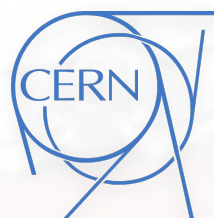




UC San Diego



Fermilab

Based on [arXiv:2405.12972](https://arxiv.org/abs/2405.12972)[\[Github\]](#) [\[Google Colab\]](#)

Sophon meets LHC: Accelerating resonance discovery via signature- oriented pre-training

Congqiao Li (李聪乔), *Peking University*

based on the work with our colleagues in the CMS Collaboration:

Antonios Agapitos¹, Jovin Drews², Javier Duarte³, Dawei Fu¹, Leyun Gao¹, Raghav Kansal³, Gregor Kasieczka²,
Louis Moureaux², Huilin Qu⁴, Cristina Mantilla Suarez⁵, Qiang Li¹

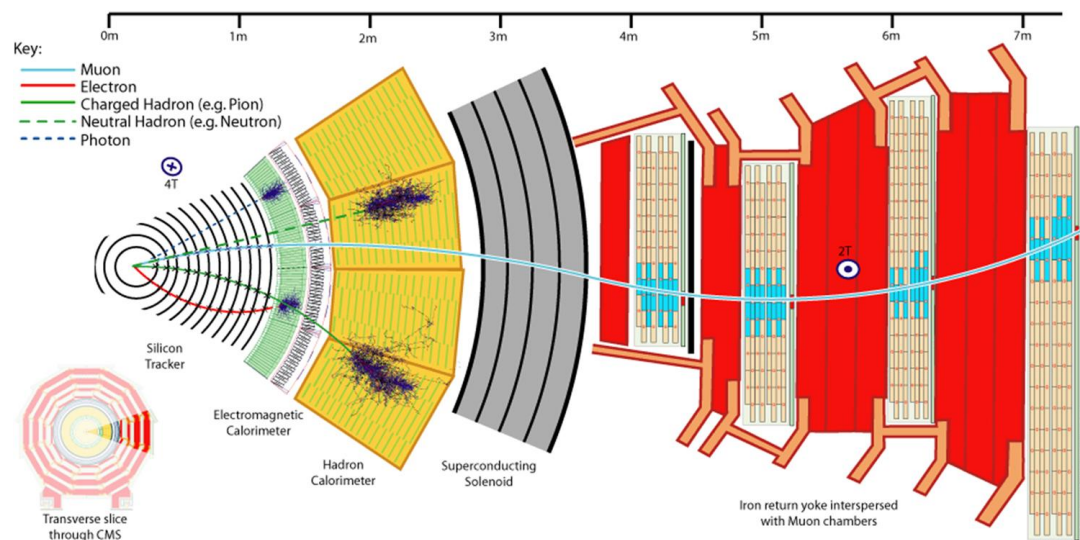
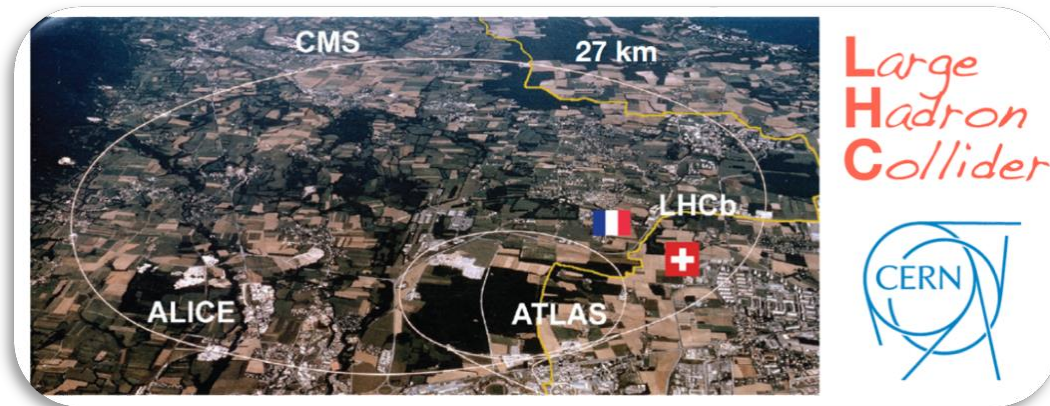
1) Peking U. 2) Hamburg U. 3) UC San Diego 4) CERN 5) FNAL

also thanks Yuzhe Zhao¹ for his contribution

Quantum Computing and Machine Learning Workshop 2024, Changchun

6 August, 2024

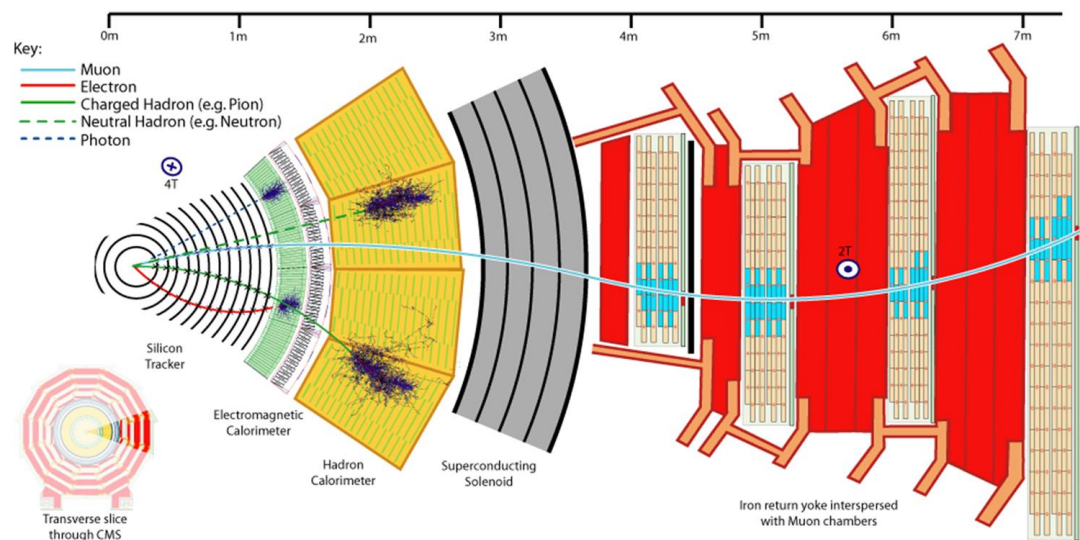
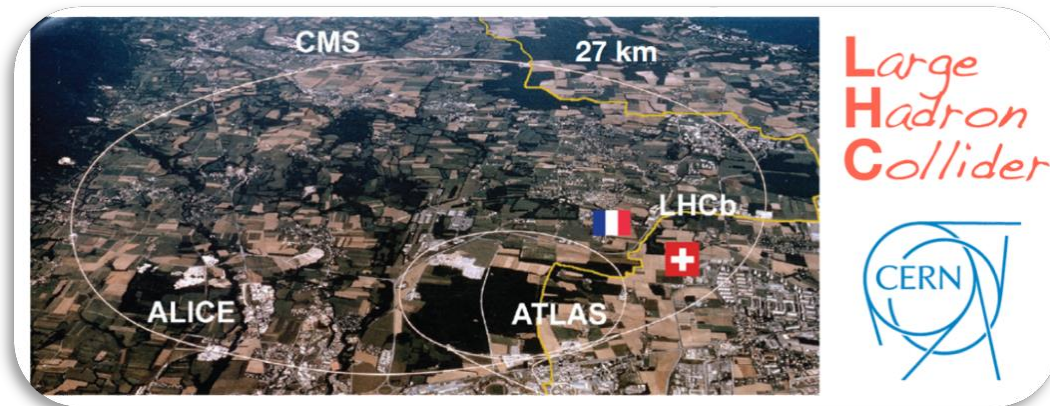
LHC physics × the era of deep learning



- The LHC: current world's largest particle collider
- Physics programs on LHC: reveal the fundamental theory of matters
- ATLAS and CMS: general-purpose detectors for precise SM measurement and **searching of BSM**

LHC physics × the era of deep learning

- The LHC: current world's largest particle collider
- Physics programs on LHC: reveal the fundamental theory of matters
- ATLAS and CMS: general-purpose detectors for precise SM measurement and **searching of BSM**

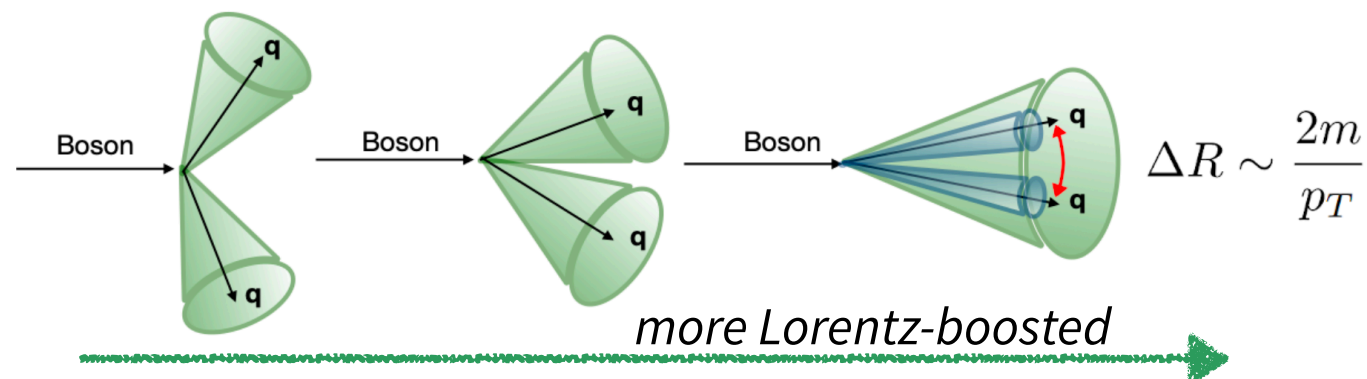


- ❁ *Deep learning (AI) is bringing a technological leap in analyzing LHC data*
- ❁ *It is intriguing to think where the future possibility lies*



Boosted topology – a booster to sensitivity for LHC physics

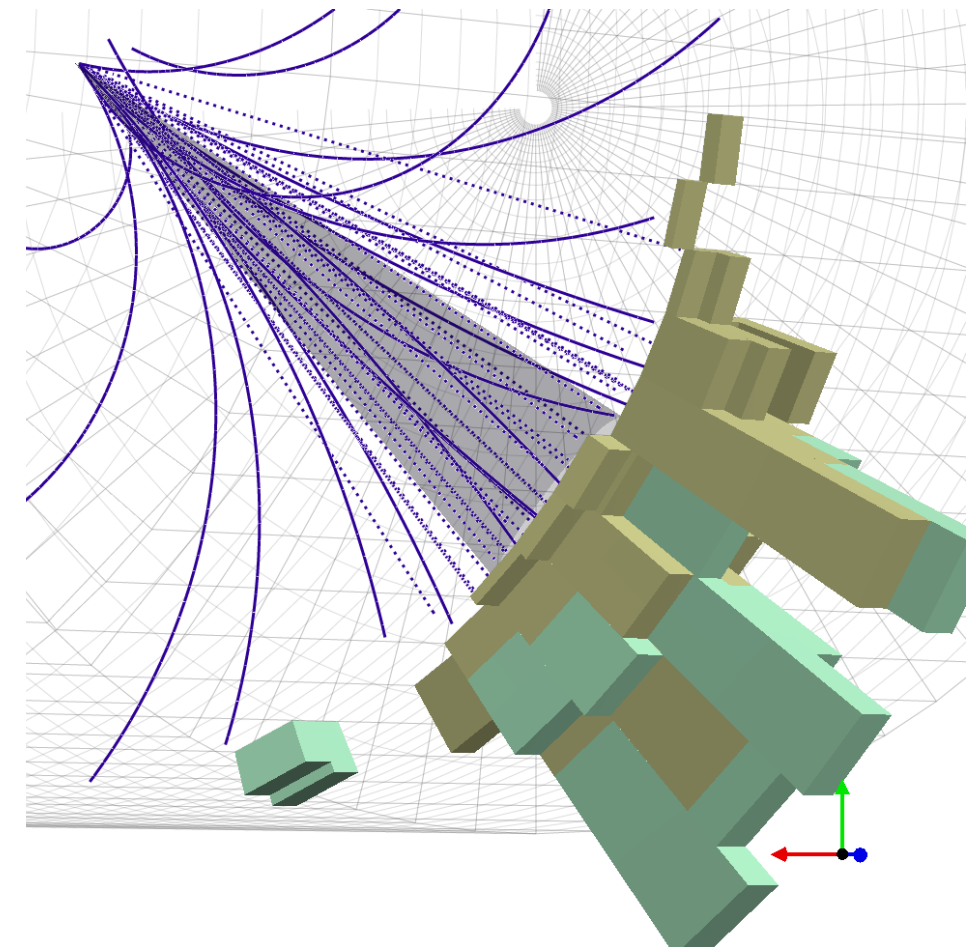
→ Large- R jets: an important handle to analyze boosted topologies at the LHC



❖ Applications to Higgs/di-Higgs/BSM searches in boosted $H(X) \rightarrow b\bar{b}/c\bar{c}$ final states have been a success

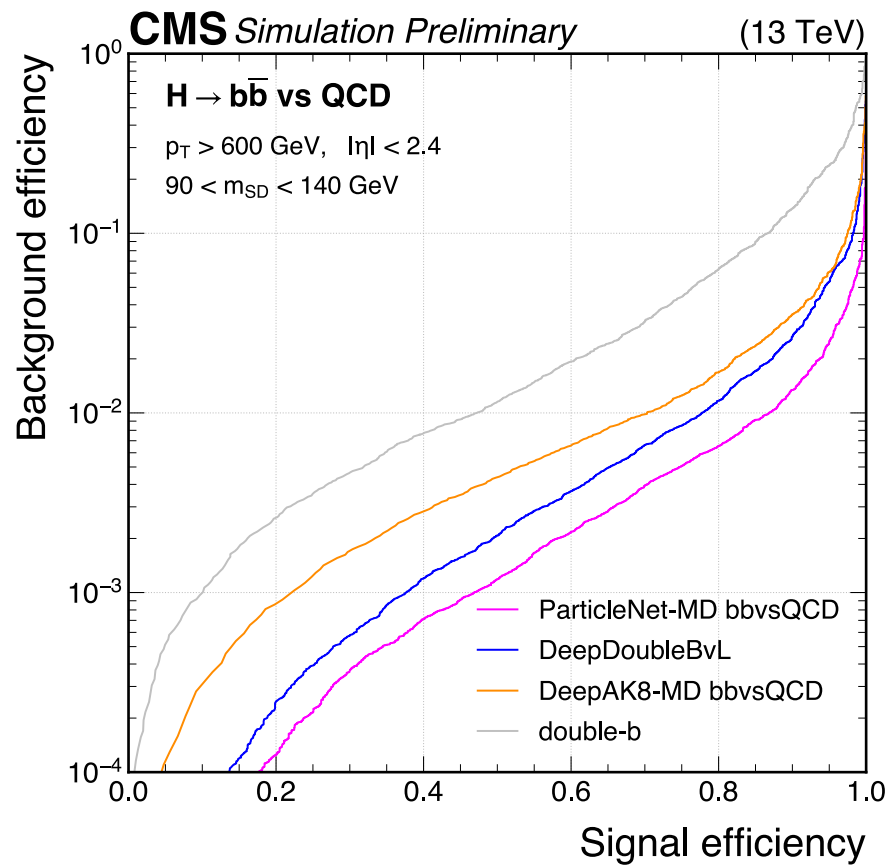
→ *Suitable for deploying cutting-edge deep learning techniques*

