

Towards a foundational jet model:

Enhancing generalization with contrastive “gen-reco” pre-training

Zixun Kou Congqiao Li Qiang Li

Peking University

Quantum Computing and Machine Learning Workshop 2025, Qingdao
23 August, 2025

Motivation

Universal model

- Pre-training
according to scaling law, pre-training large models on comprehensive datasets can push a broad range of final states towards their sensitivity frontier
- Achieve state-of-the-art performance
reach the best possible accuracy across all established tasks (e.g. A vs. B tagging, mass/ p_T regression, etc.)
- Ensure strong generalization
reach as better performances as possible for new tasks

Unsupervised approach

Collected a few examples
(sorry for missing your favourite work)

- ❖ **Contrastive learning** via pos/neg samples **JetCLR** **RS3L**
[SciPost Phys. 12, 188 \(2022\)](#) [PRD, 111 \(2025\) 3, 032010](#)
- ❖ **BERT-like**: masked XX modelling **MPM** **MPMv2** **Bumble-bee**
[2401.13537](#) [2409.12589](#) [2412.07867](#)
- ❖ **GPT-like**: next-token prediction **OmniJet-a**
[MLST, 5 035031 \(2024\)](#)
- ❖ **JEPA-like**: predict masked representation **P-JEPA** **J-JEPA** **HEP-JEPA**
[ML4Jets talk](#) [2412.05333](#) [2502.03933](#)

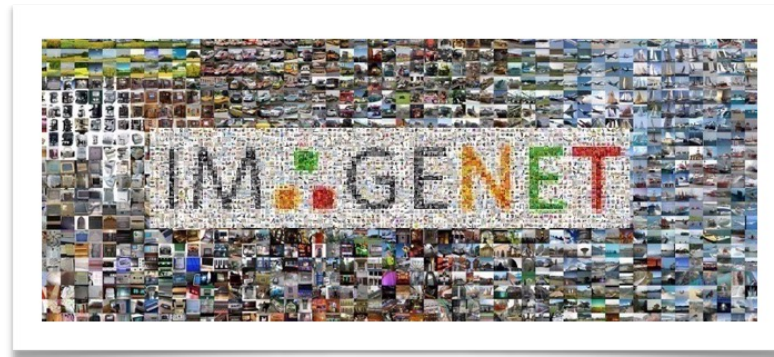
Supervised approach (+X)

Sophon
w/ **JetClass-II**
[2405.12972](#)

Motivation

Computer Vision

- models trained on ImageNet were among the **earliest pretrained CV models**. (serving to **generalize to other CV domains**)
- Modern self-supervised learning (SSL) methods (e.g. MAE, I-JEPA, ...) show strong performance and beat SL.
- but one fact is that the **supervised baseline is relatively weak** (ImageNet-1k only has 1M images)



ImageNet-1k
the cornerstone of modern CV
1M dataset; 1000 labels

HEP dataset

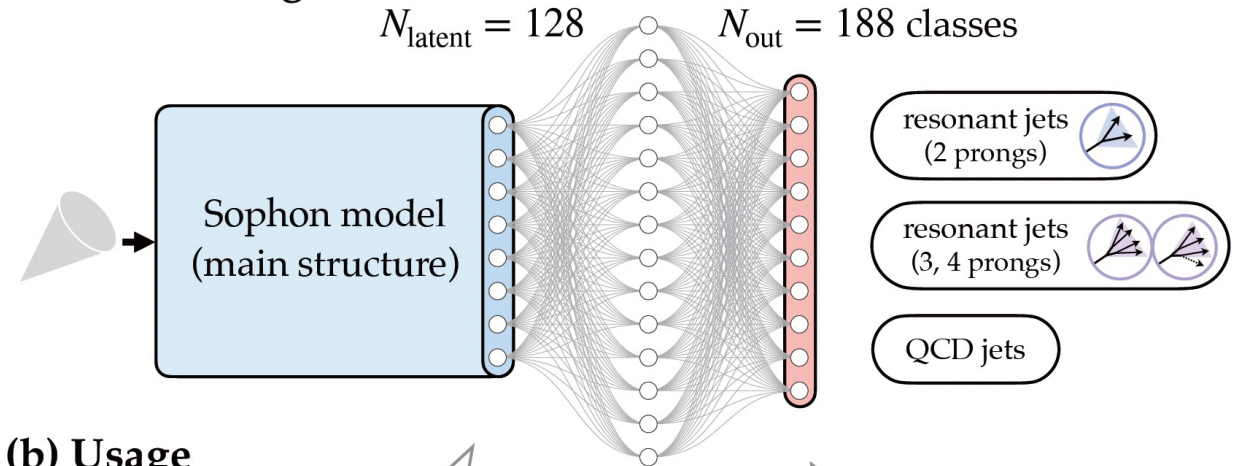
- Establish a **very strong baseline**
modern supervised models (e.g. SoTA taggers in ATLAS/CMS) are already trained on $\mathcal{O}(100\text{M})$ -level datasets.

Motivation

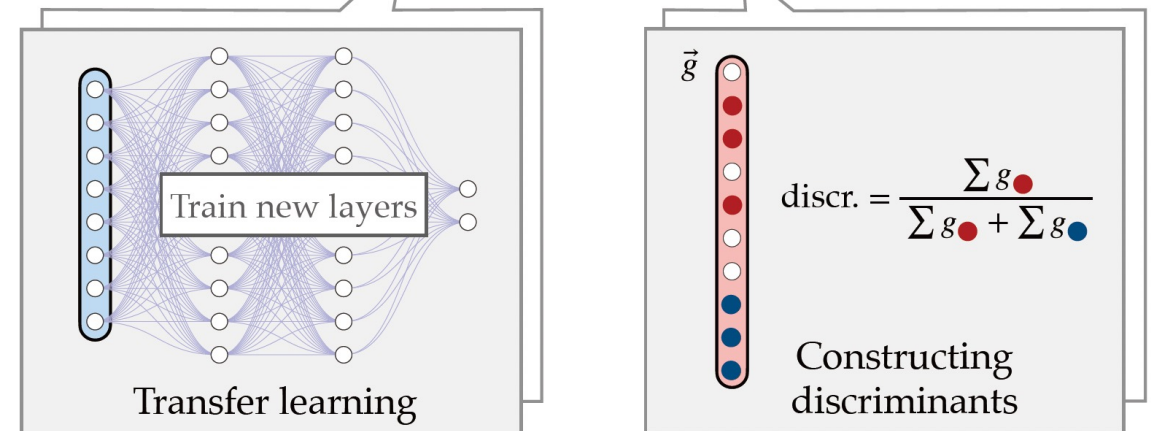
Current *Sophon* model

- *Sophon*: Signature-Oriented Pre-training for Heavy-resonance ObservationN
- Key concept: Pre-training an expressive Transformer model on a wide range of jet phase spaces on a multiclass classification task
- Target boosted-jet phase-space explored simulate datasets
QCD, resonant jets with 2, 3 or 4 prongs (188 categories)

(a) Pre-training



(b) Usage

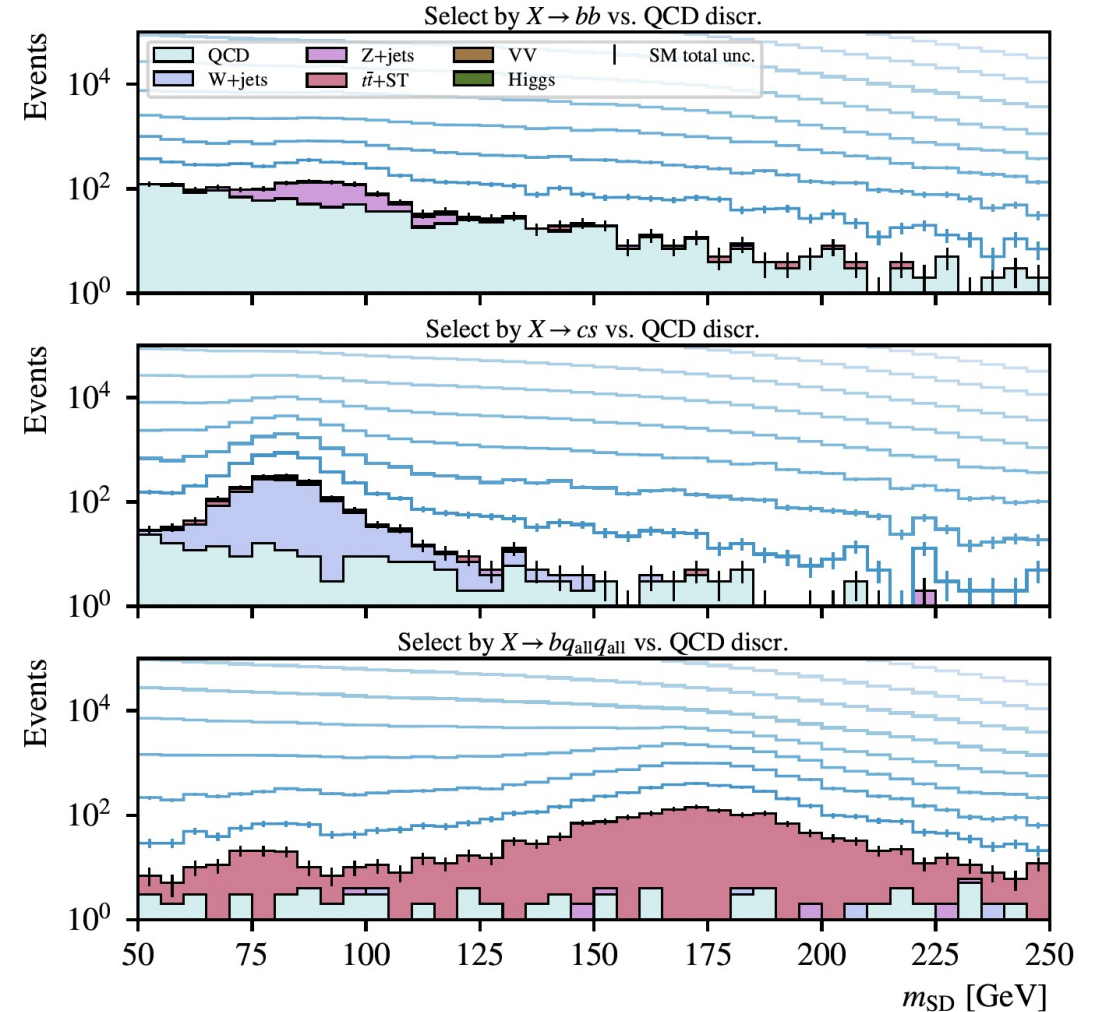


Motivation

- Provides directly usable discriminants and outperforms SOTA results
- Demonstrate broad ability to construct discriminants and sensitively probe resonances with unknown properties (QCD, V+jets, ttbar+single top, di-boson, Higgs production)

Way to improve

- A simple classification approach
- Jet signatures of initial state remain unused
- Enhance transfer-learning ability in other tasks



Next Step

Supervised approach + ***X?***



“Signature-Oriented Pretraining” ↔ *pretraining with labels*
serving as a pretraining foundation

Sophon



Further enhance model
generalization with some SSL/...

