# GIVING ML A BOOST TOWARDS RESPECTING (APPROXIMATE) SYMMETRIES

*Bay Area Particle Theory Seminar, Apr 2025*

## INBAR SAVORAY

*UC Berkeley & Lawrence Berkeley National Lab*

WITH: PRADYUN HEBBAR, THANDIKIRE MADULA, VINICIUS MIKUNI, BENJAMIN NACHMAN & NADAV OUTMEZGUINE

# INTRODUCTION

# PARTICLE PHYSICS 404

➤ Despite great theoretical and experimental effort, no evidence of New Physics has been found to date.

➤ Many dedicated searches ruled out a significant portion of the parameter space of theoretically motivated models.

➤ However, there is still much more to explore:

  ➤ New theoretical models.

  ➤ A lot of data.

**Google**

**404.** That's an error.

The requested URL /newphysics was not found on this server. That's all we know.
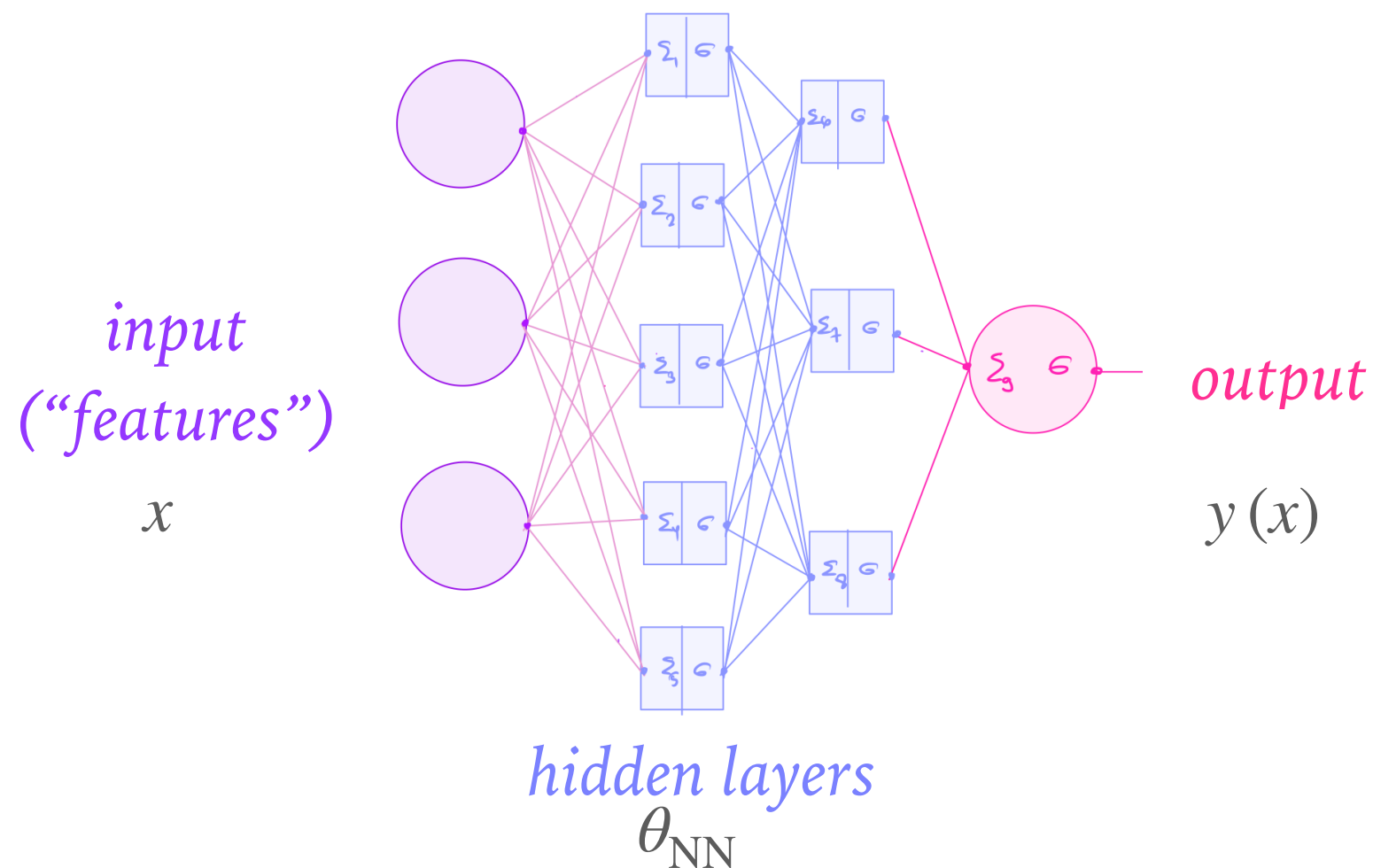
➤ A family of functions - **expressive**, **universal approximators**