- ❖ most complex object to handle at the LHC (up to ~100 constituent particles)
- ❖ advanced DNNs greatly boost analysis sensitivity

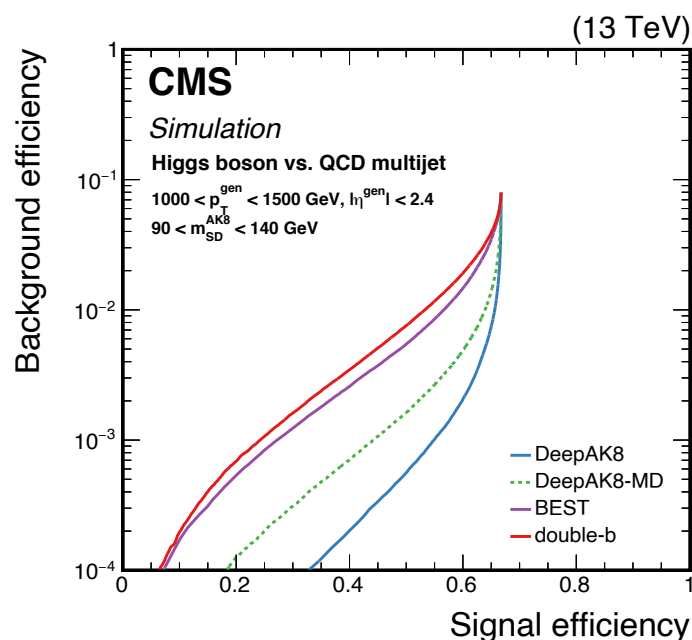


Inspiring progress on $H \rightarrow b\bar{b}/c\bar{c}$ tagging

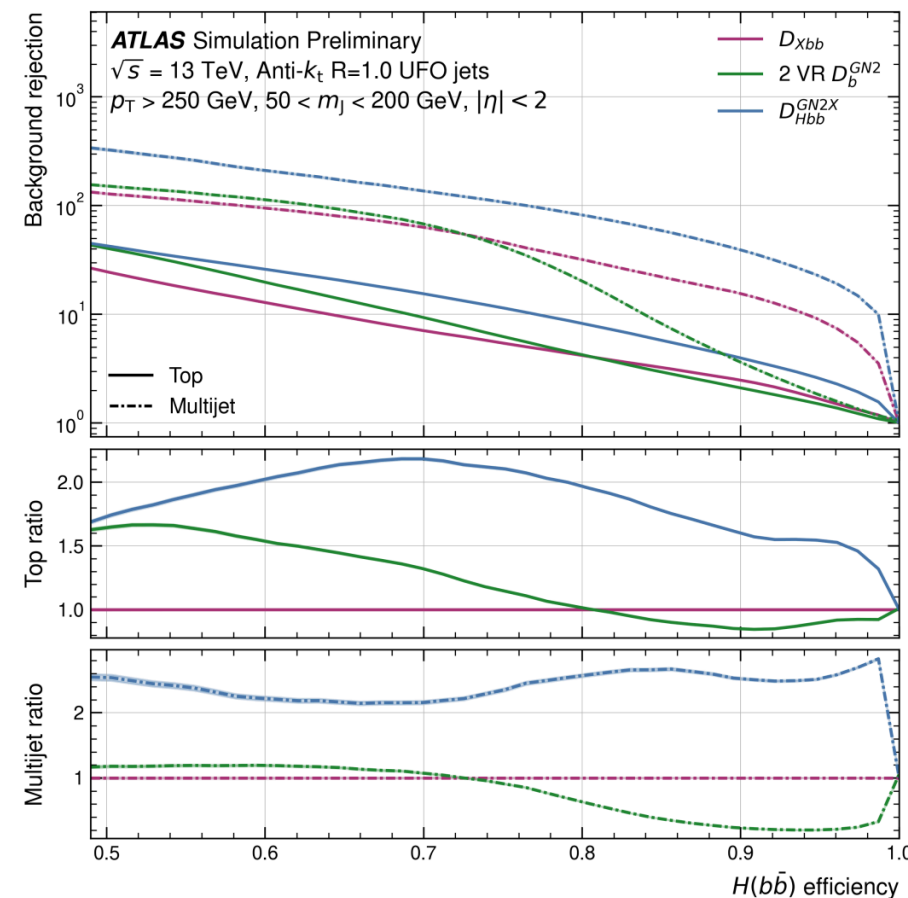
CMS-PAS-BTV-22-001



JINST 15 (2020) P06005



ATL-PHYS-PUB-2023-021



An upgrade of network

DeepAK8 \rightarrow ParticleNet:

x5 QCD background rejection

Note: back to the results 5 years ago

DeepAK8 tagger already has **~x5** improved background rejection than early methods

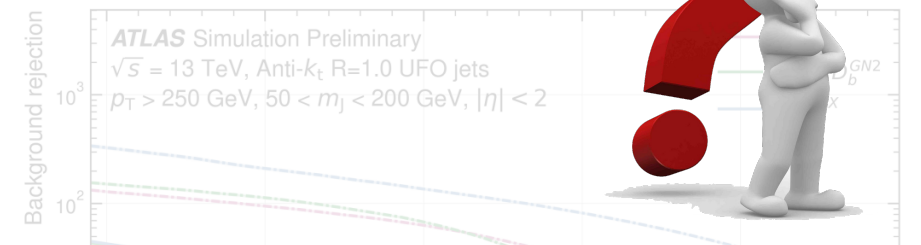
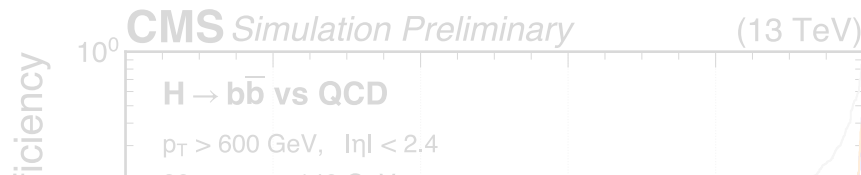
Recent GN2X tagger:

~x3 QCD and x2 top background rejection

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



Implications



Advancements in NN design are the true driving force behind the gains in sensitivity!

- **However, this tool is available only in limited phase space**
(e.g. $H \rightarrow b\bar{b}$, $W \rightarrow qq$, $t \rightarrow bqq$...)
- **Can we extend its usage to all possible boosted phase spaces?**
- **If we can boost the sensitivity by $\times 10$ in many phase spaces, will it accelerate the "next potential discovery" of a new resonance?**

DeepAK8 \rightarrow ParticleNet:

$\times 5$ QCD background rejection

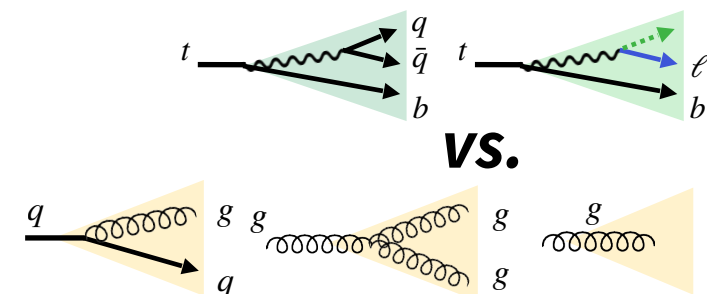
DeepAK8 tagger already has $\sim \times 5$
improved background rejection
than early methods

top background
rejection

Propose “Large model for large-scale classification”

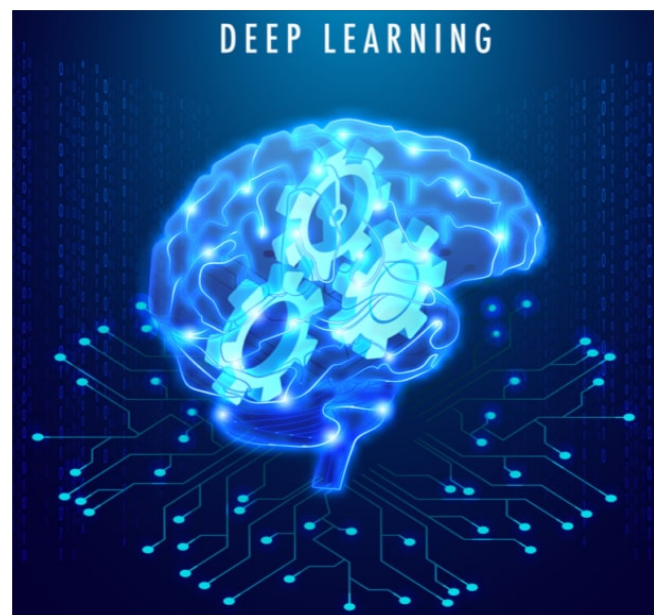
View from jet tagging

- 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 (next slide)
- 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

- We own a comprehensive jet dataset, and we hope to pre-train **a foundational 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 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

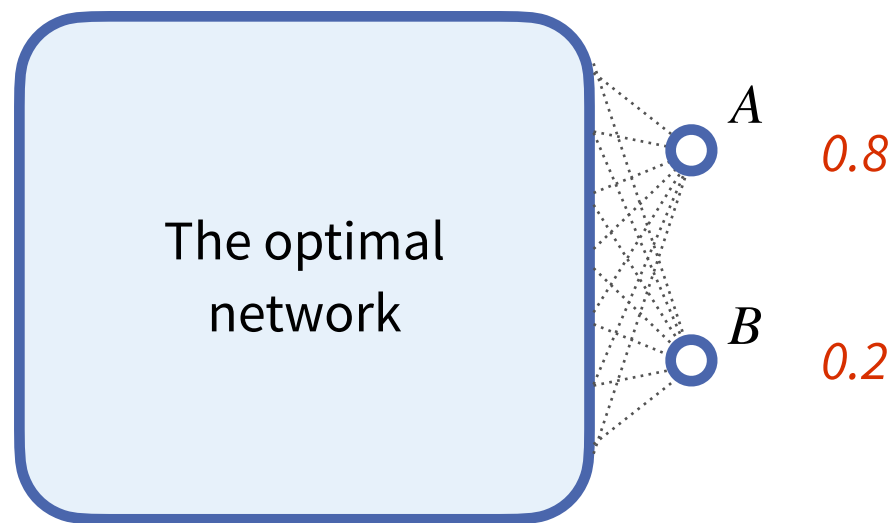
→ Statistical theory shows that:

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

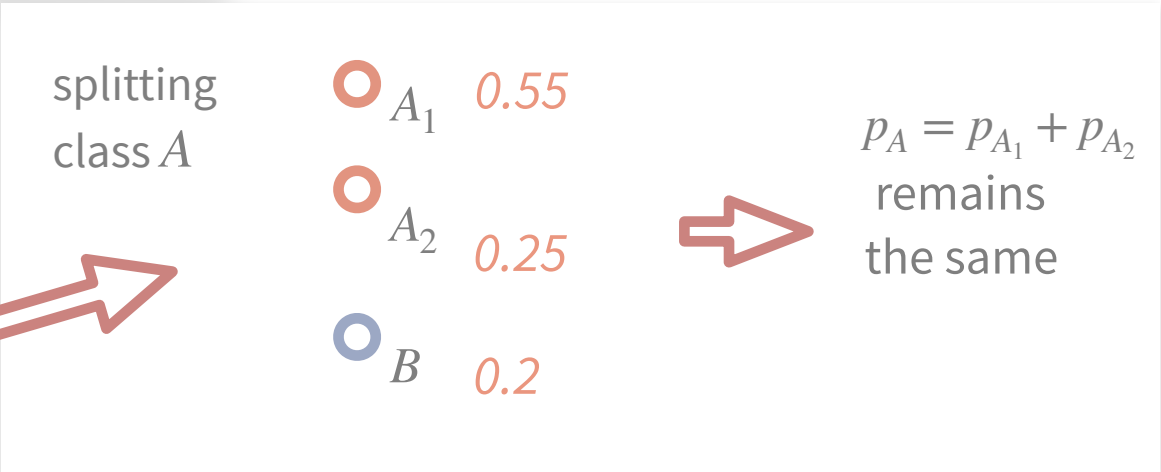
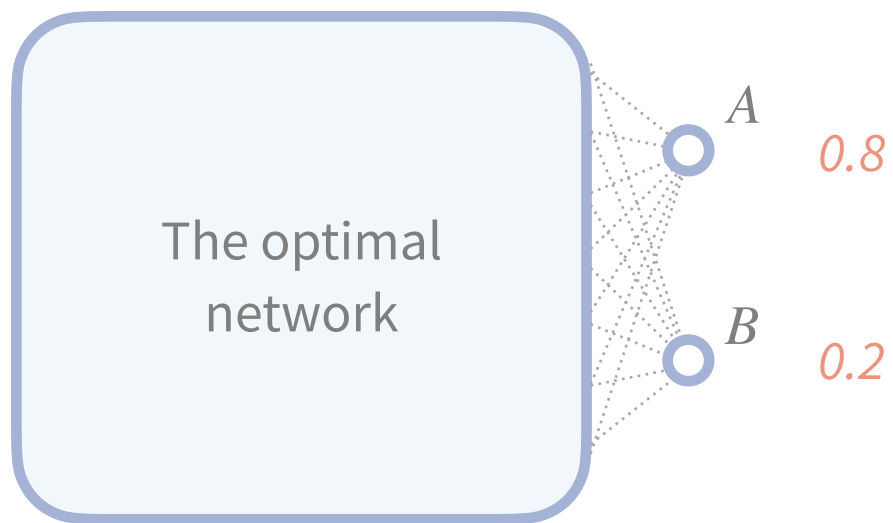
→ Statistical theory shows that:
 A **multi-class** classifier with minimum **cross-entropy loss** estimates the probability ratios on the input classes:

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

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:



Introducing Sophon

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

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

→ We explore this possibility in the CMS experiment first, and also in a recent pheno work:

❖ **Signature-Oriented Pre-training for Heavy-resonant Observation**

[H.Qu, CL, S.Qian. ICML 2022]

❖ the model is based on [Particle Transformer \(ParT\)](#) architecture



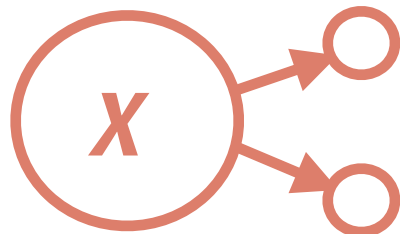
The current “state-of-the-art”
on large-scale dataset
($\mathcal{O}(100M)$), close to experimental setup)

❖ a pre-trained model on a comprehensive dataset: **JetClass-II**

▸ **finely categorized labels:**

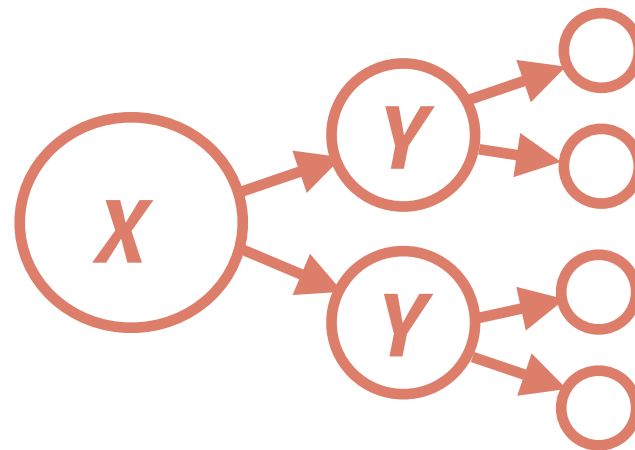
Resonant jet:

$X \rightarrow 2$ prong



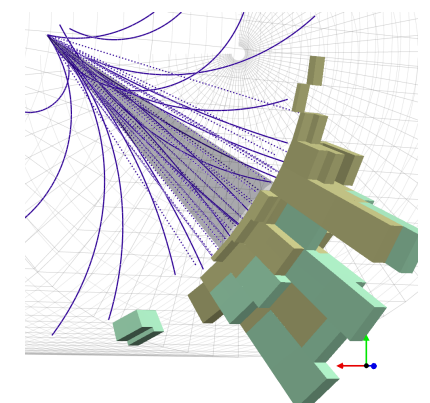
Resonant jet:

$X \rightarrow 3/4$ 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

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

→ We explore this possibility in the CMS experiment first, and also in a recent pheno work:

❖ **Signature-Oriented Pre-training for Heavy-resonant Observation**

[H.Qu, CL, S.Qian. ICML 2022]

❖ the model is based on [Particle Transformer \(ParT\)](#) architecture



The current “state-of-the-art”
on large-scale dataset
(o(100M), close to experimental setup)

❖ a pre-trained model on a comprehensive dataset: **JetClass-II**

▸ **finely categorized labels:**

Major types	Index range	Label names
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	$bbbb, bbcc, bbss, bbqq, bbgg, bbee, bb\mu\mu, bb\tau_h\tau_e, bb\tau_h\tau_\mu, bb\tau_h\tau_h, bbb, bbc, bbs, bbq, bbg, bbe, bb\mu, cccc, ccss, ccqq, ccgg, ccee, cc\mu\mu, cct_h\tau_e, cct_h\tau_\mu, cct_h\tau_h, ccb, ccc, ccs, ccq, ccg, cce, cc\mu, ssss, ssqq, ssgg, ssee, ss\mu\mu, sst_h\tau_e, sst_h\tau_\mu, sst_h\tau_h, ssb, ssc, sss, ssq, ssg, sse, ss\mu, qqqq, qqqg, qqee, qq\mu\mu, qq\tau_h\tau_e, qq\tau_h\tau_\mu, qq\tau_h\tau_h, qqb, qqc, qqs, qqg, qqe, qq\mu, gggg, ggee, gg\mu\mu, gg\tau_h\tau_e, gg\tau_h\tau_\mu, gg\tau_h\tau_h, ggb, ggc, ggs, ggq, ggg, gge, gg\mu, bee, cee, see, qee, gee, b\mu\mu, c\mu\mu, s\mu\mu, q\mu\mu, g\mu\mu, b\tau_h\tau_e, c\tau_h\tau_e, s\tau_h\tau_e, q\tau_h\tau_e, g\tau_h\tau_e, b\tau_h\tau_\mu, c\tau_h\tau_\mu, s\tau_h\tau_\mu, q\tau_h\tau_\mu, g\tau_h\tau_\mu, b\tau_h\tau_h, c\tau_h\tau_h, s\tau_h\tau_h, q\tau_h\tau_h, g\tau_h\tau_h, qqqb, qqqc, qqqs, bbcq, ccbs, ccbq, ccsq, sscq, qqbc, qqbs, qqcs, bcsq, bcs, bcq, bsq, csq, bce\nu, cse\nu, bqev, cqev, sqev, qqev, bc\nu\nu, cs\nu\nu, bq\nu\nu, cq\nu\nu, sq\nu\nu, qq\nu\nu, bct_e\nu, cst_e\nu, bqte\nu, cqte\nu, sqte\nu, qqte\nu, bct_\mu\nu, cst_\mu\nu, bqtm\nu, cqtm\nu, sqtm\nu, qqtm\nu, bct_h\nu, cst_h\nu, bqth\nu, cqth\nu, sqth\nu, qqth\nu$
QCD jets	161–187	$bbccss, bbccs, bbcc, bbcss, bbcs, bbc, bbss, bbs, bb, bccss, bccs, bcc, bcss, bcs, bc, bss, bs, b, ccss, ccs, cc, css, cs, c, ss, s, \text{others}$

resulting to 4-prong or 3-prong

All final states!

Key property: we do not focus on any specific X and Y masses

Their masses are variables: ranges from 20-500 GeV

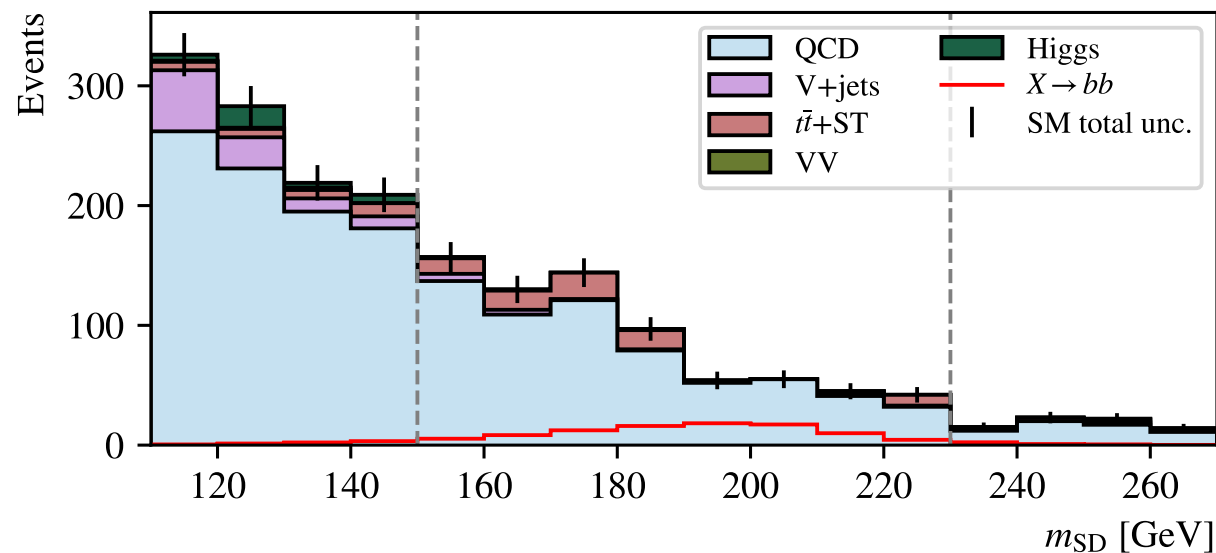
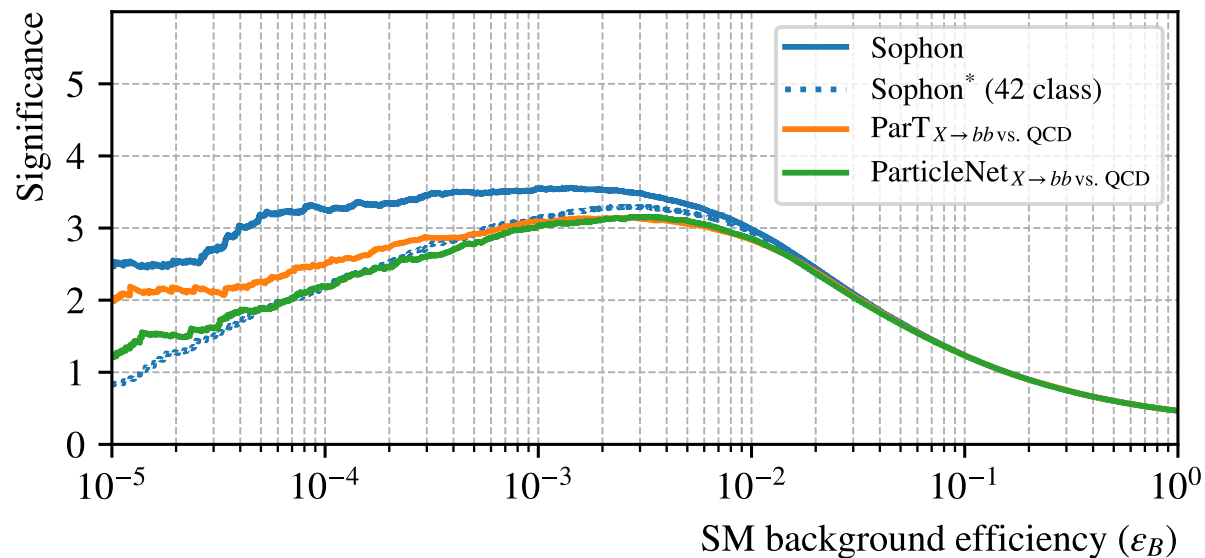
Sophon: performance benchmark

[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)

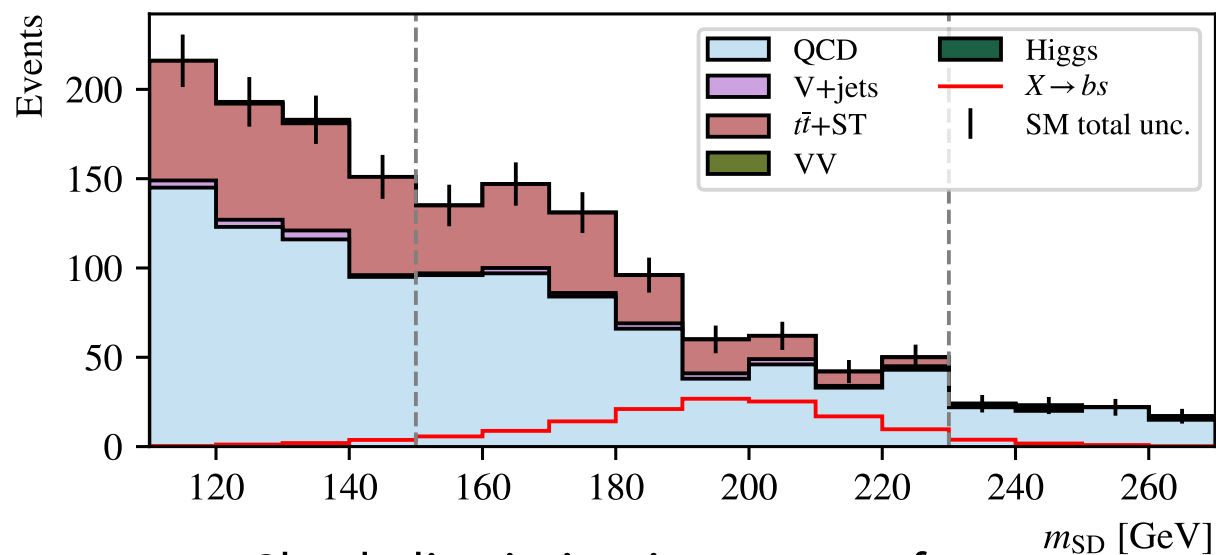
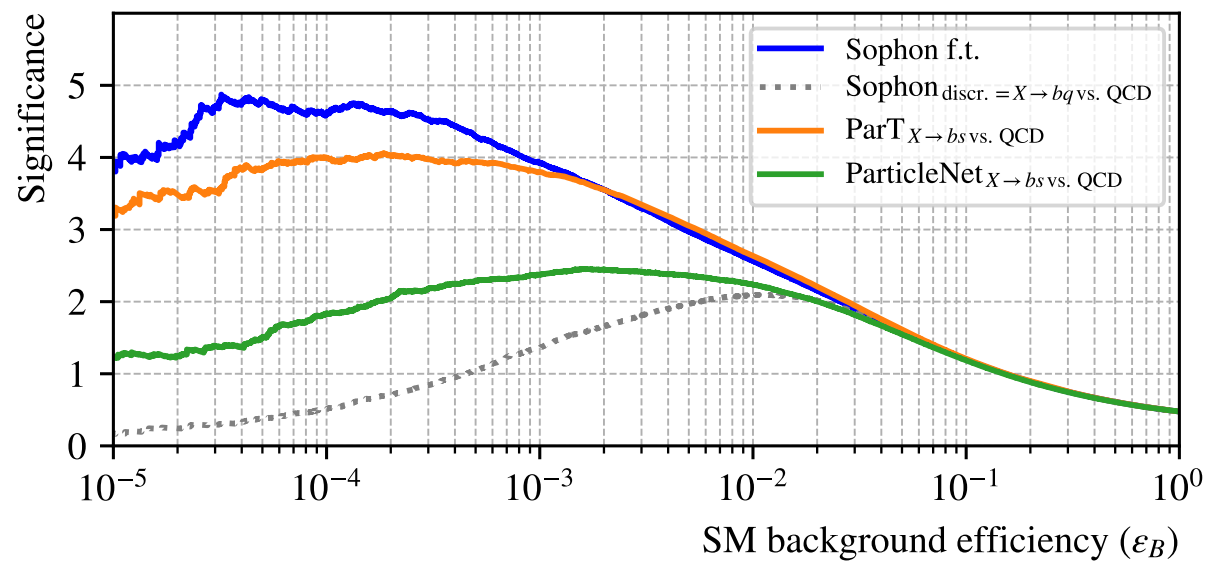
Search significance:

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



Transfer learning ability

(adapt it to a brand new 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?

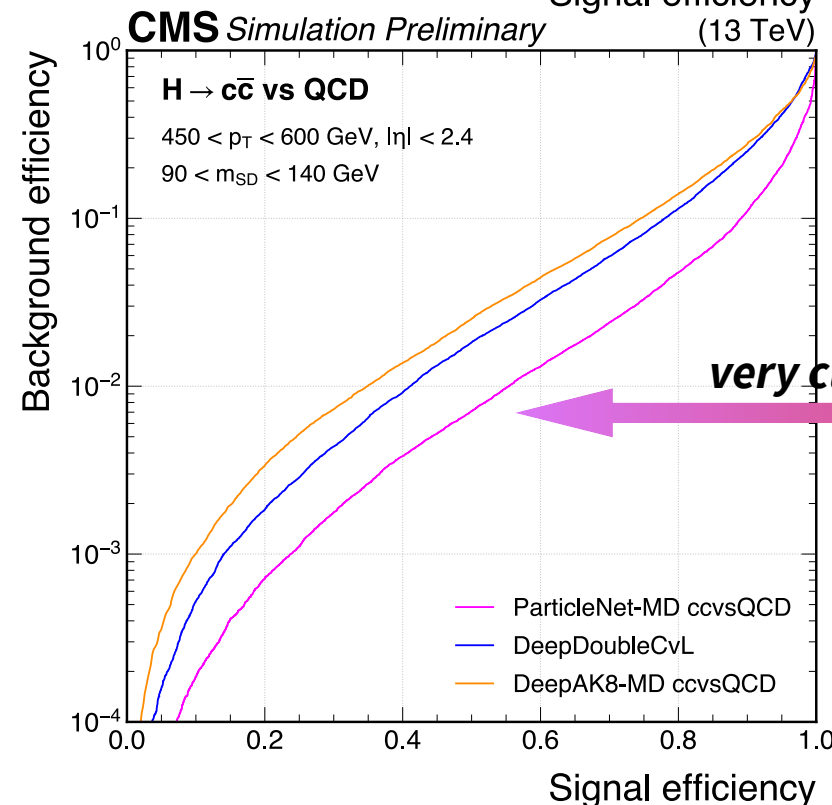
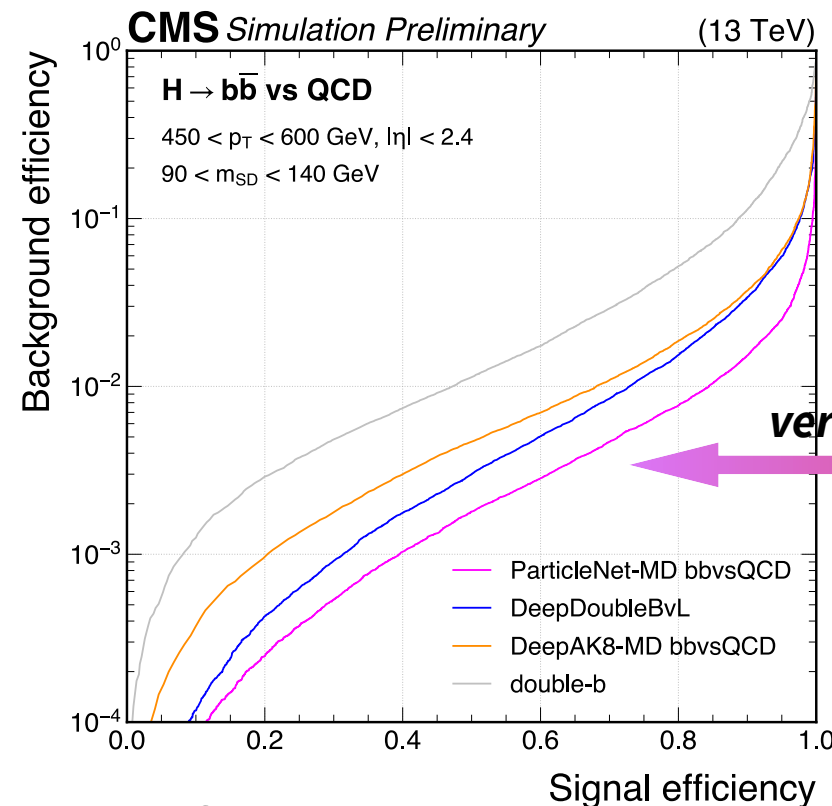
Thanks Yuzhe for his input