Next Step

Supervised approach + ***X?***



*“Signature-Oriented Pretraining” ↔ pretraining with labels
serving as a pretraining foundation*

Sophon



*Further enhance model
generalization with some SSL/...*

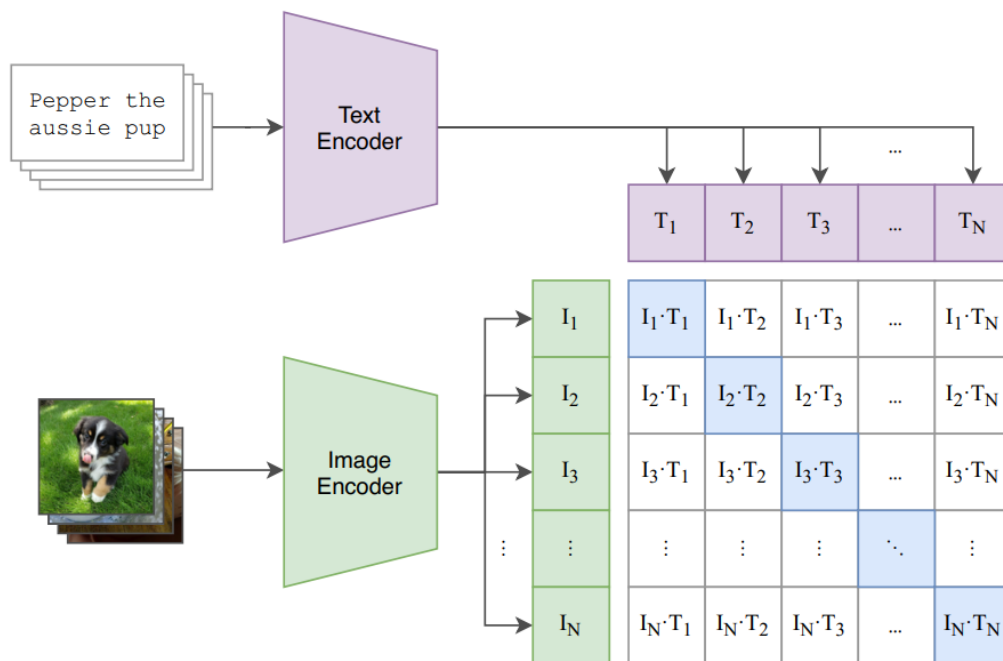
Sophon++

Next Step

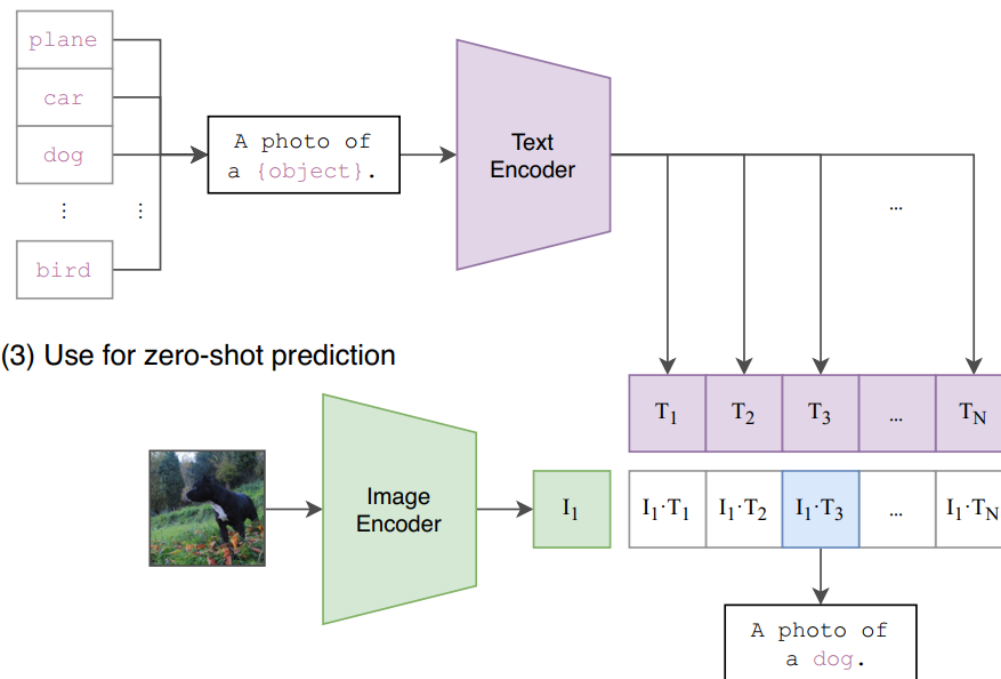
CLIP

- Contrastive Language-Image Pre-training
- Multi-Modality(language, image)

(1) Contrastive pre-training



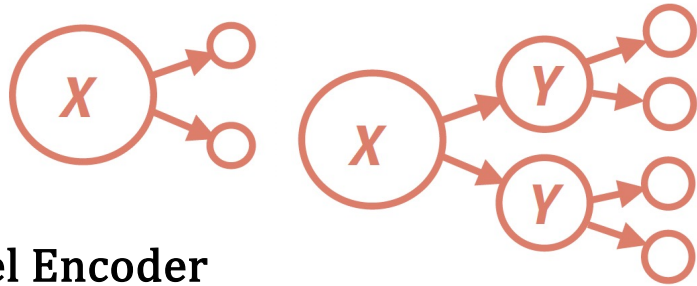
(2) Create dataset classifier from label text



Pre-training Setup

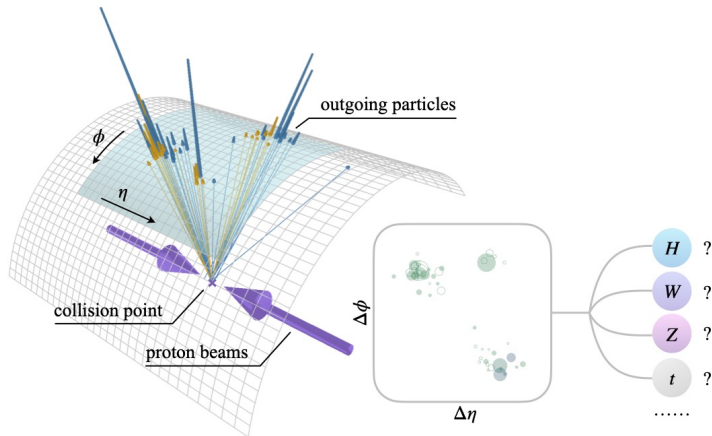
- **Gen-Level Encoder**

Input: features of dedicated final-state quarks and leptons



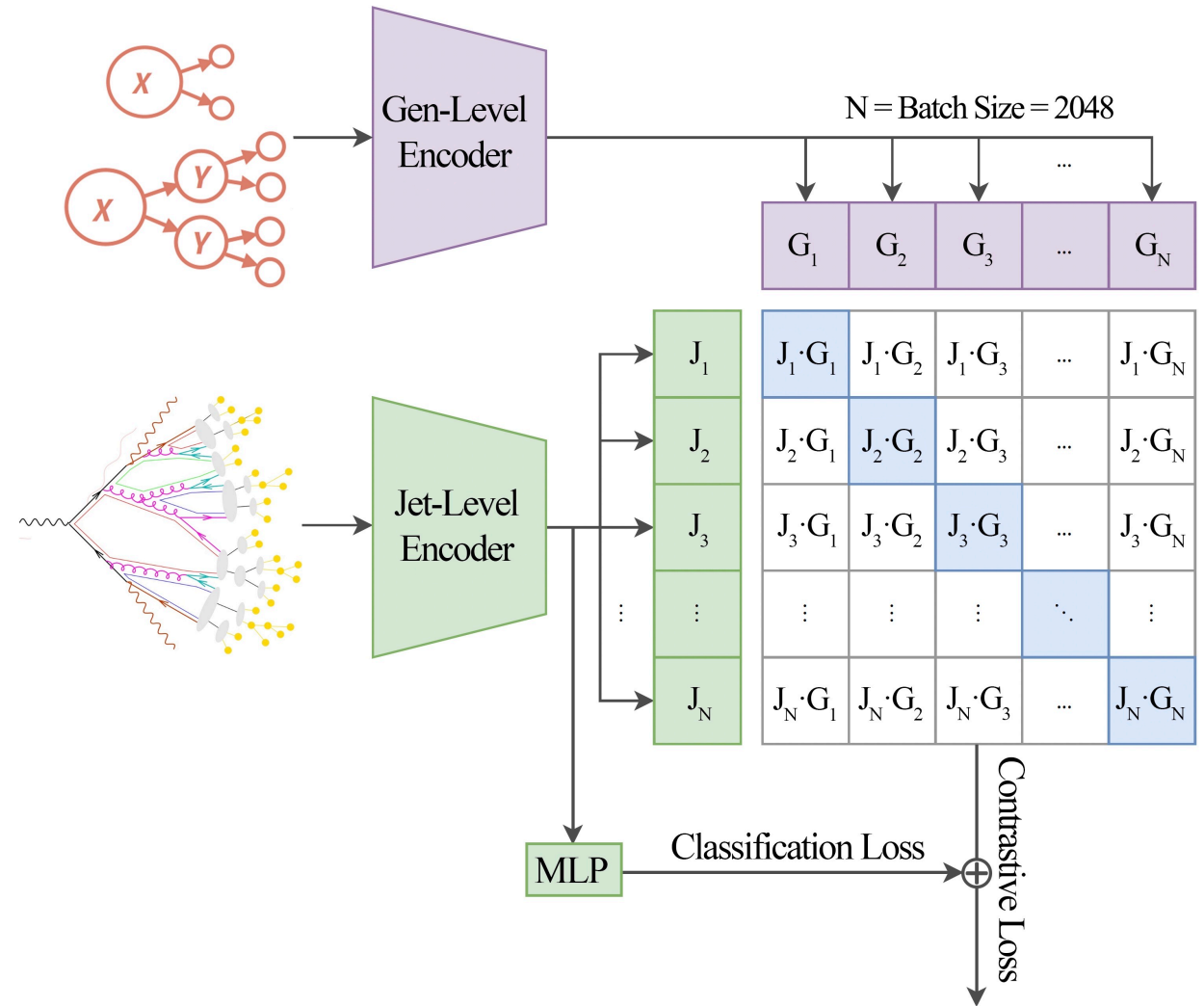
- **Jet-Level Encoder**

Input: features of jet constituents (same as original *Sophon*)



Class tokens: separated token for each task
Network structure: ParT for both

2025/8/22



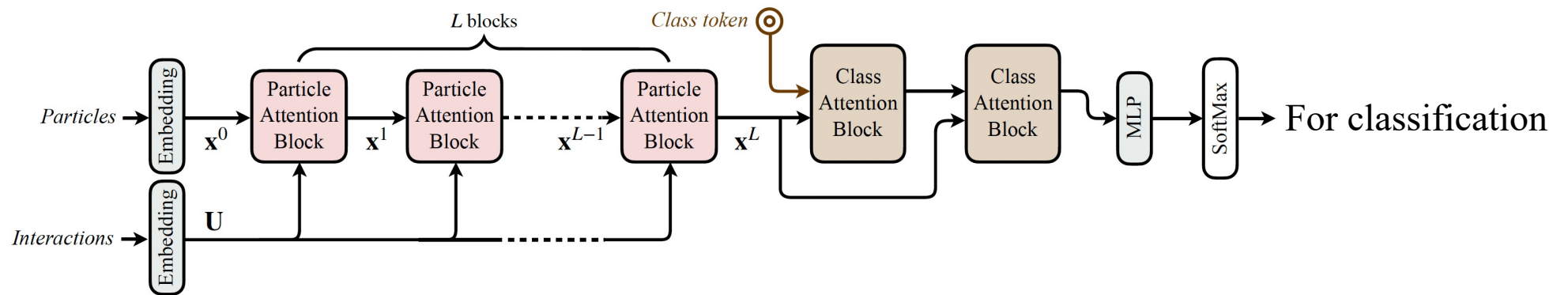
$$\text{loss} = \alpha * \text{contrastive loss} + \beta * \text{classification loss}$$

Pre-training Setup

Network parameters:

Encoder	Particle Embedding	Pairwise Embedding	Particle Attention Block	Heads	Class Attention Block	FC
Jet-level	(512, 128, 512)	(64, 64, 64, 8)	8	8	2 (for classification) 2 (for contrastive)	(512, 188) (512, 512, 512)
gen-level	(64, 64, 64)	(32, 32, 32, 4)	4	4	2	(256)

Structure of
Jet-level:



Training parameters(*Sophon++* dev):

$\alpha=0.1, \beta=1$

Batch size=2048

Learning rate= 5×10^{-3}

Steps per epoch=5000 (1024 * 10,000/2048)

Epoch=180

Use NCCL on 4 GPU

2025/8/22

Training parameters(original *Sophon*):

