



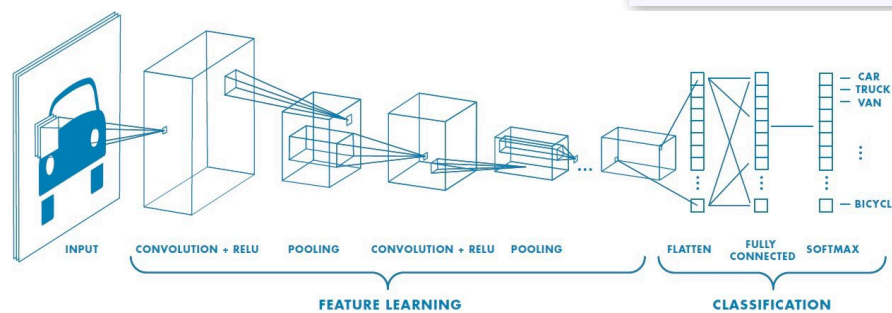
Deep NN with low-level inputs

- Using particle-level features
- Input data structure determines the type of networks
  - jet as a *image* (fixed-grid data structure)
  - jet as a *sequence* → 1D CNN or RNN

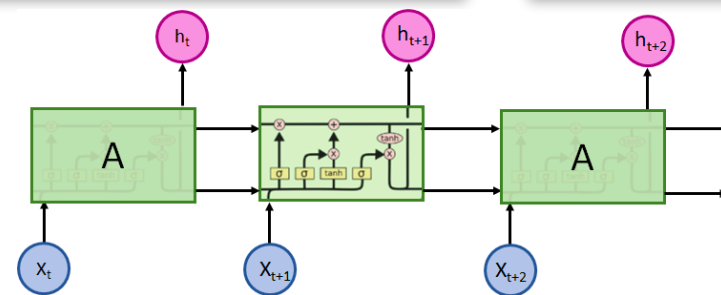


Graph structure

- Graph neural networks
  - treat a jet as a permutational-invariant set of particles (or, point cloud)
  - build “edges” between particles
- Transformer networks
  - modern architectural designs; like a full-connected graph



Typical CNN



Typical RNN



Typical graph

# GNNs and Transformers

- **Modern architectures done right:** (which types of DNNs better suit the particle-format data?)
  - ❖ **inductive bias:** particle-format data has their intrinsic symmetries
    - **permutational-invariant symmetry:** GNN is better than CNN/RNN; native Transformer (w/o positional encoding)
    - Lorentz symmetry: adding “pairwise particle masses” to input features
  - ❖ **let particles interact:**
    - “message passing” in GNNs and attention mechanism in Transformers
  - ❖ **scale better with data and model size**
    - Transformers!