Take recent **CMS performance plots** as an example

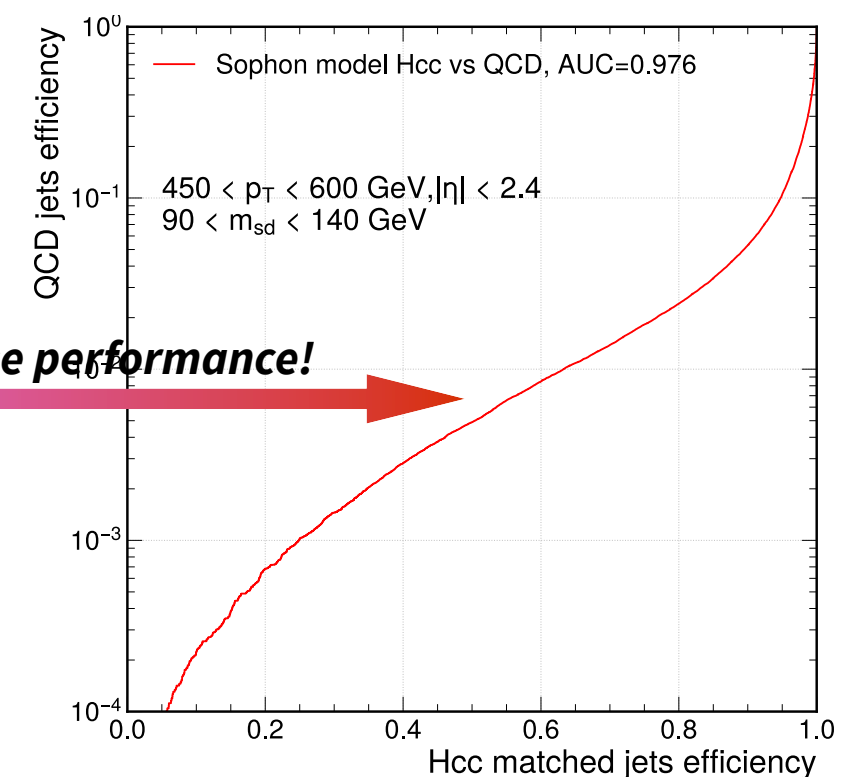
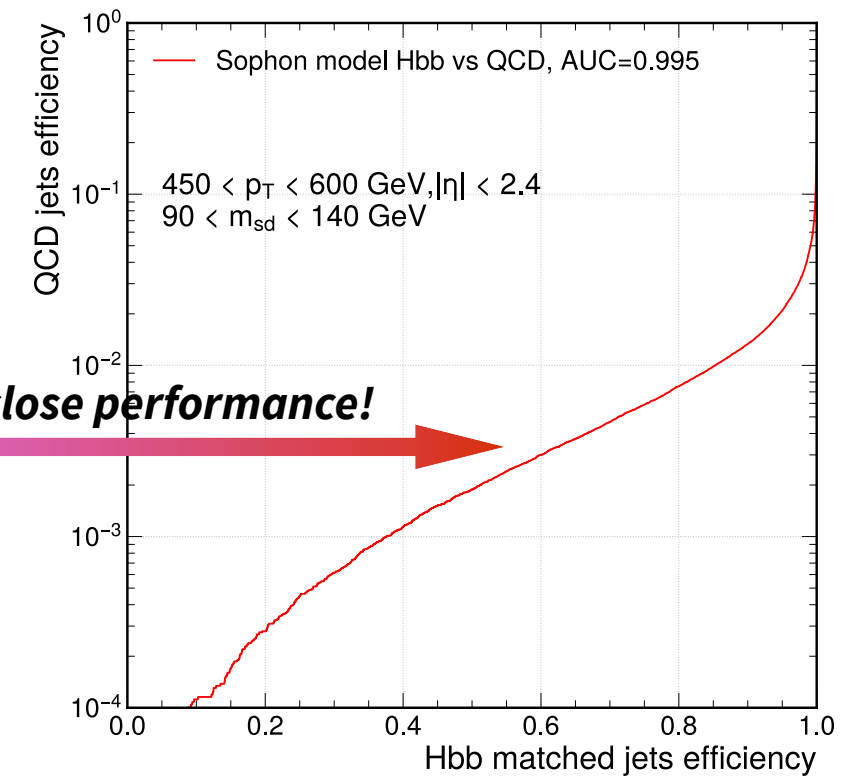
Benchmark performance on SM $H \rightarrow b\bar{b}/c\bar{c}$ jets vs QCD background jets

- Demonstrate that **applying “Sophon” on Delphes dataset** for pheno study is pretty realistic

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



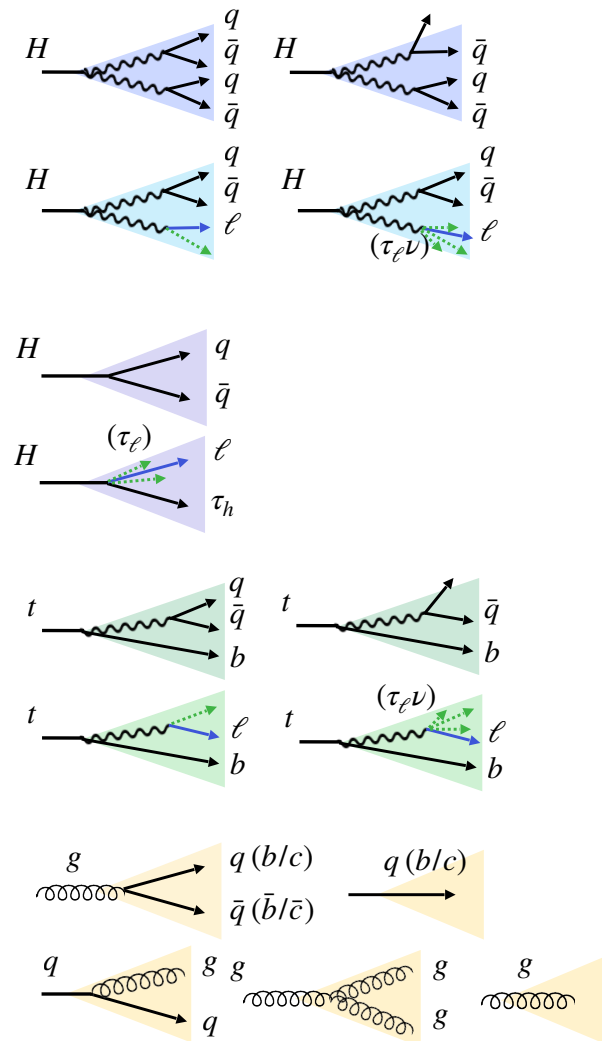
Sophon results (on Delphes dataset)



Highlight: CMS's Global ParT tagger

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 _e νqq		2
	τ _μ νqq		2
	τ _h νqq		2
H→qq		bb	1
		cc	1
		ss	1
		qq (q=u/d)	1
H→ττ	τ _e τ _h		1
	τ _μ τ _h		1
	τ _h τ _h		1
t→bW (hadronic)	bqq	1b + 0c/1c	2
	bq		2
t→bW (leptonic)	bēν	1b	1
	bμν		1
	bτ _e ν		1
	bτ _μ ν		1
	bτ _h ν		1
QCD		b	1
		bb	1
		c	1
		cc	1
		others (light)	1



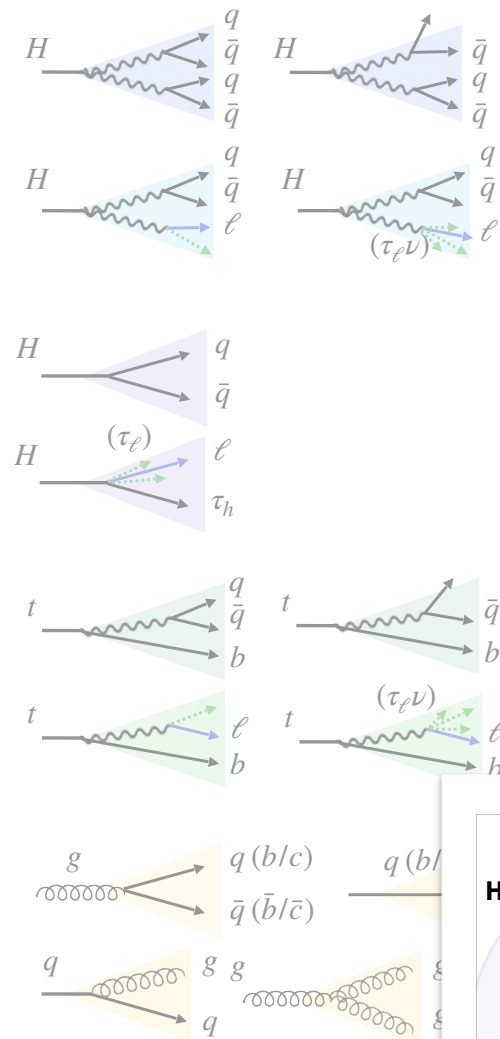
A global large-*R* mass-decorrelated tagger for **37-category classification**

- First time identifying the H→WW→4q signature with a jet tagger
- set a strong limit to κ_{2V} in the search of HH→bbVV signal

Highlight: CMS's Global ParT tagger

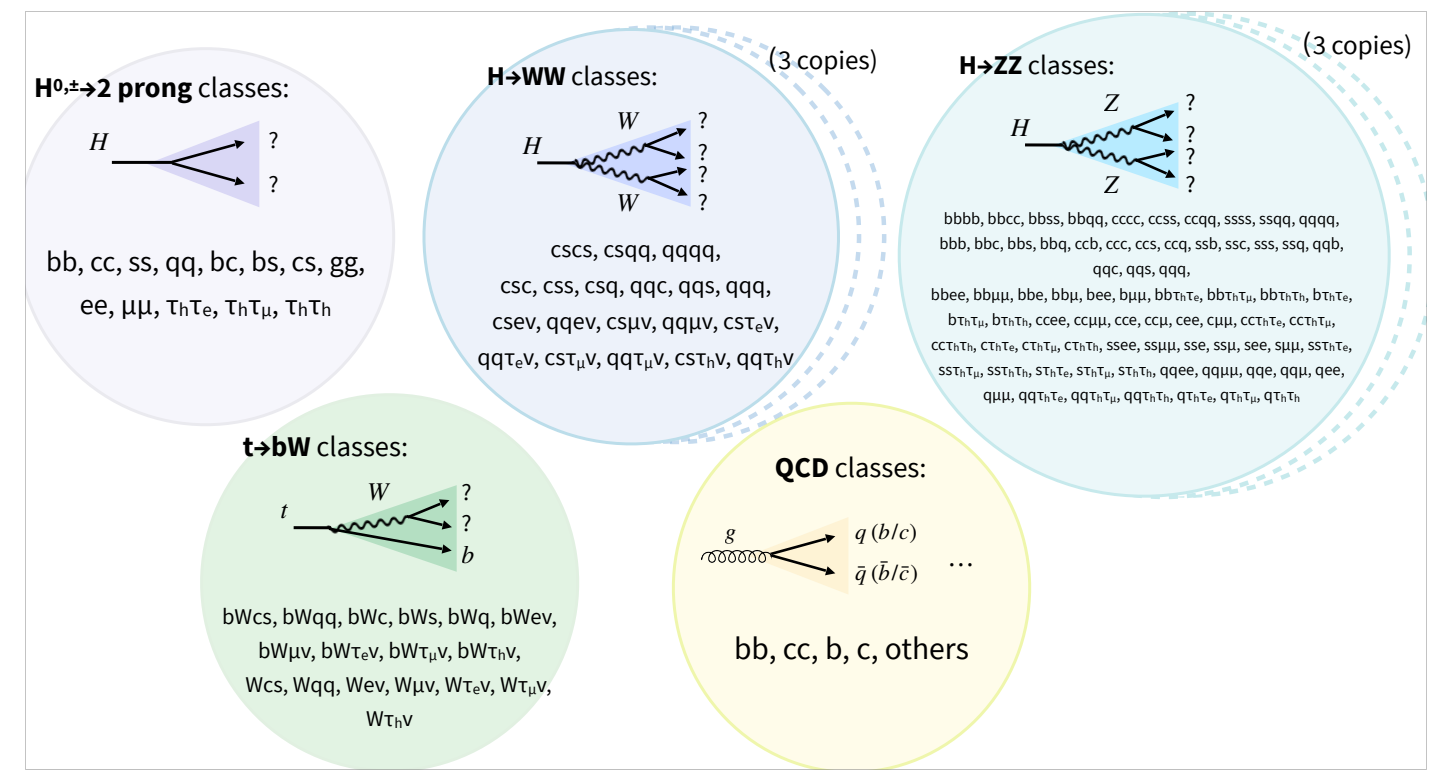
CMS-PAS-HIG-23-012

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 _e νqq		2
	τ _μ νqq		2
	τ _h νqq		2
H→qq		bb	1
		cc	1
		ss	1
		qq (q=u/d)	1
H→ττ	τ _e τ _h		1
	τ _μ τ _h		1
	τ _h τ _h		1
t→bW (hadronic)	bqq	1b + 0c/1c	2
	bq		2
t→bW (leptonic)	b _e ν	1b	1
	b _μ ν		1
	b _{τ_e} ν		1
	b _{τ_μ} ν		1
	b _{τ_h} ν		1
QCD		b	1
		bb	1
		c	1
		cc	1
		others (light)	1



A global large- R mass-decorrelated tagger for **37-category classification**

- First time identifying the $H \rightarrow WW \rightarrow 4q$ signature with a jet tagger
- set a strong limit to κ_{2V} in the search of $HH \rightarrow bbVV$ signal



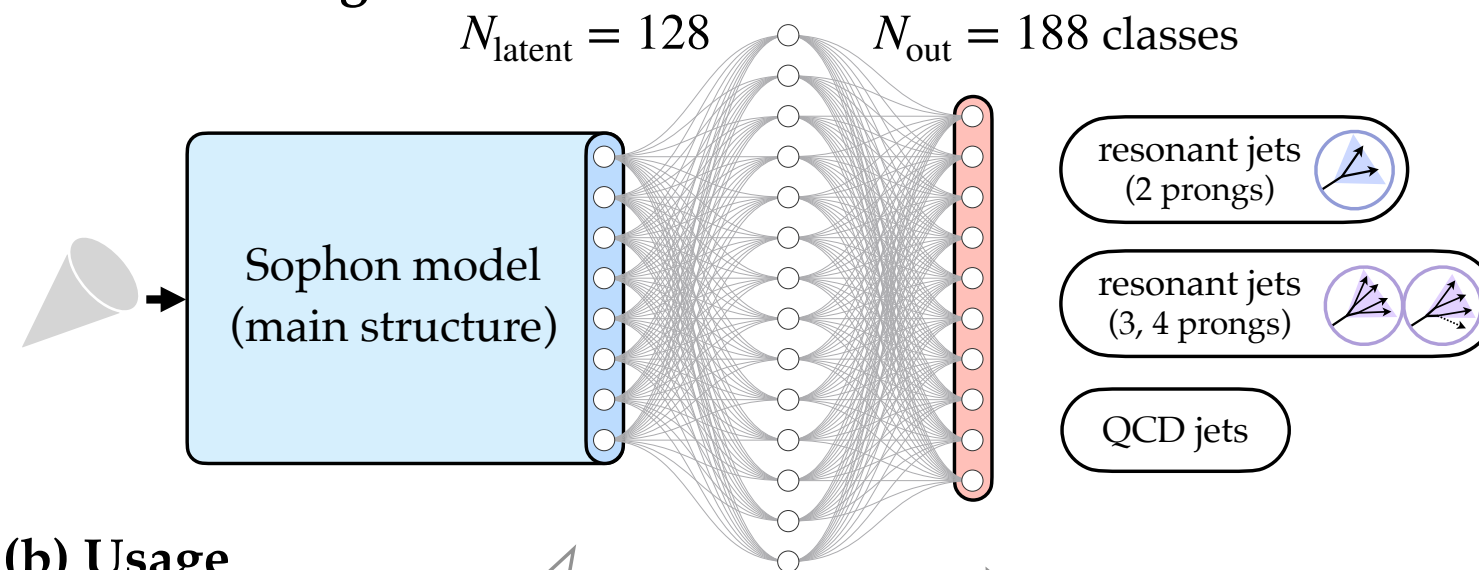
and we plan to update the CMS tagger under Sophon's philosophy

please stay tuned!