$\alpha=0, \beta=1$

Batch size=2048

Learning rate= 3×10^{-3}

Steps per epoch=5000

Epoch=160

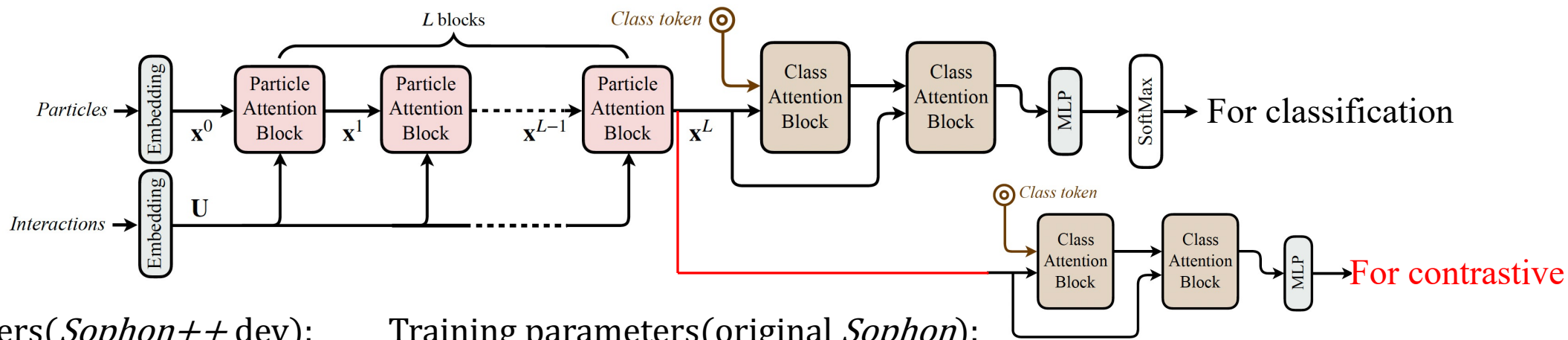
Use NCCL on 4 GPU

Pre-training Setup

Network parameters:

Encoder	Particle Embedding	Pairwise Embedding	Particle Attention Block	Heads	Class Attention Block	FC
Jet-level	(512, 128, 512)	(64, 64, 64, 8)	8	8	2 (for classification) 2 (for contrastive)	(512, 188) (512, 512, 512)
gen-level	(64, 64, 64)	(32, 32, 32, 4)	4	4	2	(256)

Structure of
Jet-level:



Training parameters(*Sophon++* dev):

$\alpha=0.1, \beta=1$

Batch size=2048

Learning rate= 5×10^{-3}

Steps per epoch=5000 (1024 * 10,000/2048)

Epoch=180

Use NCCL on 4 GPU

2025/8/22

Training parameters(original *Sophon*):

$\alpha=0, \beta=1$

Batch size=2048

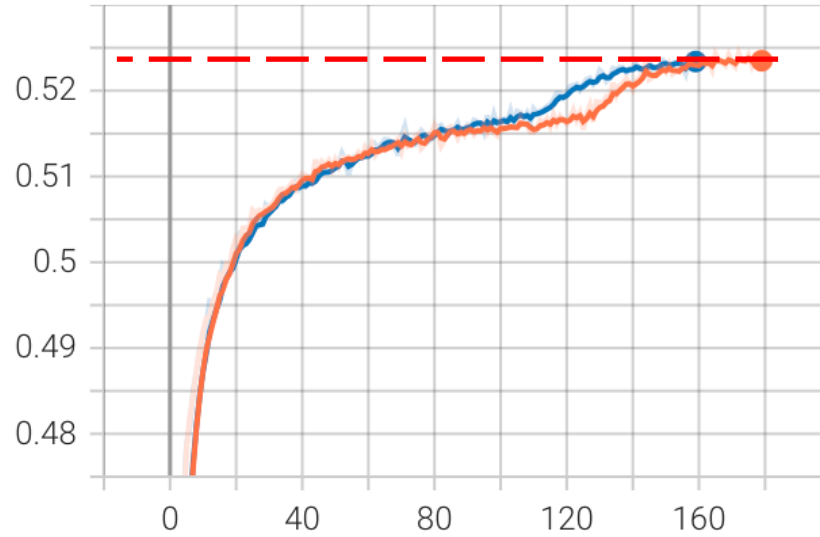
Learning rate= 3×10^{-3}

Steps per epoch=5000

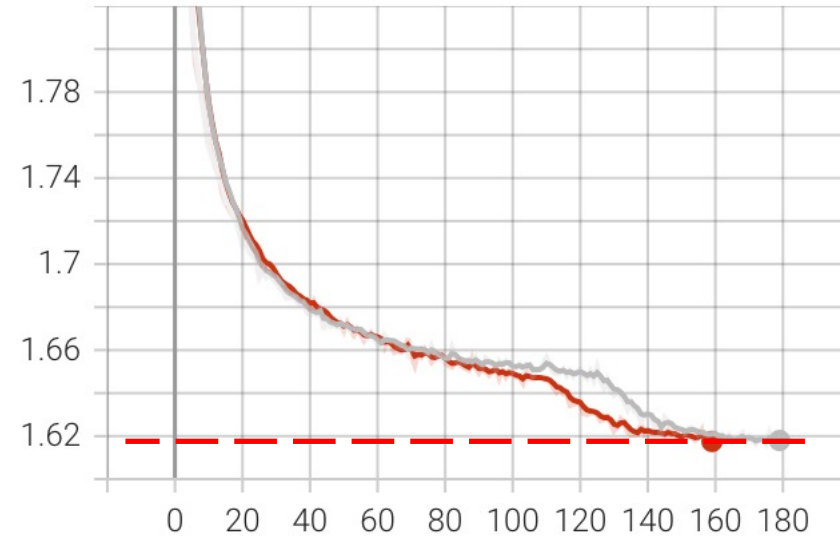
Epoch=160

Use NCCL on 4 GPU

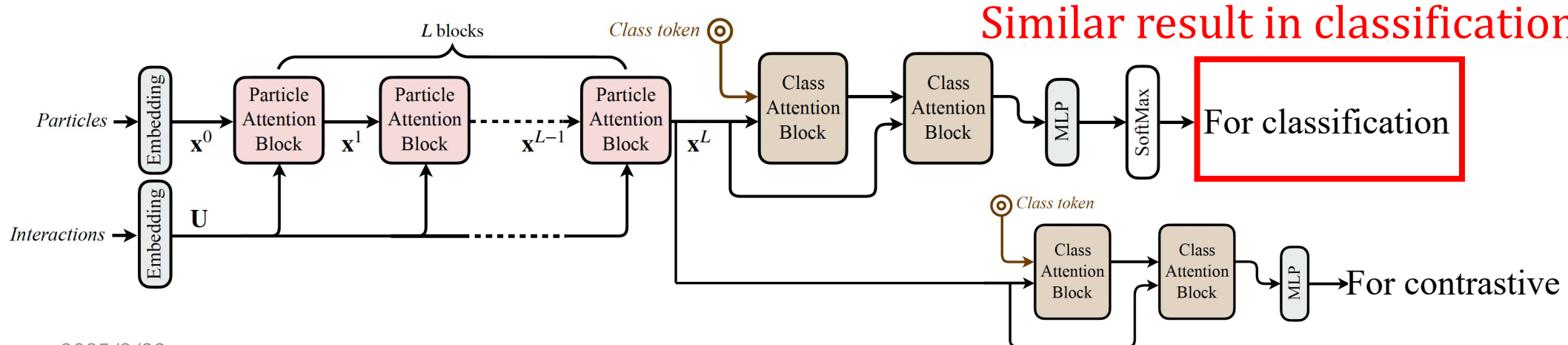
Result of classification (188 categories)



Accuracy



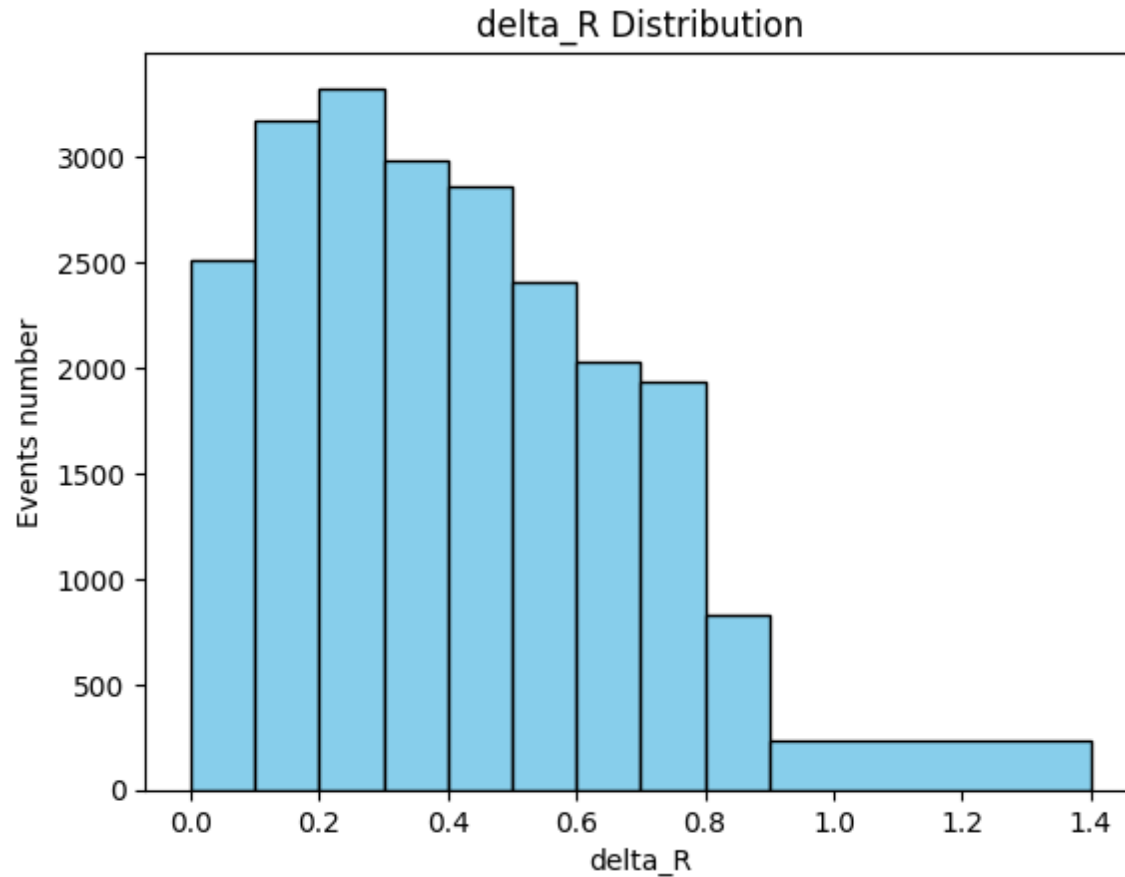
Classification Loss



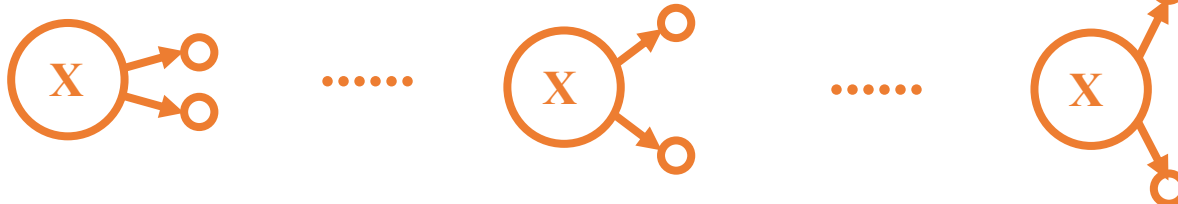
Verification Experiments

(To testify global optimization of *Sophon++*)

Verification Experiment 1



Distribution of angle between bb



Modifying data config(yaml card):

- Only input events labeled as 'X to bb'
- Define new variables to calculate bb_deltaR
- Assign indices for each event according to bb_deltaR
- Set the indices as the new label

0.0~0.1: 0

0.1~0.2: 1

...

0.8~0.9: 8

>0.9: 9



A 10-class classification
model

Verification Experiment 1

Universal model selection:

Model	Epoch Selection	Accuracy	Classification Loss	Contrastive Loss
<i>Sophon</i>	158	0.5241	1.619	-
<i>Sophon++</i> dev	179	0.5230	1.618	0.4922

Training parameters:

Batch size=1024

Learning rate= 1×10^{-3}

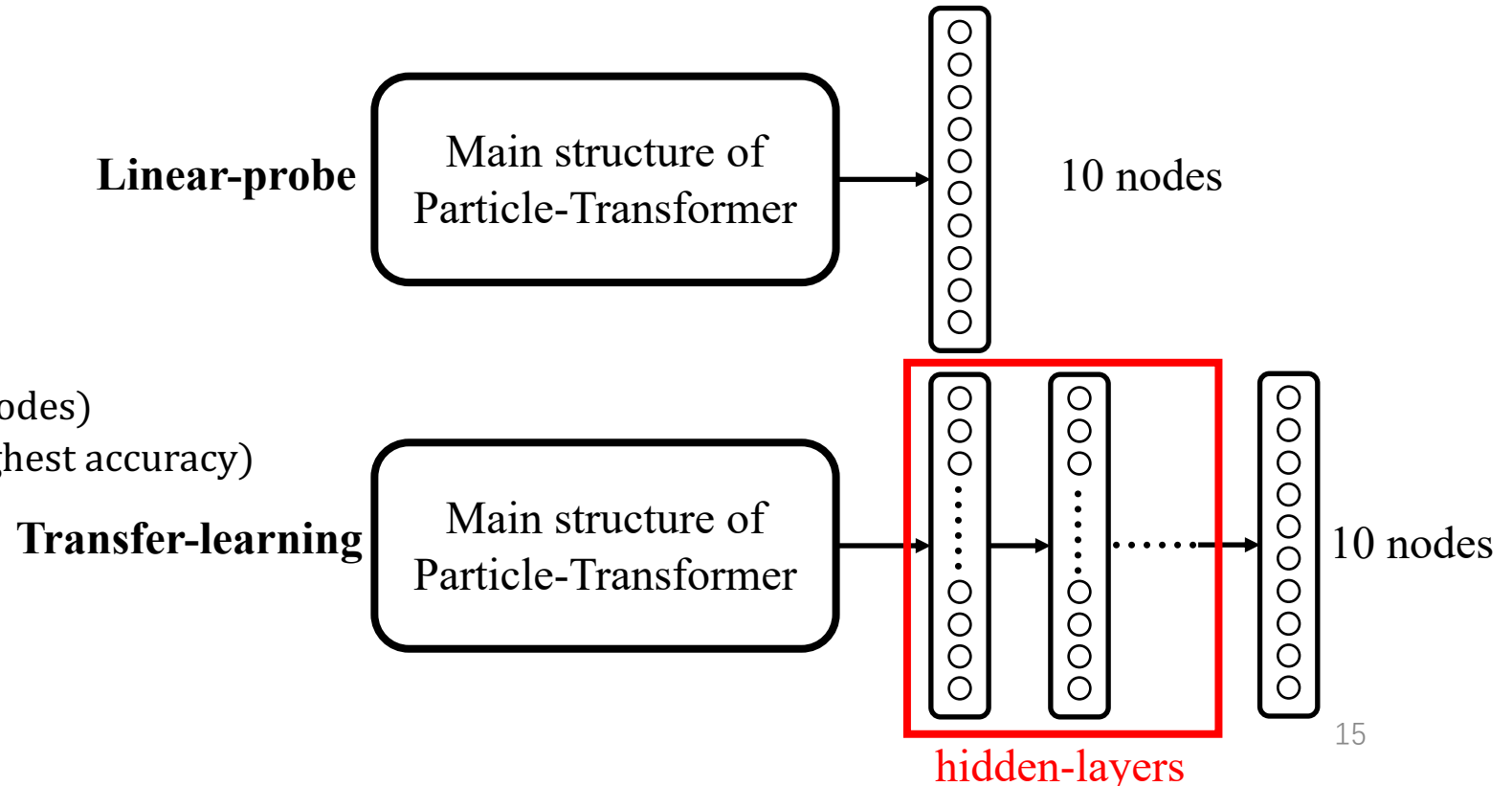
Steps per epoch=5000 ($1024 * 5,000 / 1024$)

Epoch=20 (10 for Linear-probe)

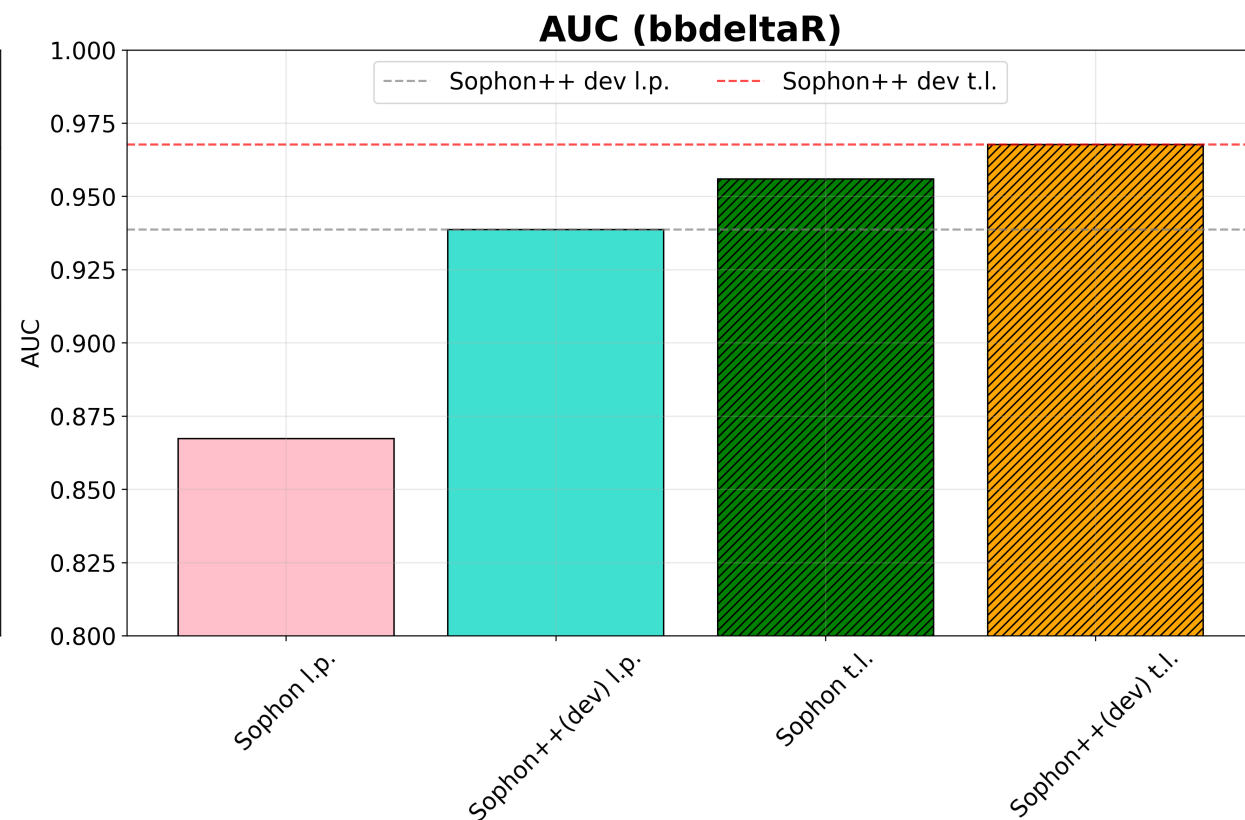
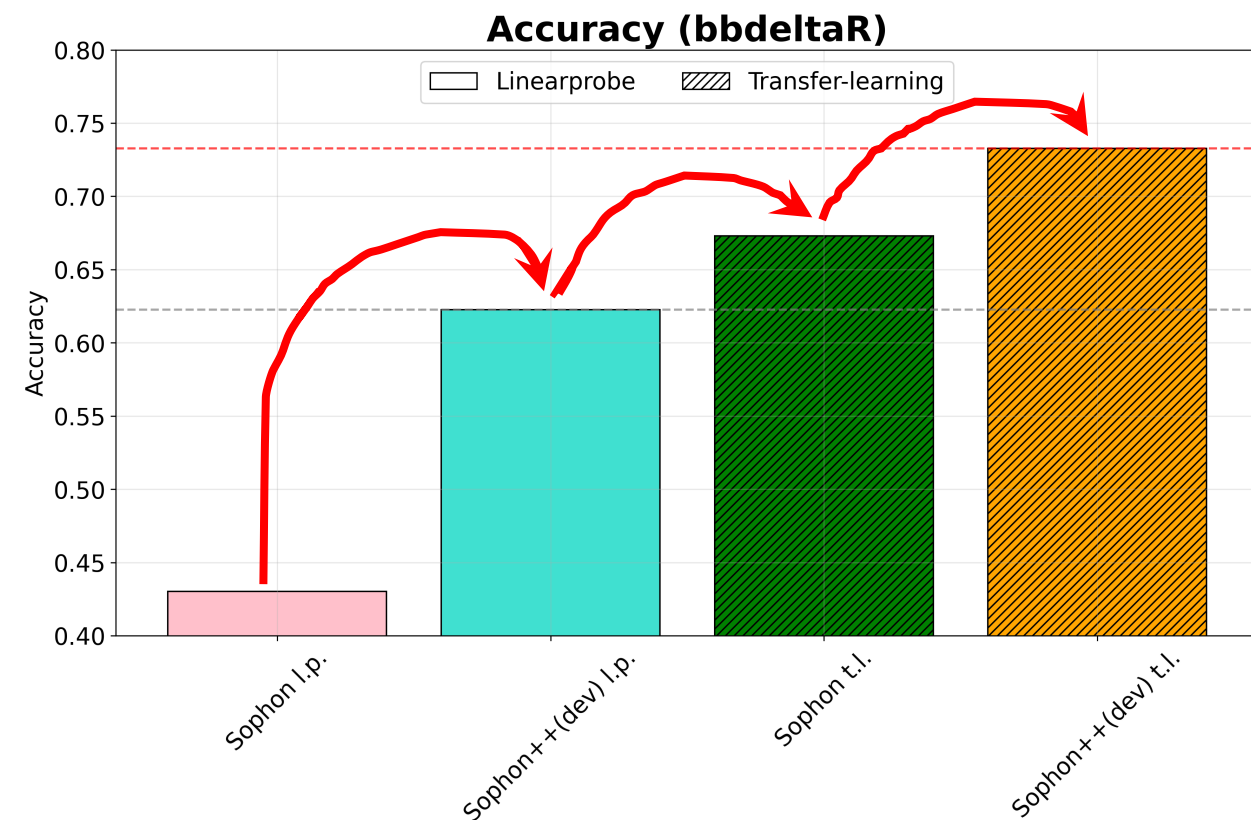
Linear-probe (train a one-layer FC with 10 nodes)

Transfer-learning (train a best FFN to get highest accuracy)

1 GPU per training

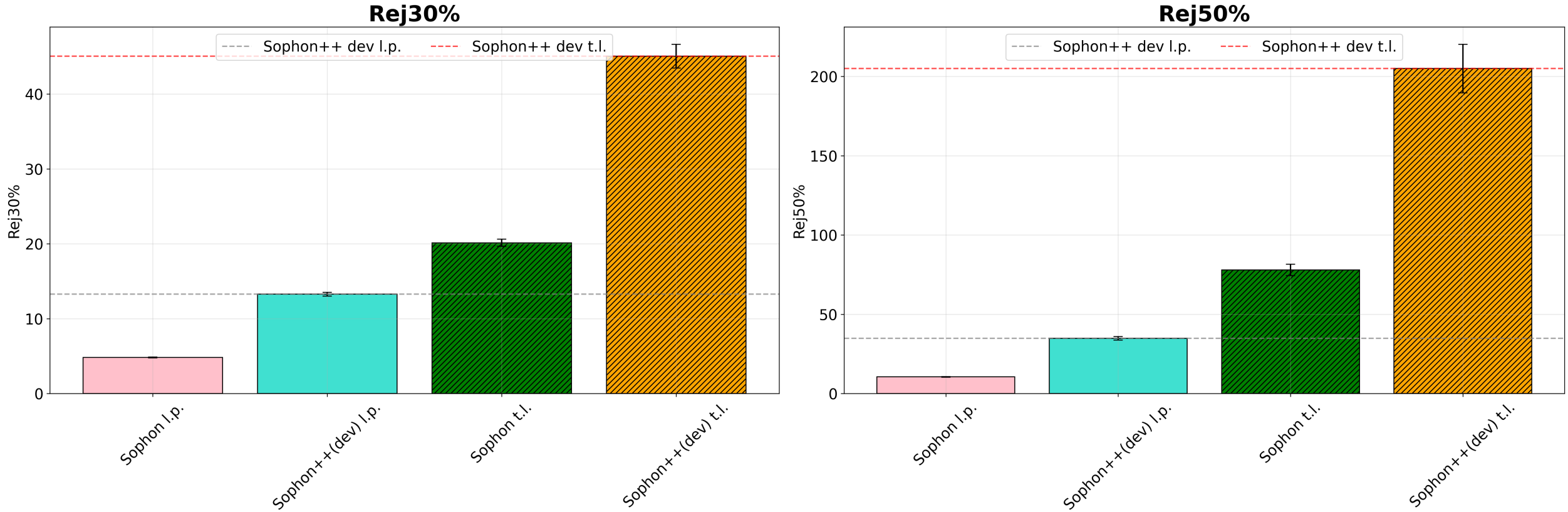


Verification Experiment 1



performance of *Sophon++* is much better than Original Sophon!