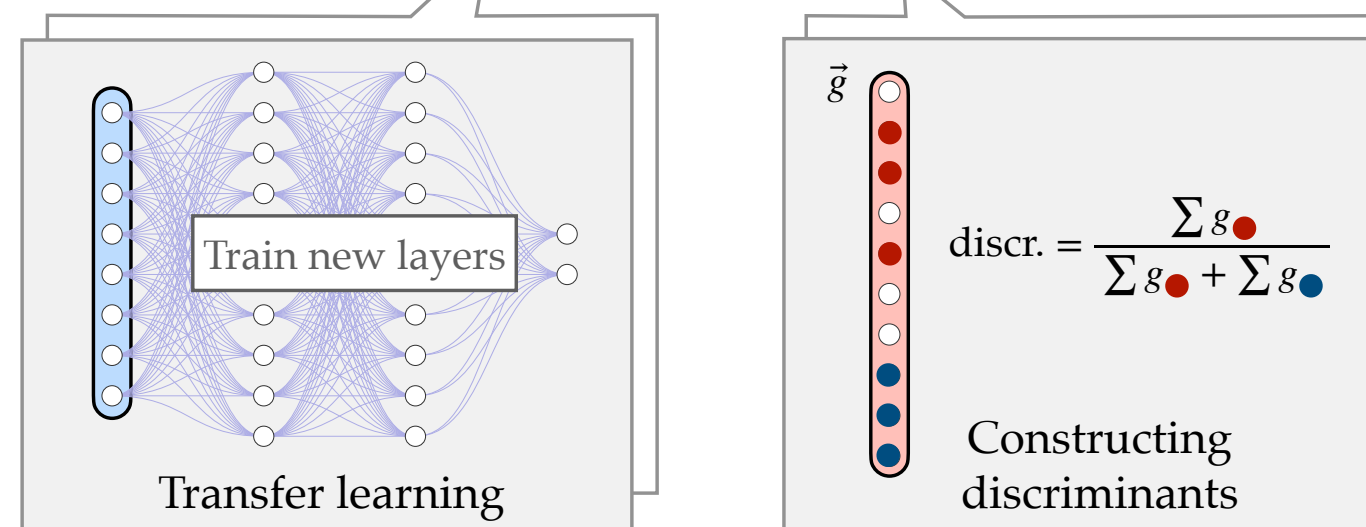
Implications for LHC resonance search

Using Sophon

(a) Pre-training

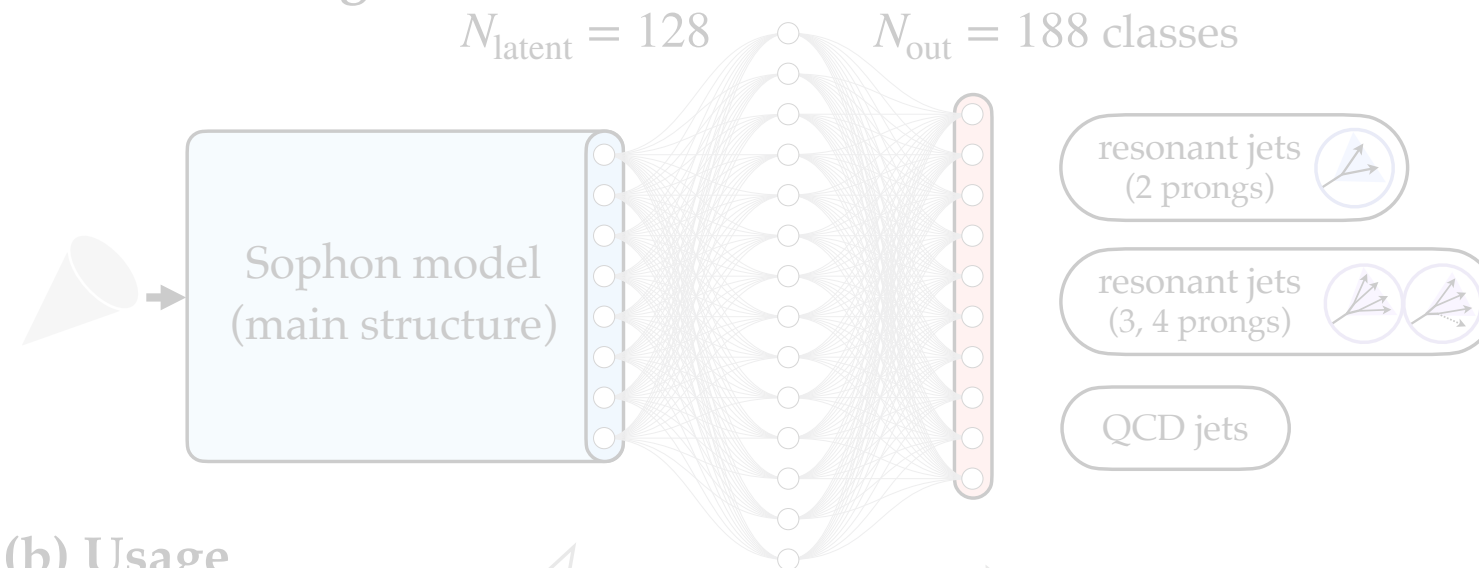


(b) Usage

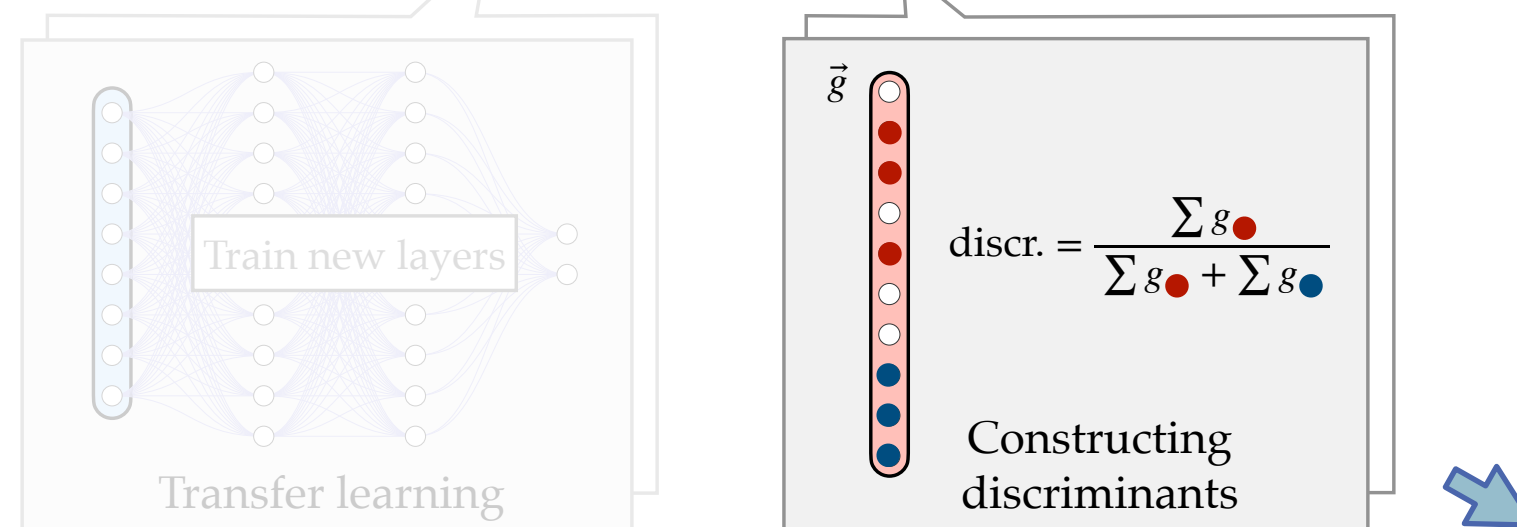


Using Sophon

(a) Pre-training



(b) Usage

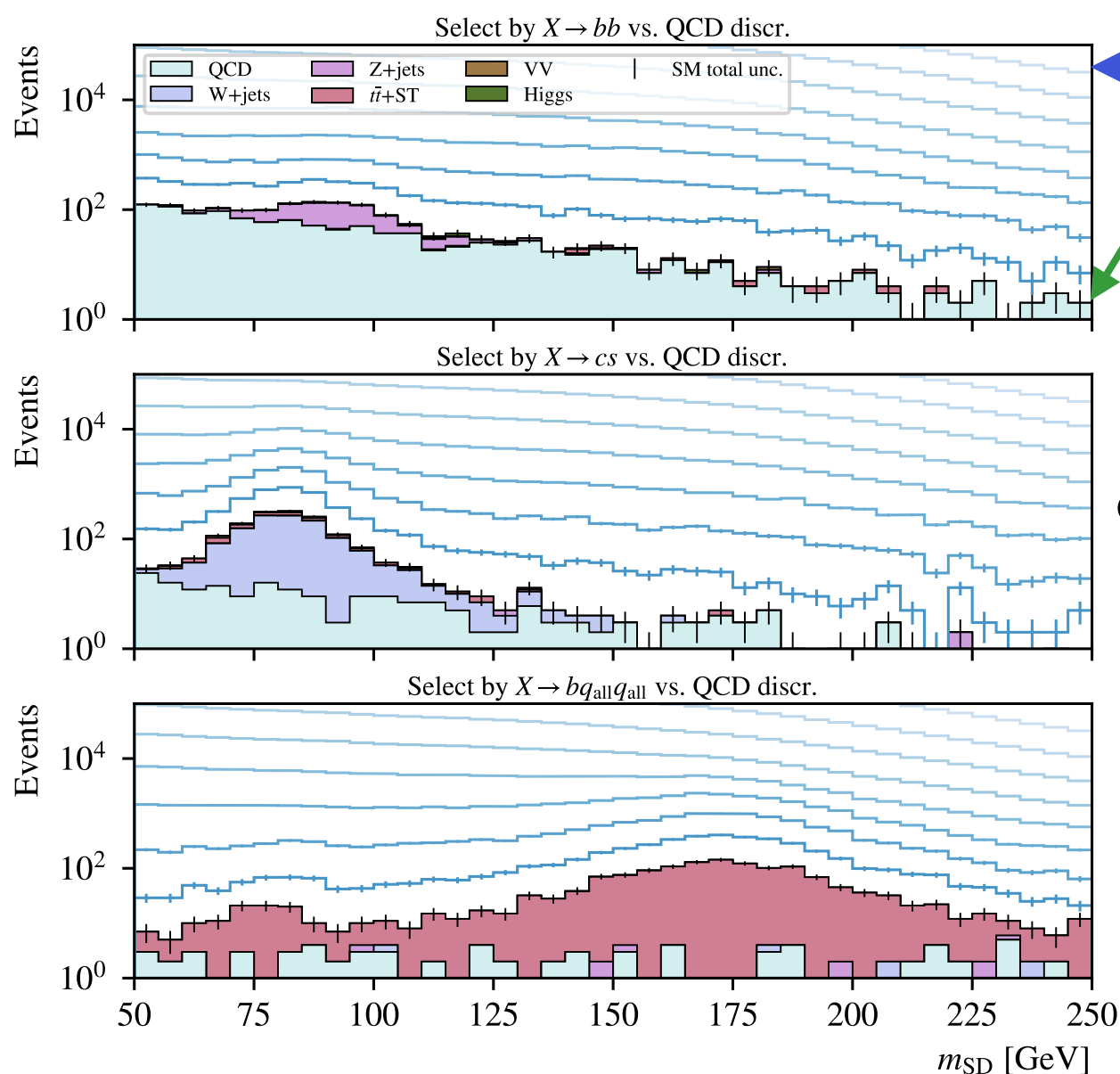


Use it out of the box!

Construct a dedicated discr.
→ perform a bump hunt

Can we rediscover the SM particles?

- Simulate 40fb^{-1} LHC collision events, $\sqrt{s} = 13\text{ TeV}$, $n\text{PU}=50$
- ❖ focus on the large- R jet trigger (triggered with Σp_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

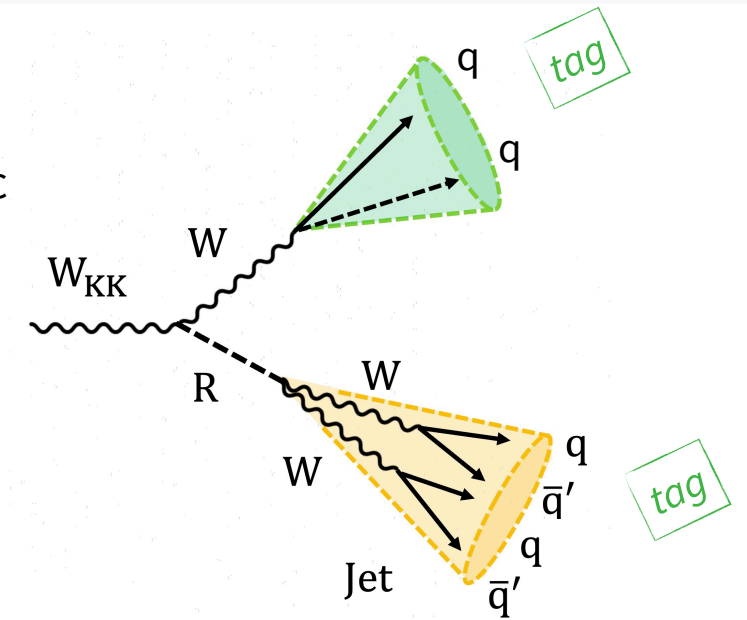
→ Consider triboson signal:

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.