Verification Experiment 1



bin 4 vs bin 5

performance of *Sophon++* is much better than Original Sophon!

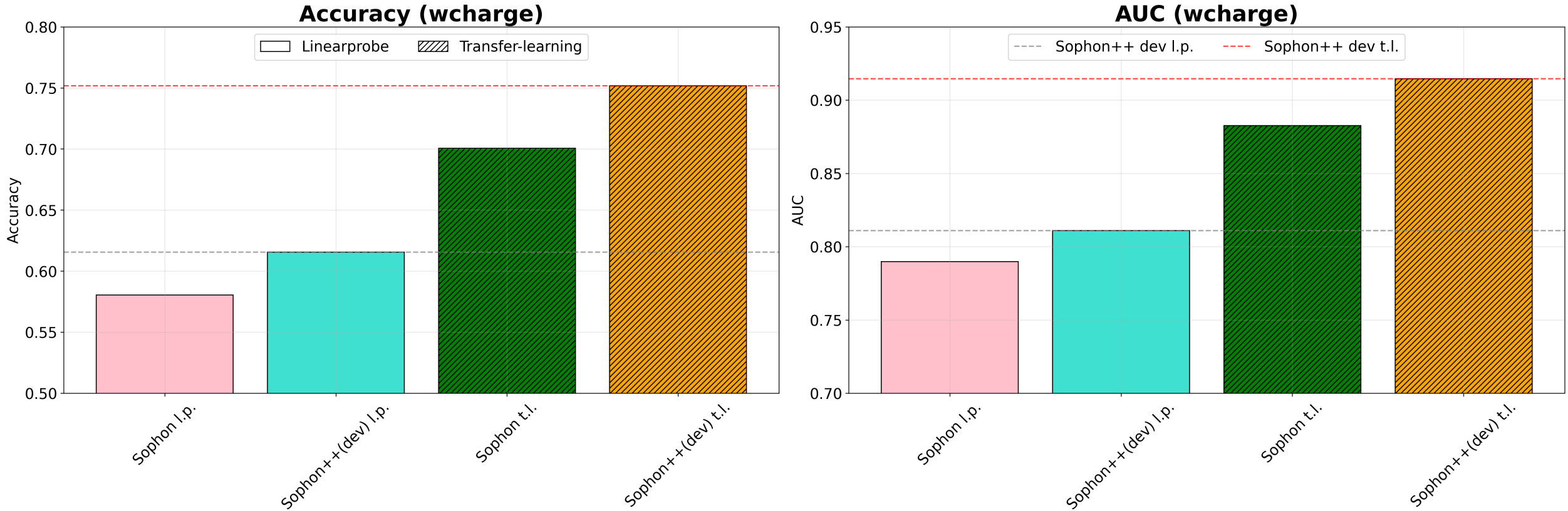
Verification Experiment 2: hadronic W^\pm tagging

- While Lorentz-boosted hadronic W/Z jet tagging has been considered in Run 1, recent deep learning advances now enable discrimination of the W boson charge
 - traditional approach: uses jet charge (p_T weighted particle charge) as a discriminant
 - opportunity: integrate charge determination into the standard deep learning jet tagging framework

e.g. as explored in [\[PRD 101, 053001 \(2020\)\]](#)



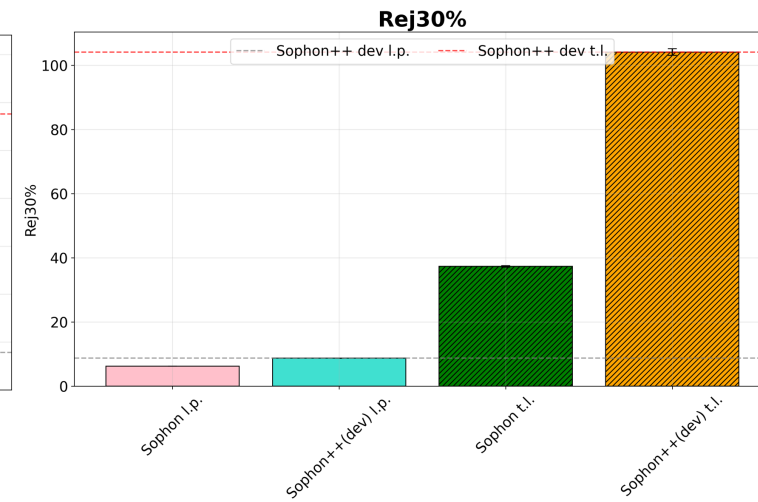
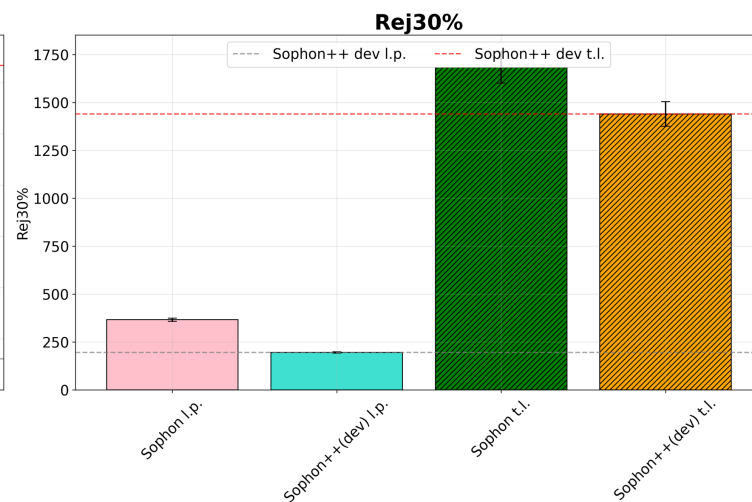
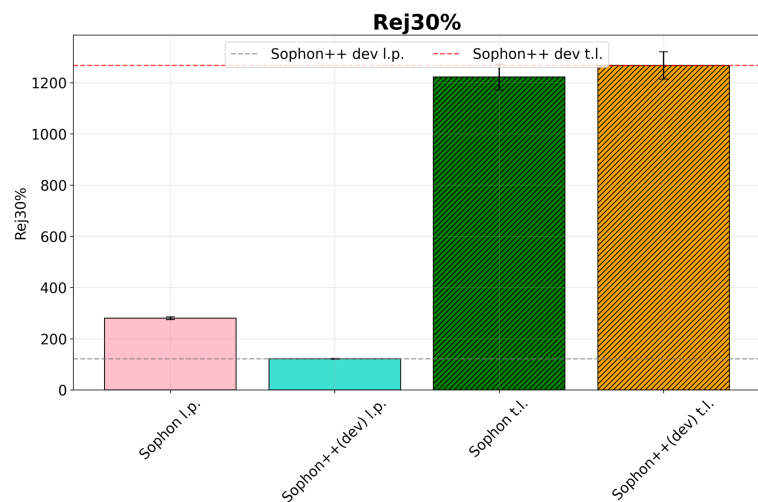
Verification Experiment 2: hadronic W^\pm tagging



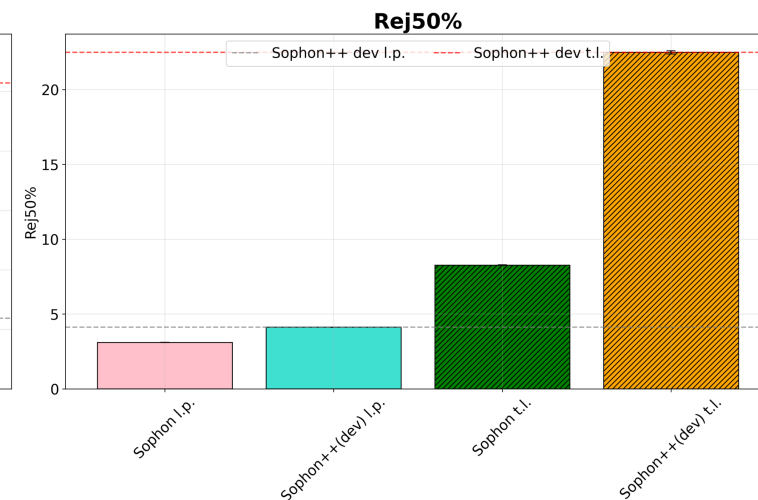
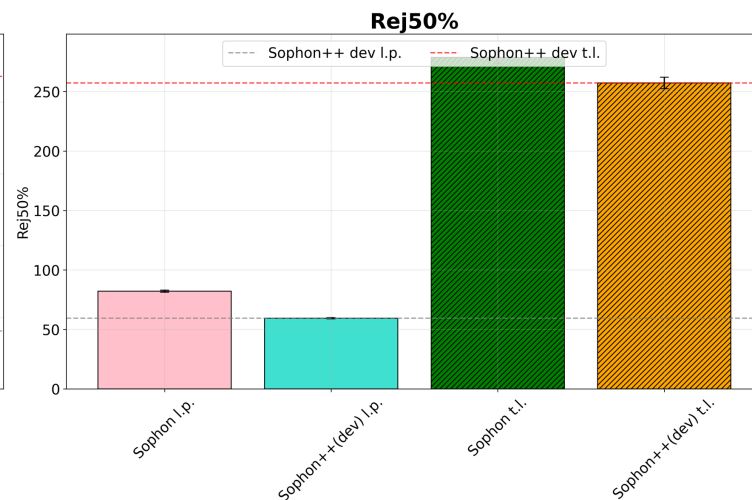
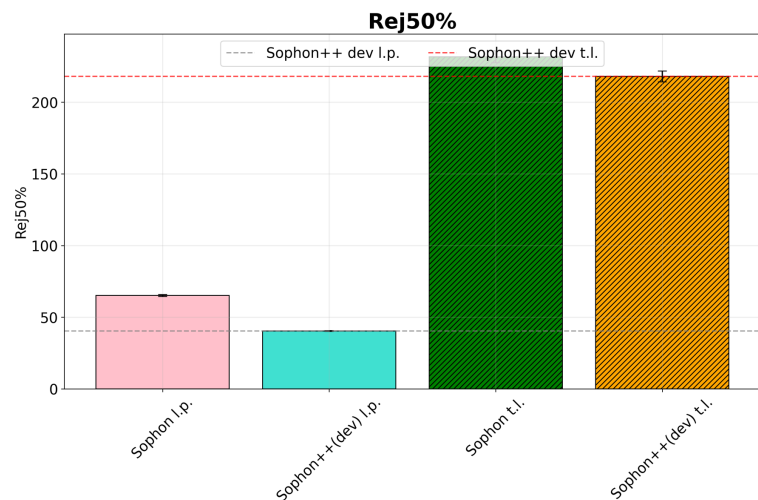
performance of *Sophon++* is better than Original Sophon