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

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

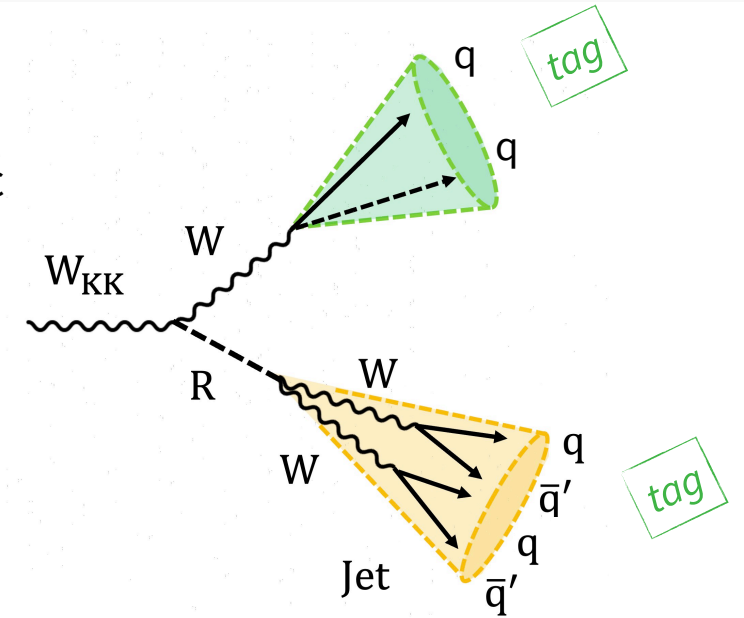


More heavy resonances

→ 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.



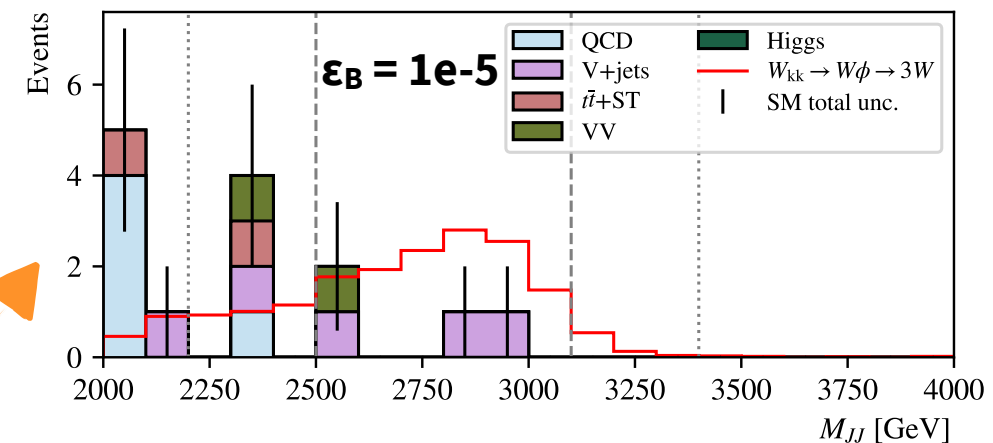
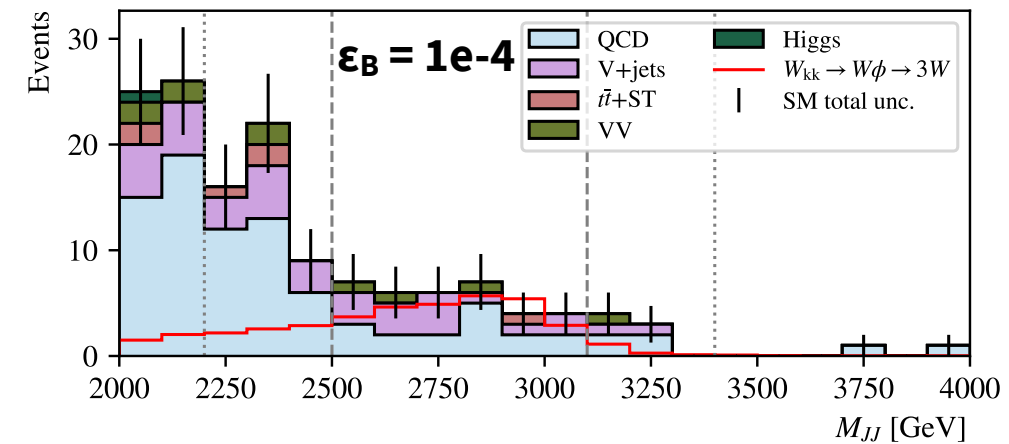
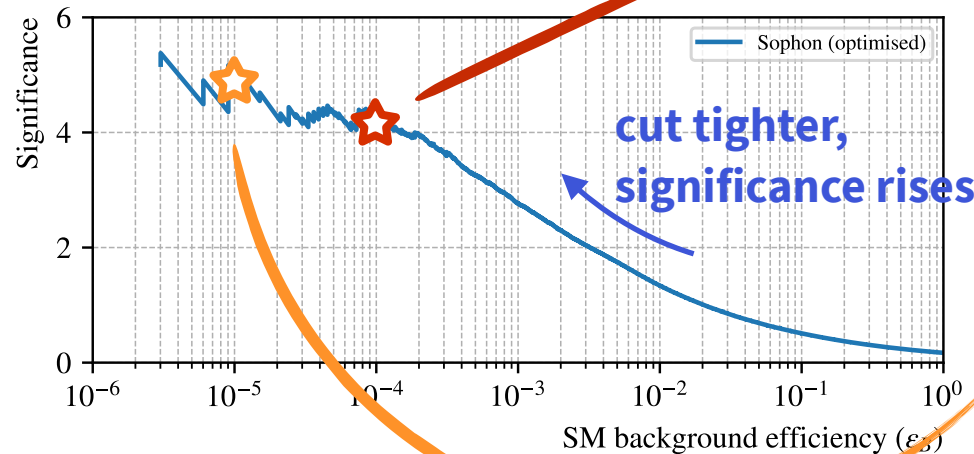
$$\text{discr} = \sum_{\text{jet}=1,2} \frac{g_{A,\text{jet}}}{g_{A,\text{jet}} + \sum_{l=1}^{27} g_{\text{QCD},l,\text{jet}}} \quad (\text{sum for jets 1, 2})$$

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

Search significance

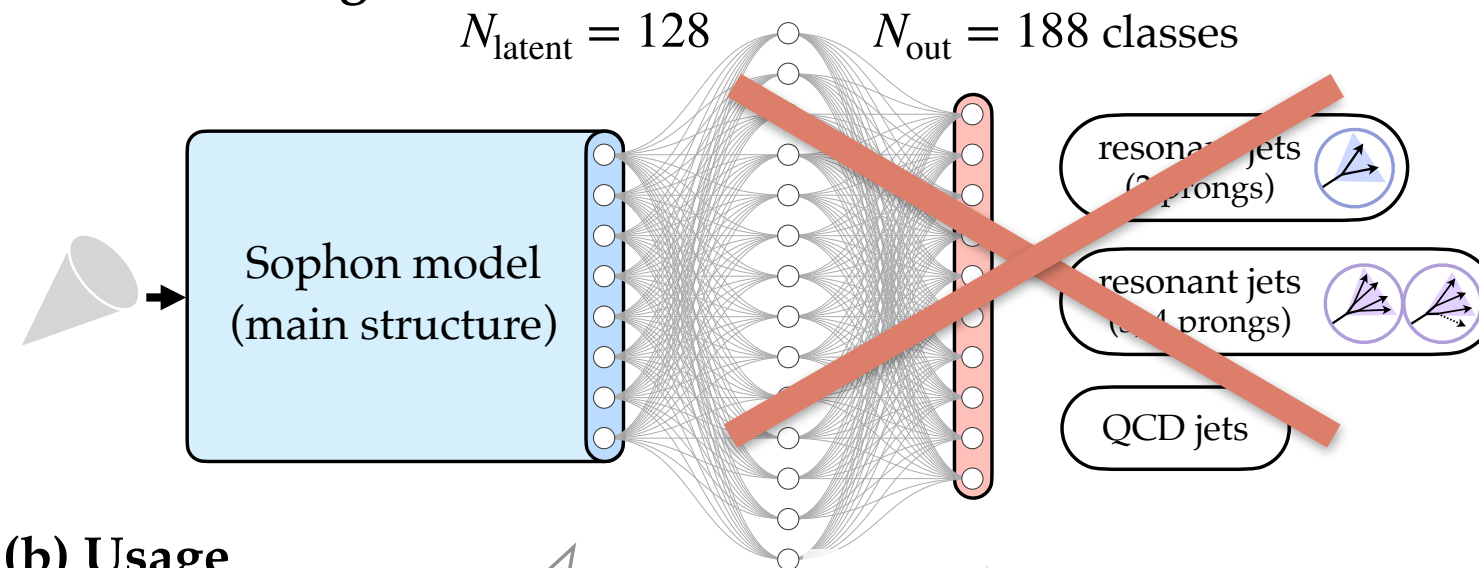
$$Z = \sqrt{2((s+b)\log(1+s/b) - s)}$$

in dijet inv. mass window
2500–3100 GeV

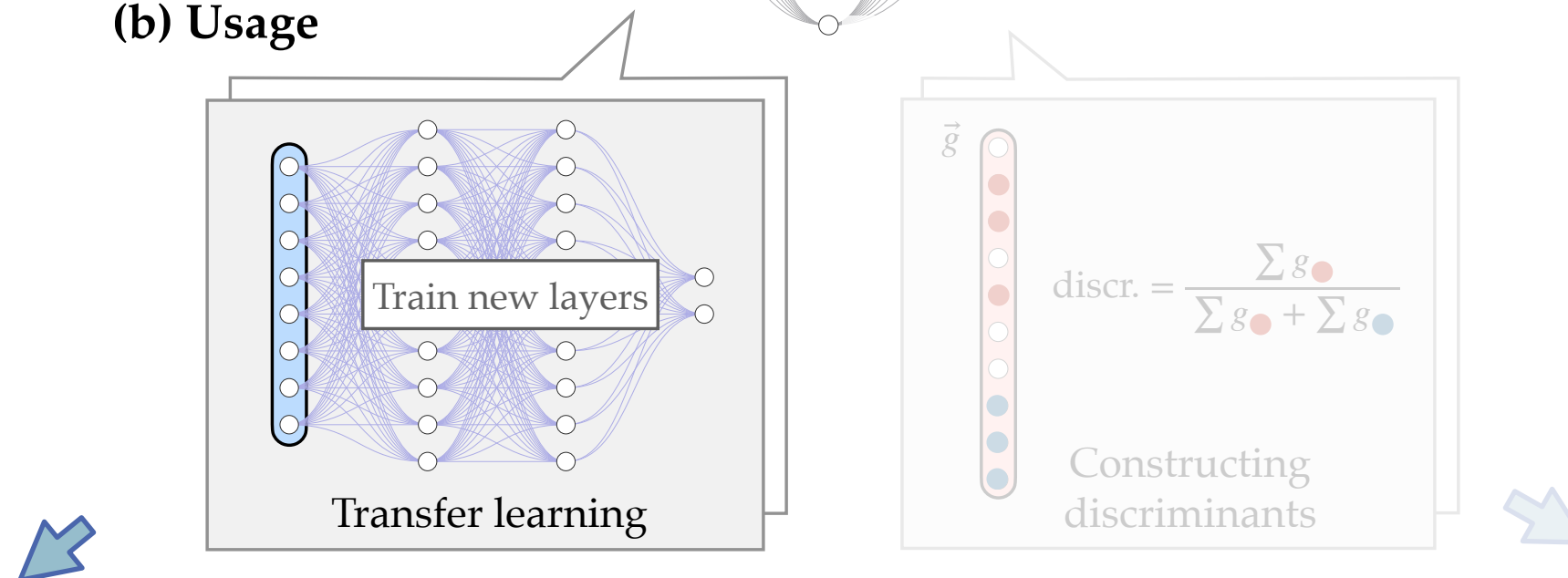


Sophon's transfer learning

(a) Pre-training



(b) Usage



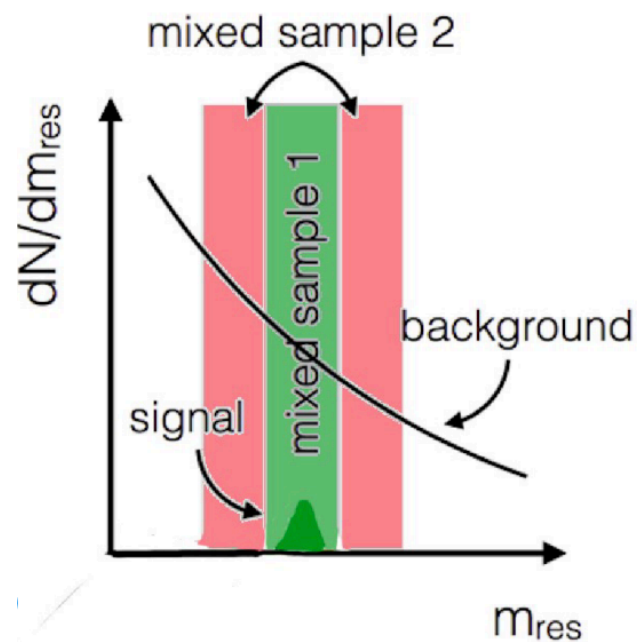
Use it out of the box!