Verification Experiment 2: hadronic W^\pm tagging

Rej30%



Rej50%



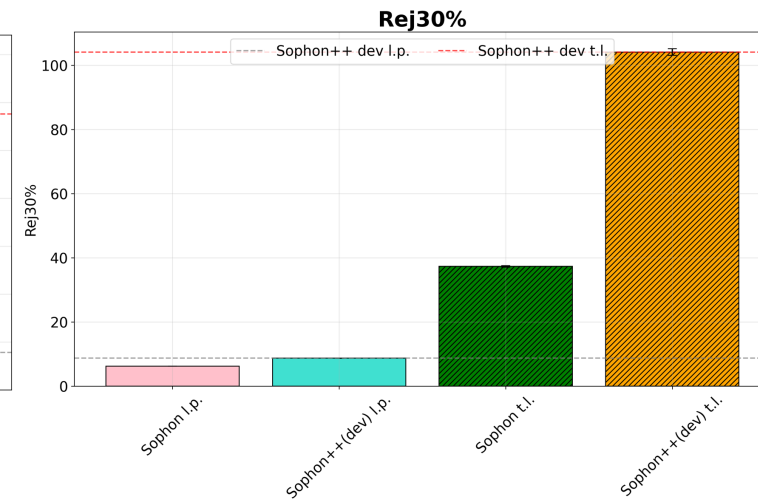
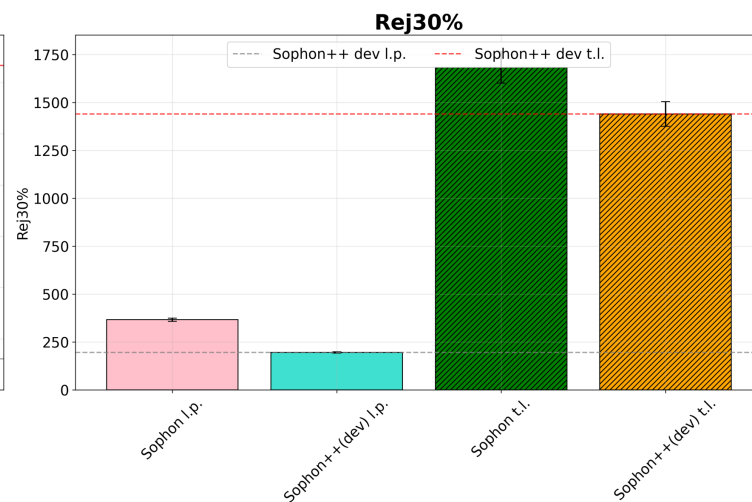
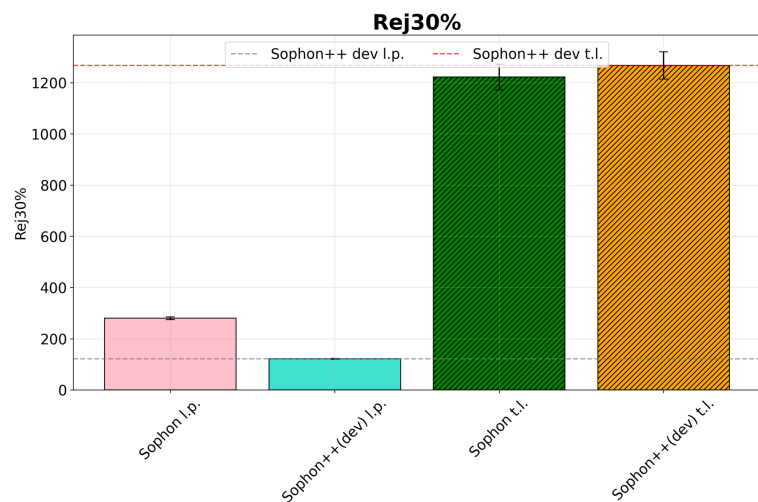
W^+ vs QCD

W^- vs QCD

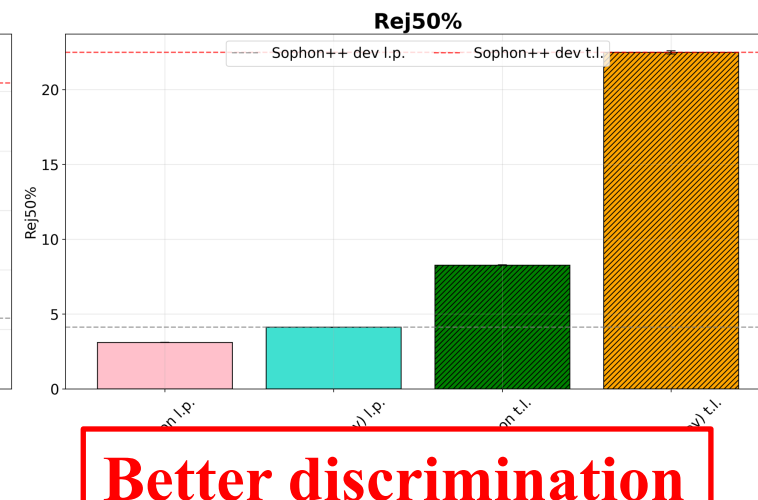
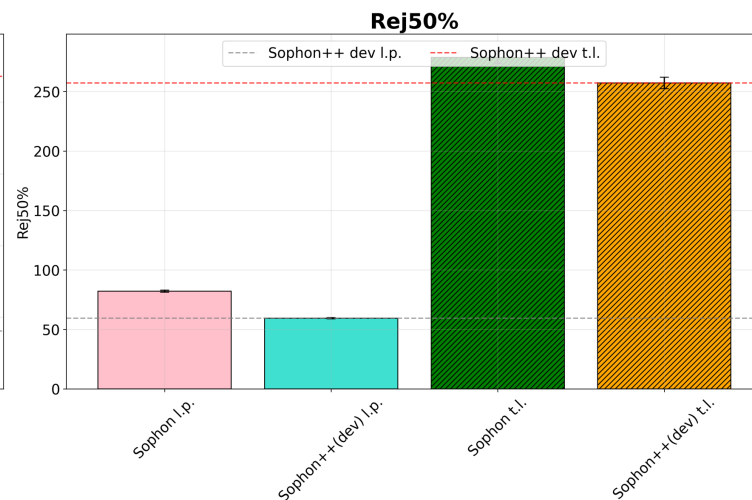
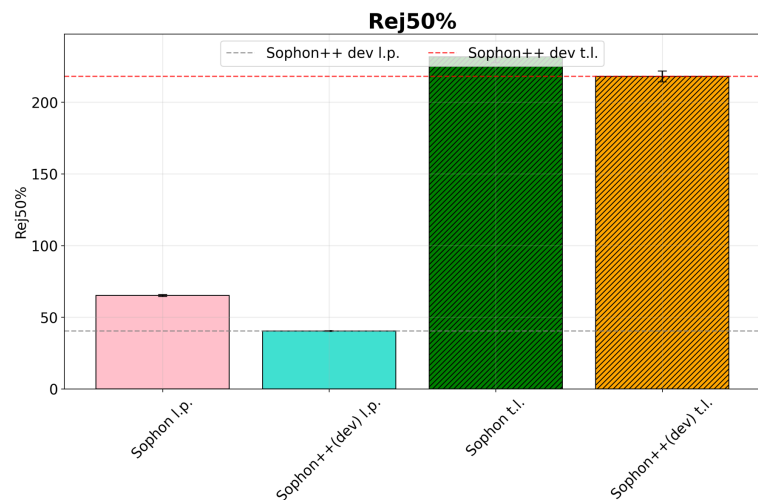
W^+ vs W^-

Verification Experiment 2: hadronic W^\pm tagging

Rej30%



Rej50%



Better discrimination

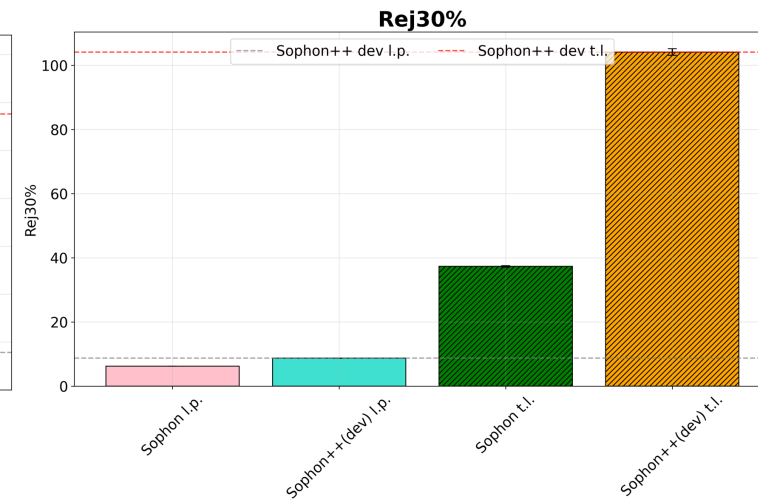
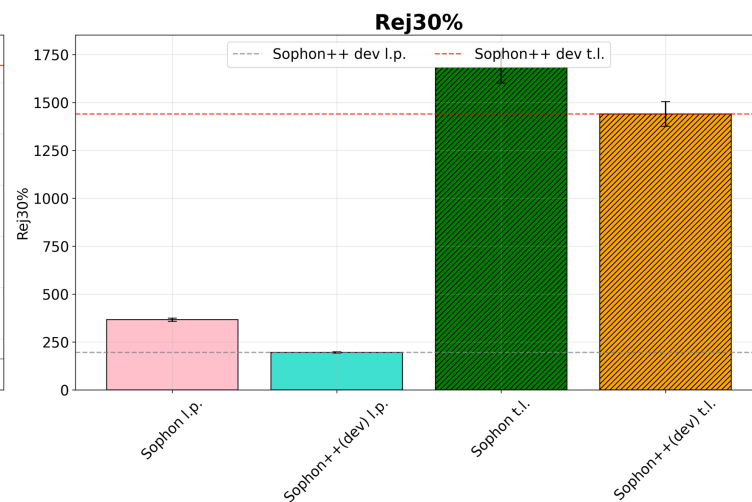
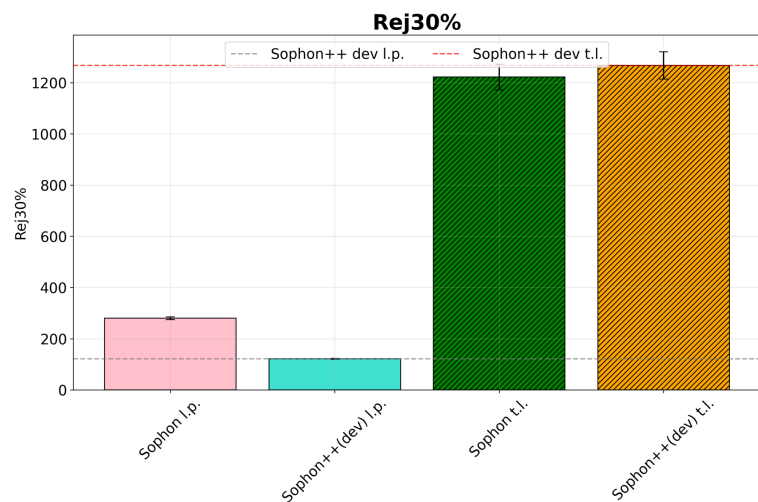
W^+ vs QCD