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

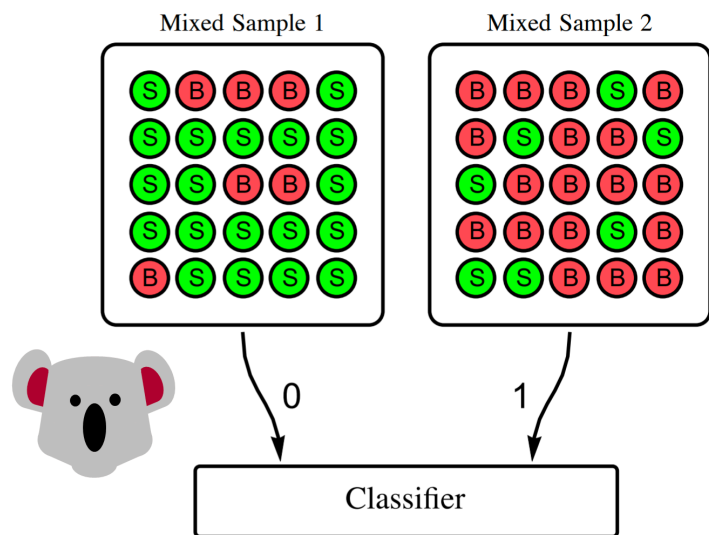
Construct a dedicated discr.
→ perform a bump hunt

Background: anomaly detection in weakly-supervised approach

JHEP 10 (2017) 174



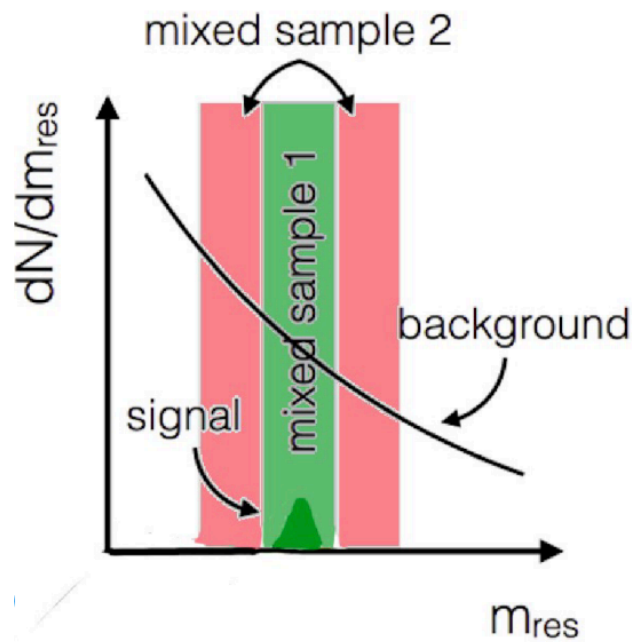
- Recall the early work: CWoLa (classification without labels) Hunting
- ❖ allow to detect anomalies purely from data
 - ❖ train a classifier for mass window vs mass sideband (mixed sample 1 vs 2)
 - ❖ many improved approaches in recent years → very active field



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

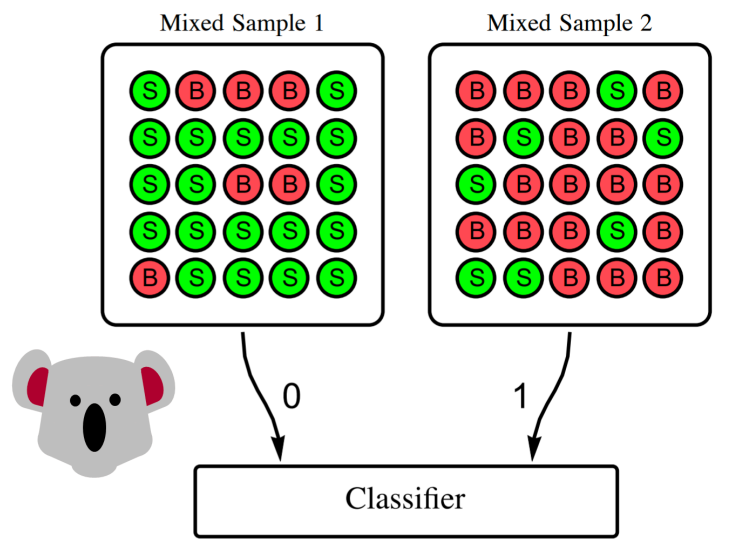
Background: anomaly detection in weakly-supervised approach

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

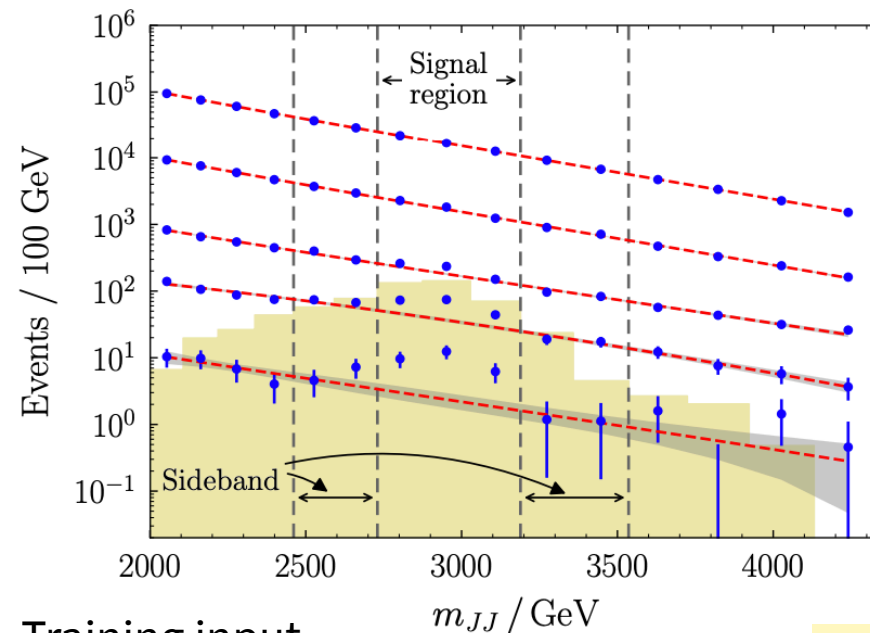


→ Recall the early work: CWoLa (classification without labels) Hunting

- ❖ 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



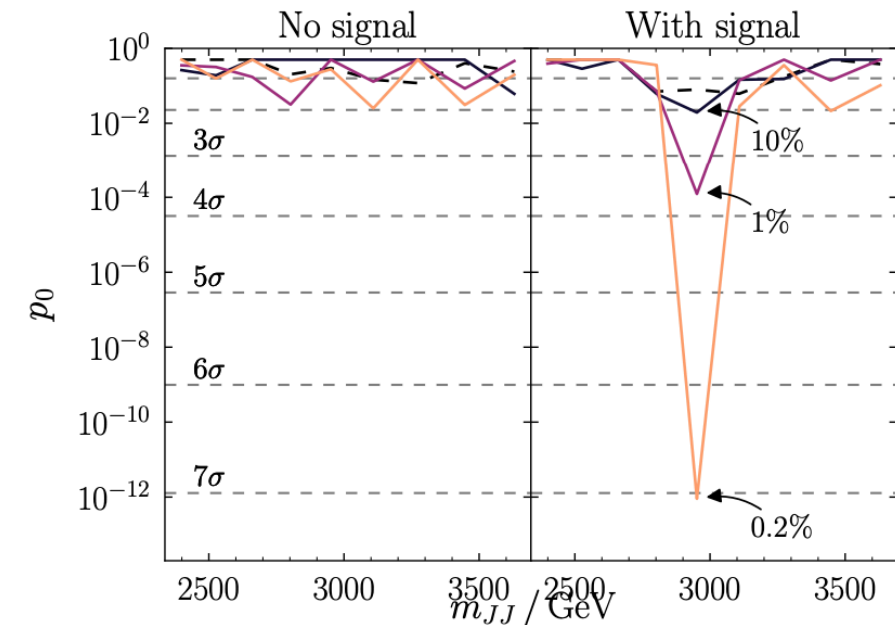
Equivalent effect for training **S** vs **B**



Training input

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

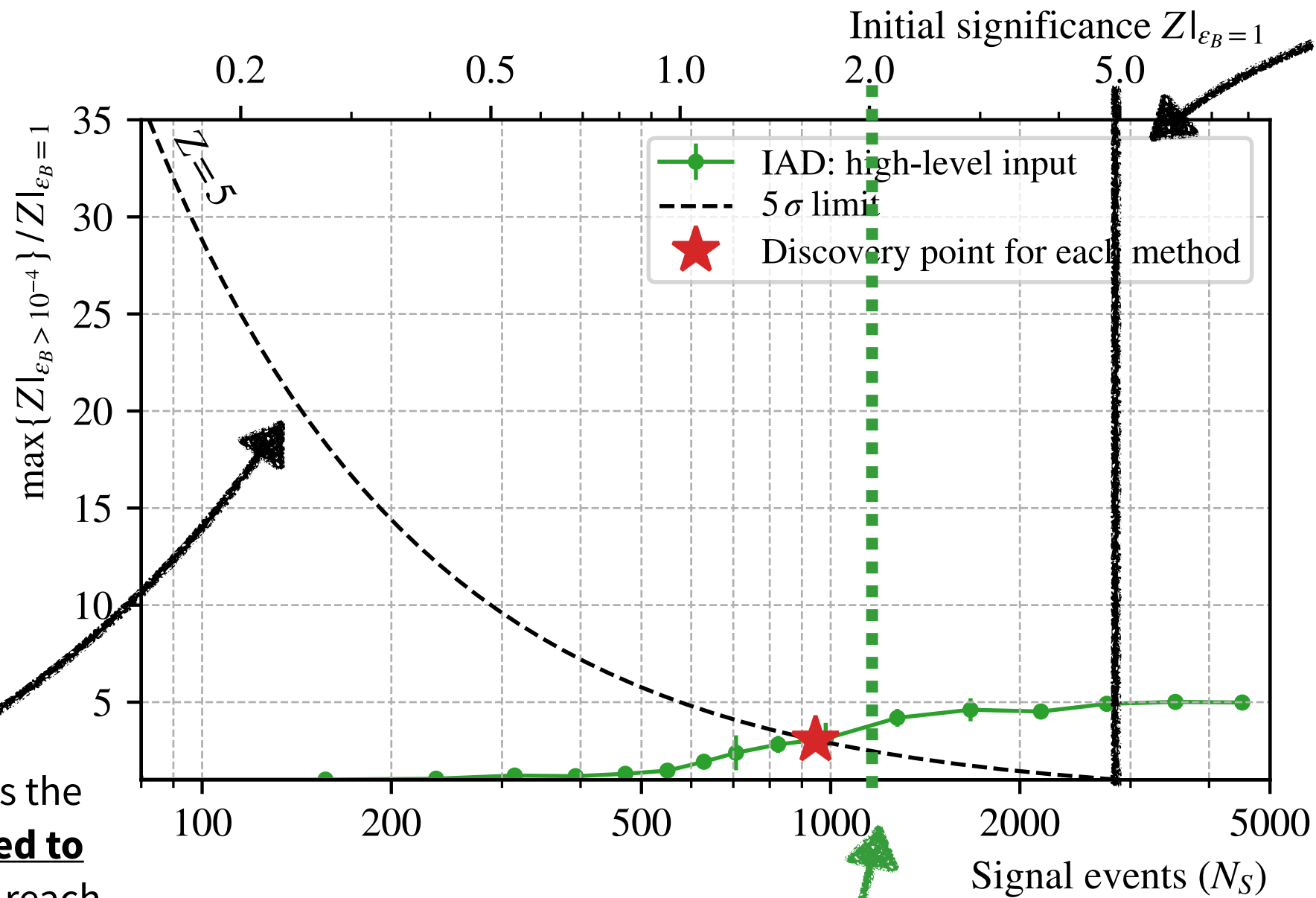
can discover $W' \rightarrow W\phi \rightarrow WWW$ signals
see $2\sigma \rightarrow 7\sigma$ improvement



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

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

Dijet search capabilities



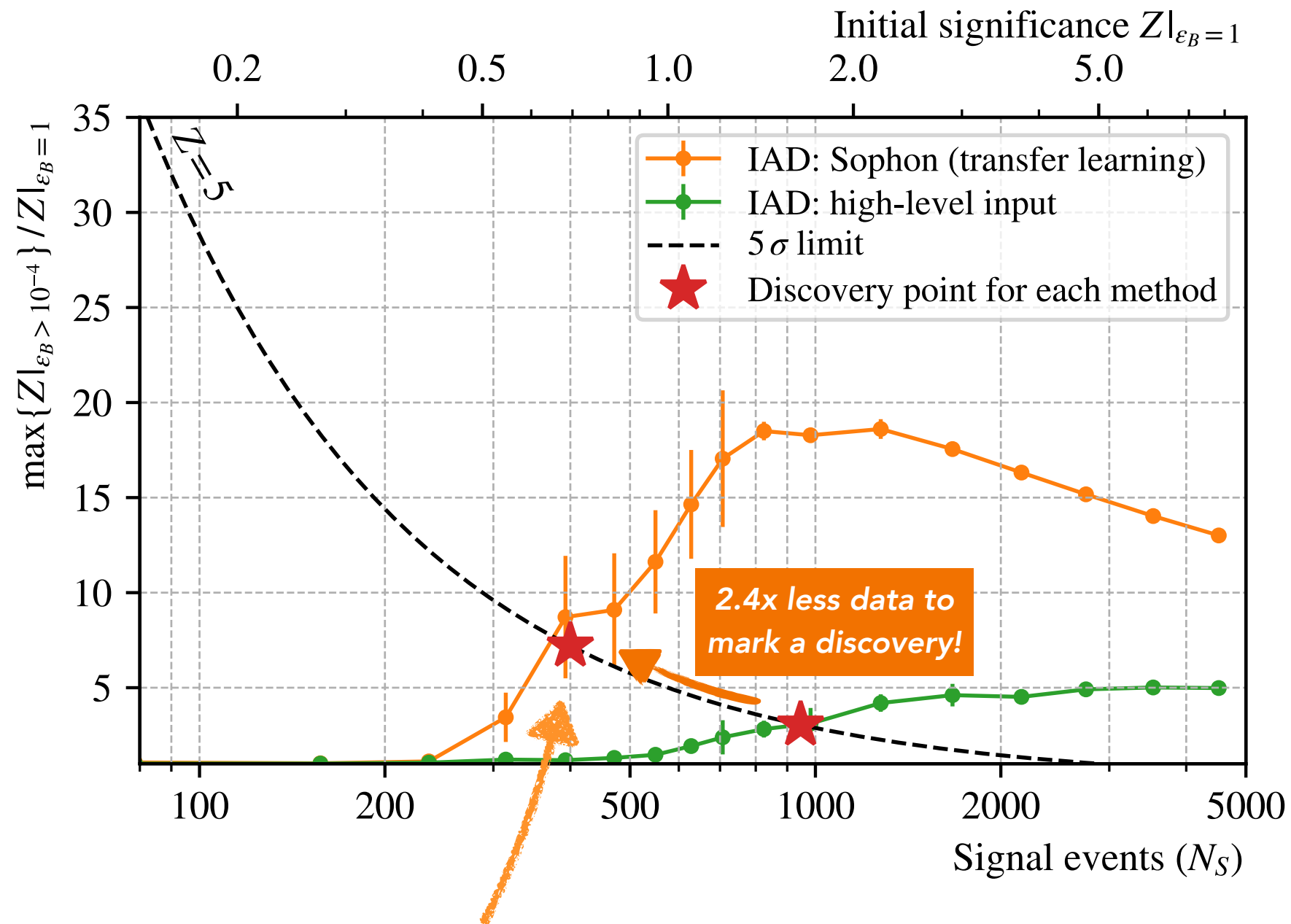
“If signal events reach this point, **with initial $Z=5$** , then we have already discovered the signal without needing to make a cut”

“How much does the **significance need to be increased** to reach the 5σ discovery”

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

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

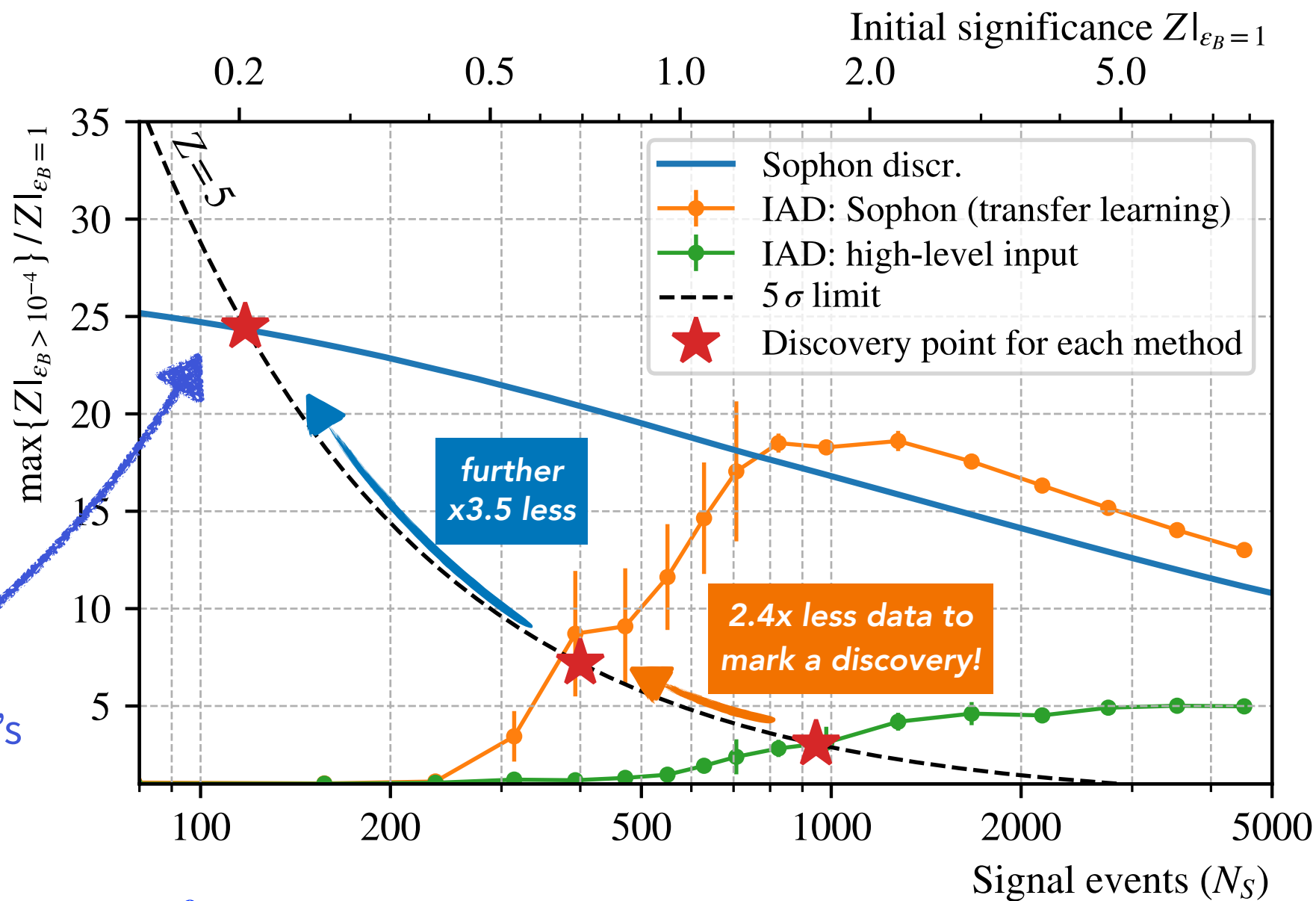
Dijet search capabilities



Combining Sophon's transfer learning (using Sophon's "knowledge") with AD marks a success

- More sensitive to low signal (**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

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

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

Summary and outlook

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

- Sophon releases a lot of new opportunities for future LHC experiments
 - ❖ simply viewed as a “**global large- R jet tagger**” → should bring benefits of the advanced NN to ~all hadronic final-state searches
 - ❖ also viewed as a pre-trained jet model: a foundation model tailored for LHC analyses
- Proposed the **JetClass-II** dataset and the **Sophon** model
 - ❖ **JetClass-II** covers more comprehensive phase spaces and can be a good playground to develop future foundation models
 - ❖ the **Sophon** model can be helpful in delivering future LHC pheno research! [\[see implementation details on our Github repo\]](#)
 - optimizing sensitivity for dedicated searches/anomaly detection/novel paradigms...
 - ❖ this work demonstrates that it can be a great booster to LHC’s broad resonance search programs
- Stay tuned to their applications to real LHC experiments!

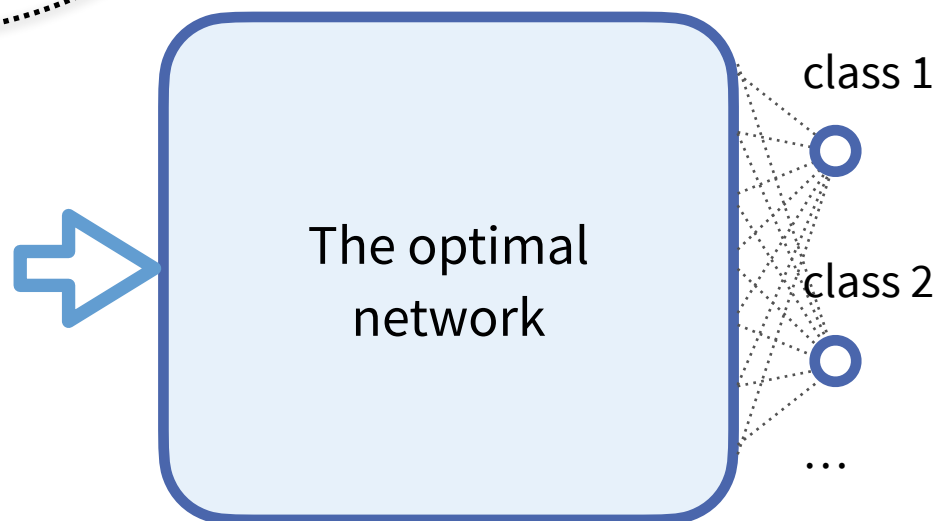
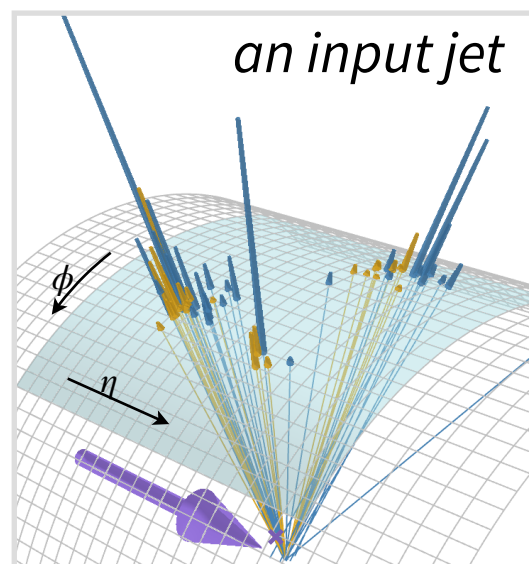
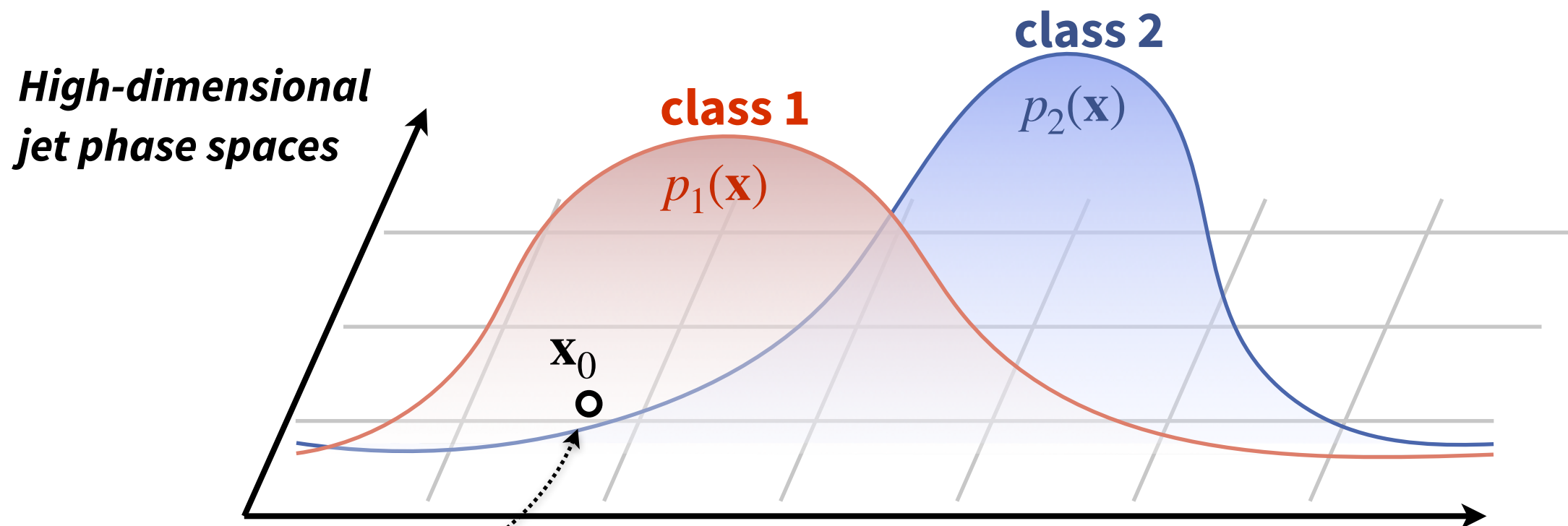


Backup

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**

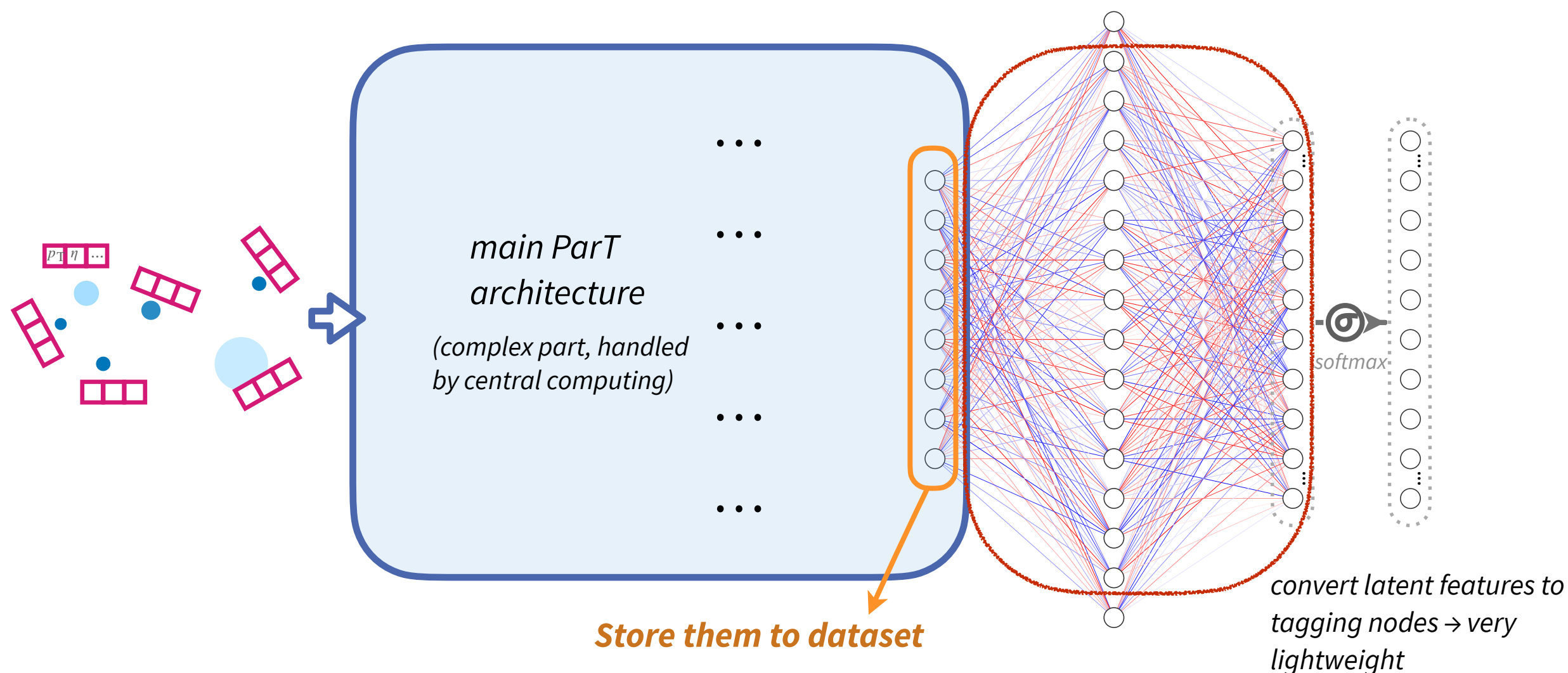


- ❖ Ideal classifier network results in
 $g_1 : g_2 : \dots = p_1(\mathbf{x}_0) : p_2(\mathbf{x}_0) : \dots$
- ❖ It is a direct estimation of p
- ❖ The **network capacity** decides how close the estimation is

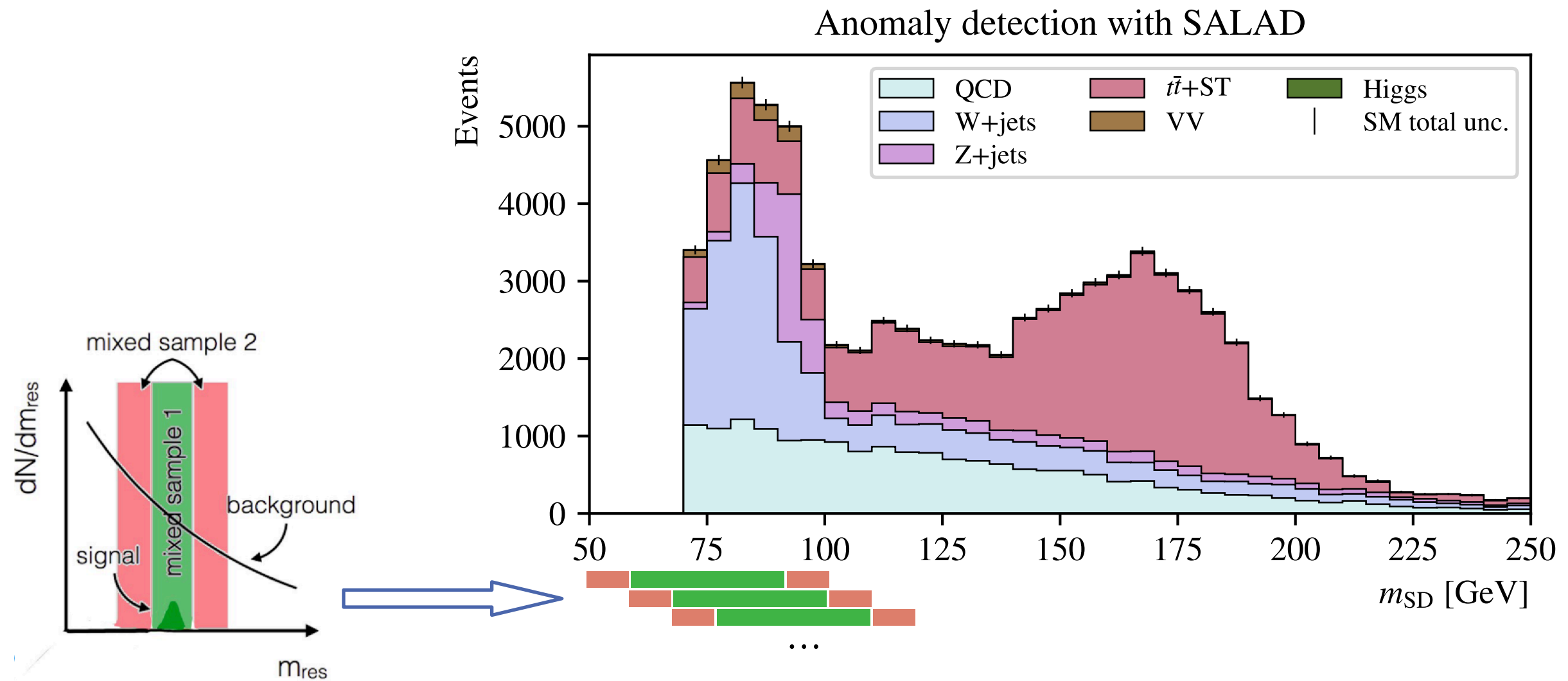
How to deploy the model to LHC experiments?

→ Implies how we can do future analysis

- ❖ hidden layer neurons values are stored in official sample
- ❖ analysis can use them for fine-tuning (equivalently, just think that they are special jet variables)
- ❖ easy to implement & integrate into existing workflow



Sophon's transfer learning × anomaly detection



- Do SALAD (similar to CWOLA Hunt) in each sliding window
 - ❖ purify those peculiar jets in that mass window
- Sophon's latent space has encoded fruitful knowledge on “final-state properties”