W^- vs QCD

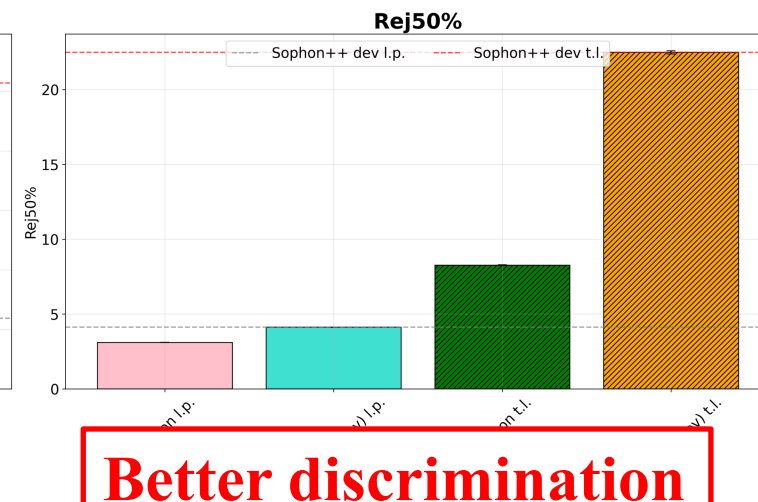
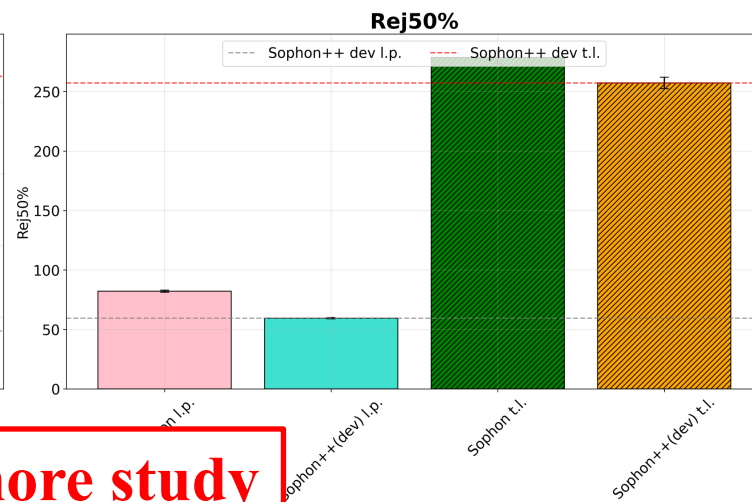
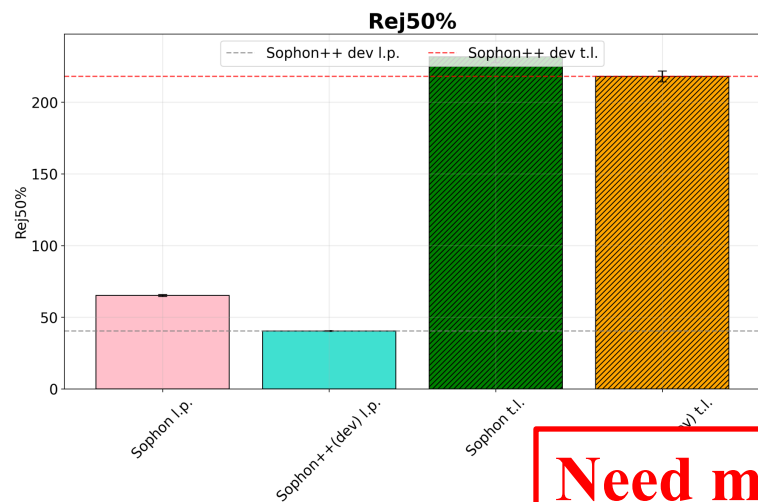
W^+ vs W^-

Verification Experiment 2: hadronic W^\pm tagging

Rej30%



Rej50%



Need more study

Better discrimination

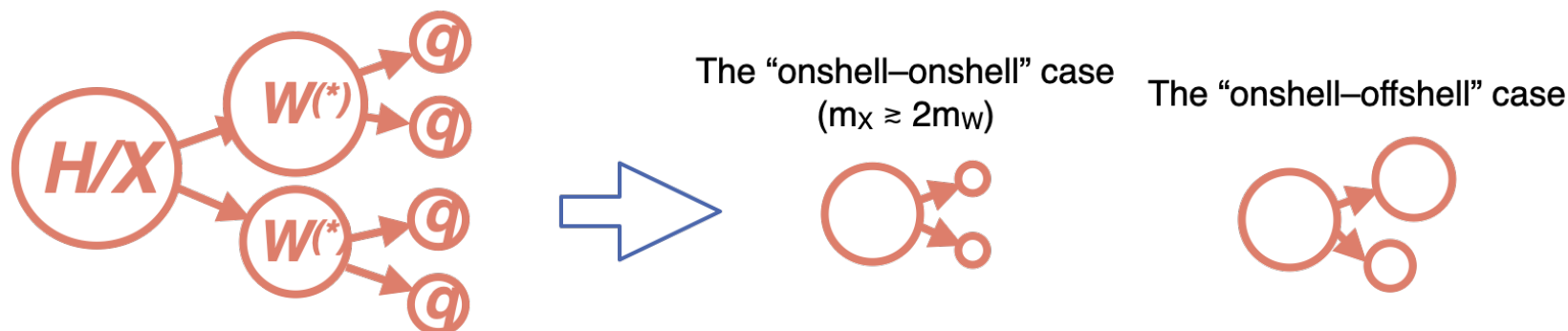
W^+ vs QCD

W^- vs QCD

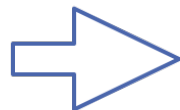
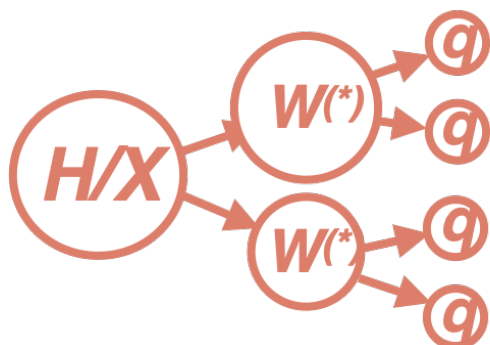
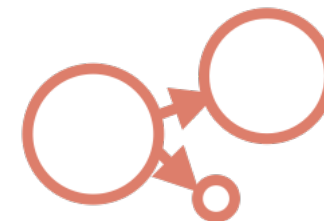
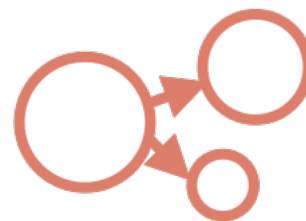
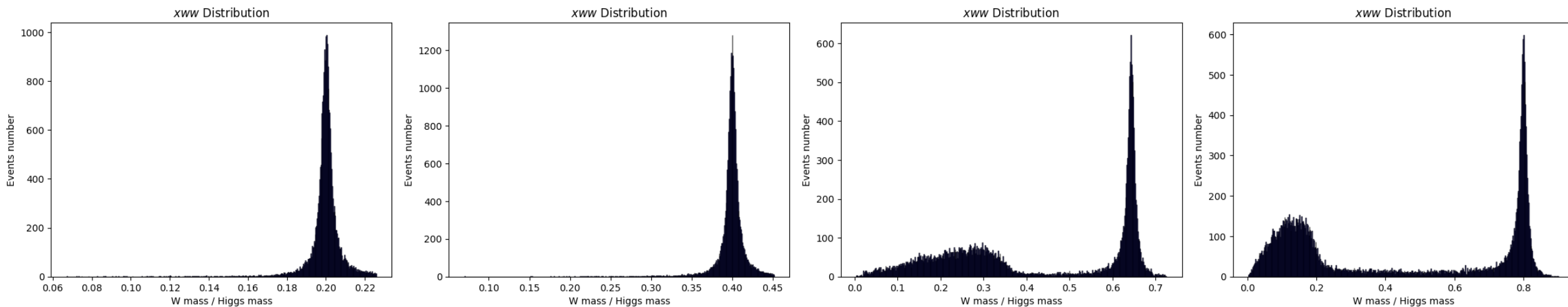
W^+ vs W^-

Verification Experiment 3: $X \rightarrow WW^{(*)}$ tagging

- Discriminating $X \rightarrow WW^{(*)}$ signals across varying m_W/m_X and against QCD backgrounds
- Exploring $H/X \rightarrow WW^{(*)}$ in a fully-merged topology is an emerging LHC direction
 - multiple CMS results now available: SM (single/di-Higgs), and BSM searches
SM $HH \rightarrow bbWW^{*}(4q)$ [CMS-PAS-HIG-23-012], $H \rightarrow WW^{*}$ [CMS-PAS-HIG-24-008]
Resonant $X \rightarrow H(bb)Y(WW)$ [CMS-PAS-B2G-23-007]
- Different m_W/m_X ratios probe distinct phase space regions \rightarrow require separately optimized discriminants
- We consider four benchmark scenarios
 - $m_W/m_X = 0.2, 0.4, 0.6432$ (SM case), 0.8



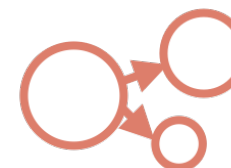
Verification Experiment 3: $X \rightarrow WW(*)$ tagging



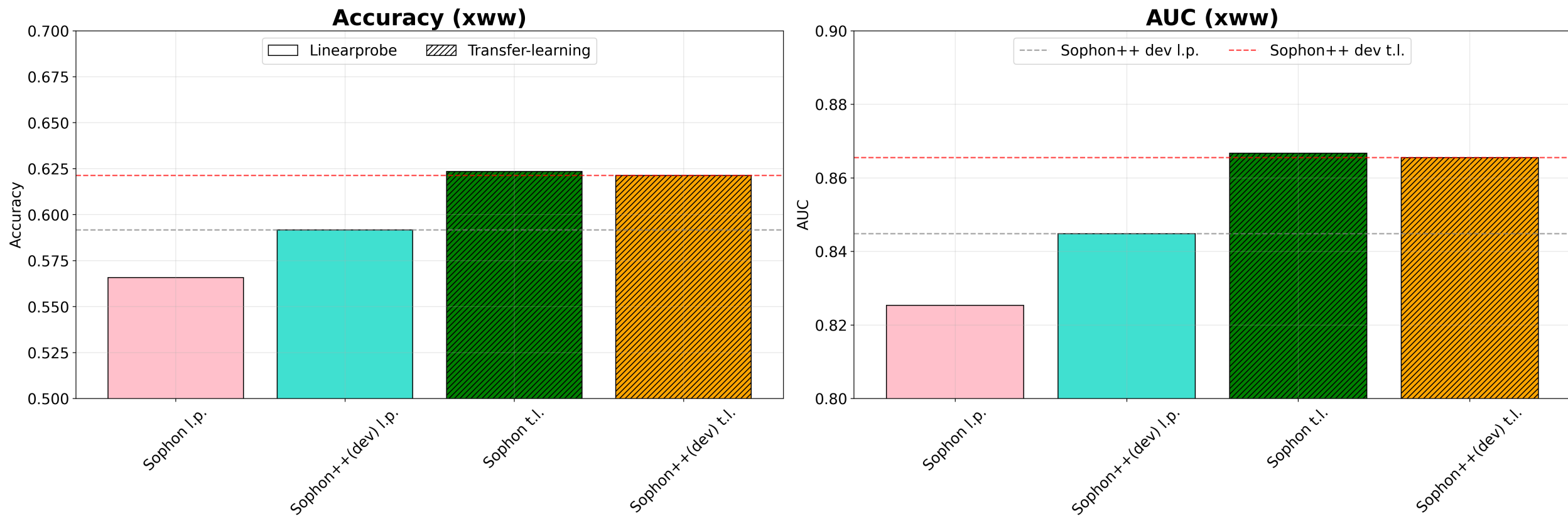
The "onshell–onshell" case
($m_X \geq 2m_W$)



The "onshell–offshell" case



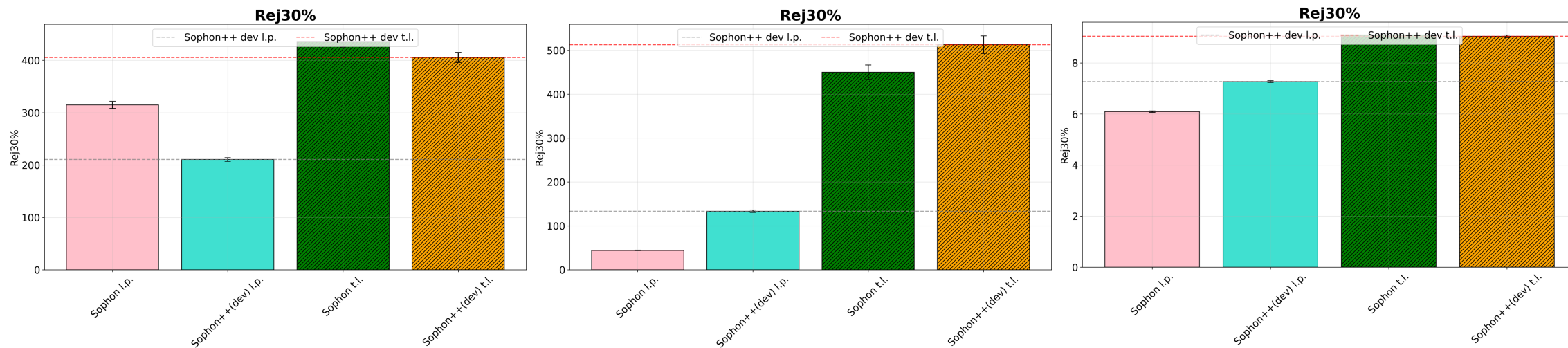
Verification Experiment 3: $X \rightarrow WW(*)$ tagging



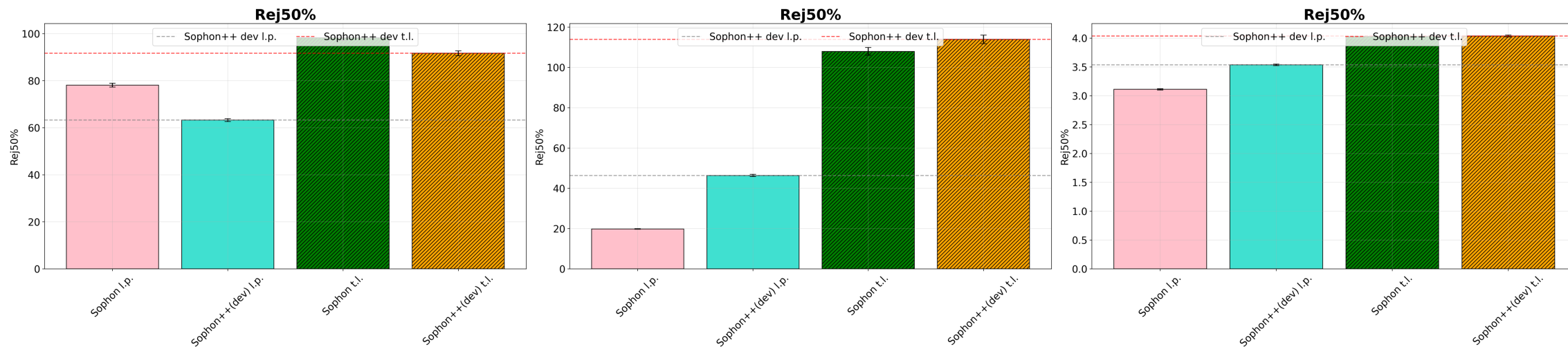
performance of *Sophon++* is better than Original Sophon in Linear-probe

Verification Experiment 3: $X \rightarrow WW(*)$ tagging

Rej30%



Rej50%



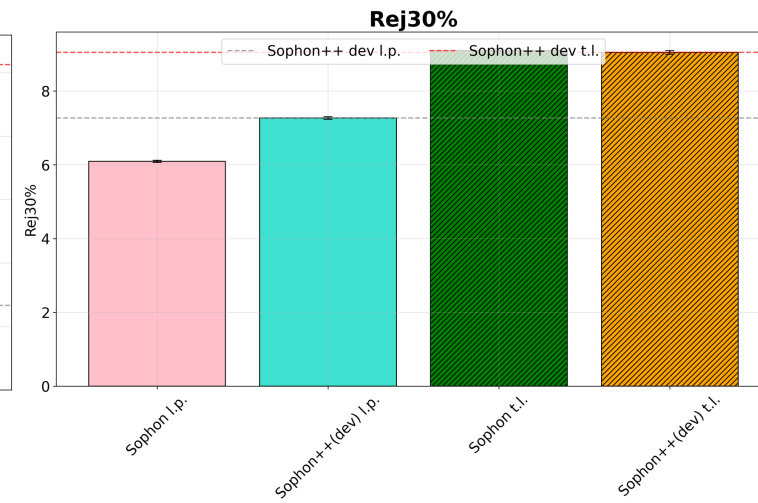
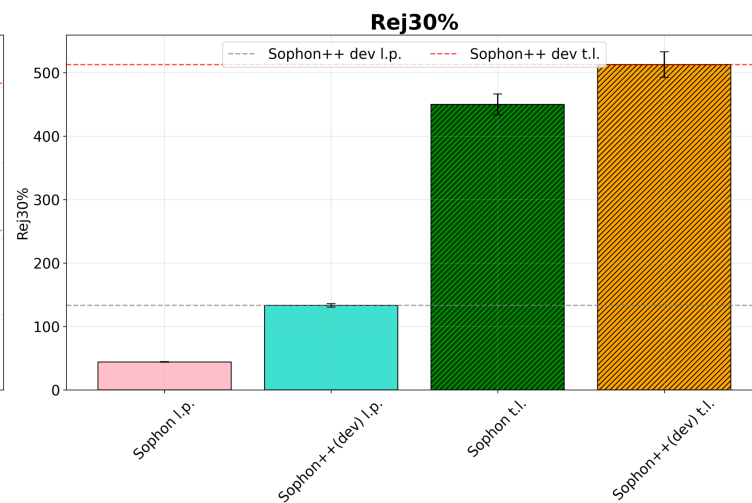
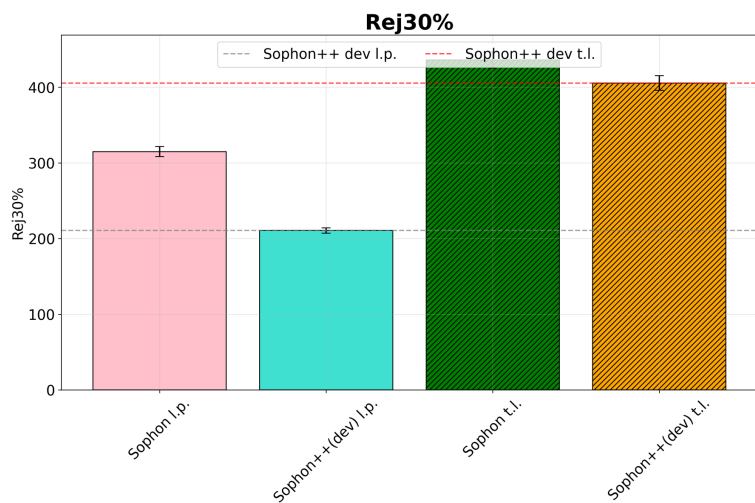
0.6432 vs QCD

0.6432 vs 0.2

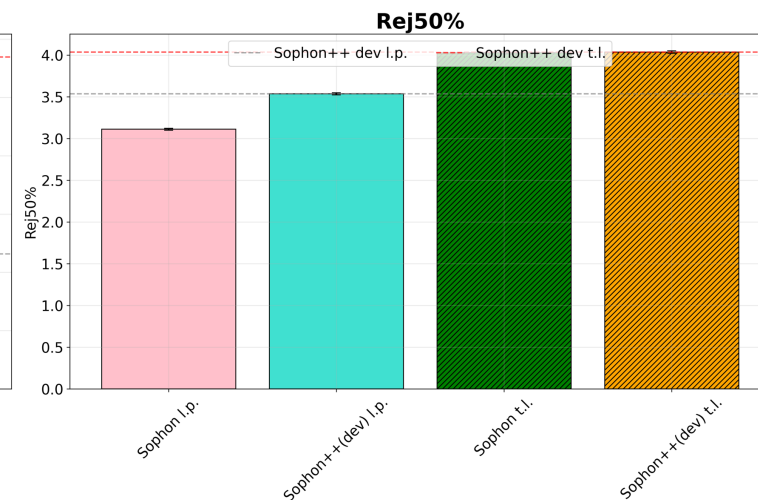
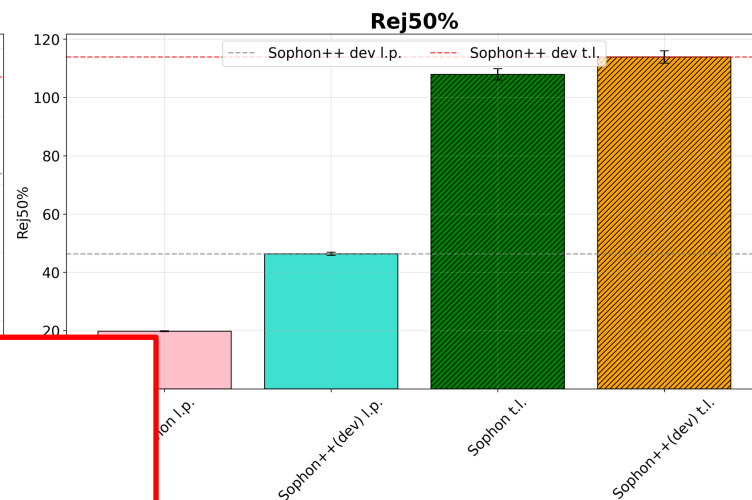
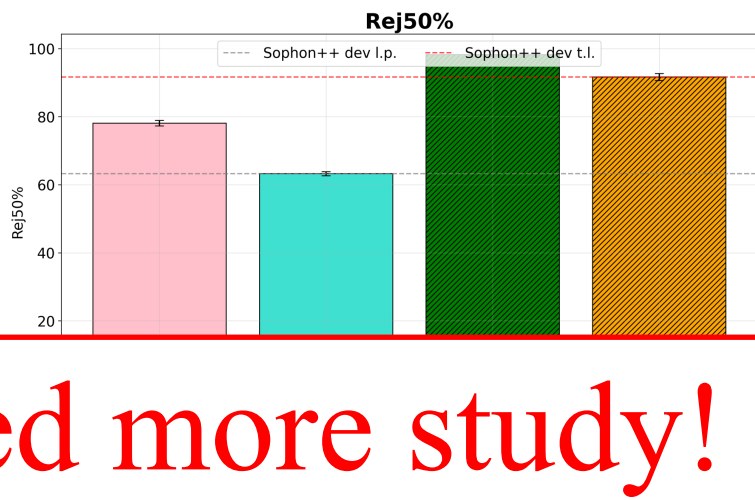
0.8 vs 0.6432

Verification Experiment 3: $X \rightarrow WW(*)$ tagging

Rej30%



Rej50%



need more study!

0.6432 vs QCD

0.6432 vs 0.2

0.8 vs 0.6432

Conclusion

- ***Sophon***, as a large-scale multi-classifier, has achieved optimal performance across various specific tagging tasks.
 - note: It is implemented in CMS as part of GloParT and has been utilized for various boosted jet analyses in Run 3.
- **How can we enhance its generalization capabilities?**
 - we are exploring contrastive "gen-reco" matching.
 - while pure unsupervised models have been developed in recent years to create jet foundation models, we propose insights derived from computer vision history.
 - a “supervised + X approach” can be developed.
 - “X” refers to a methodology akin to OpenAI’s CLIP: it continuously encodes different configurations of gen-level information into *Sophon*’s latent space, thereby achieving stronger generalization.
 - demonstrated superior performance in specific fine-tuning tasks; further optimizations are underway.
- ***Sophon++* provides a promising pathway to upgrade GloParT in CMS!**



Thanks for your attention!

Backup: Core value of Sophon/GloParT at LHC physics

- *Sophon* method implies a pre-training philosophy: “large models for large-scale classification”
 - It has been applied in CMS as GloParT
- A brief summary: GloParT's role for CMS's joint Run2+Run3 analyses in the coming years:
 - **Greater sensitivity improvements for the planned analyses**
 - H/HH/BSM search related to $H \rightarrow b\bar{b}/c\bar{c}/WW^*/\tau\tau\dots$, analyses requiring W/Z/t tagging, ...(several works now using GloParTv1/2/3)
 - **Expanding the landscape of boosted-topology search**
 - Novel final states to be explored! - a reminder to investigate calibration feasibility and discuss with BTV+JME
 - **Creating new paradigms in exploiting the AK8 jet model**
 - Fine-tune GloParT at the analyses level (design custom taggers); support broad MC-free searches
 - Facilitate anomaly detection via GloParT fine-tuning will be a great complement to current BSM programs