    ➤ <u>Architecture</u> - the specific family of functions

Flexible

➤ A family of functions - **expressive**, **universal approximators**

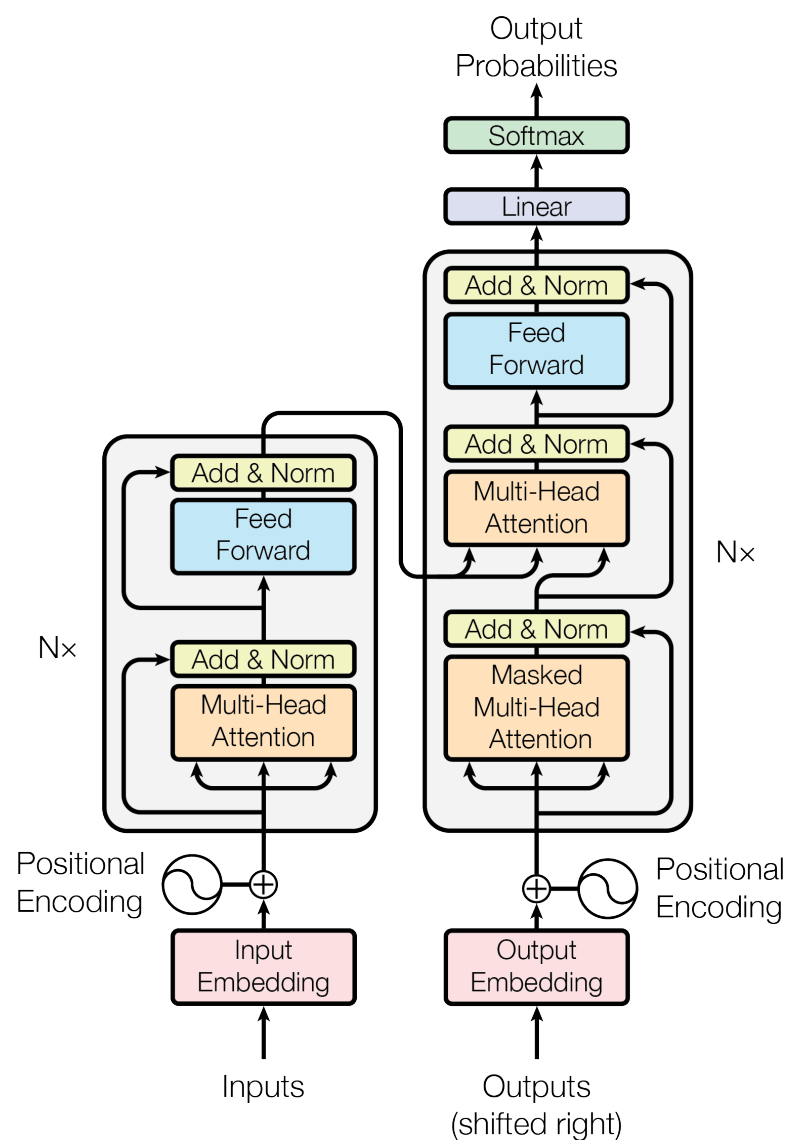➤ <u>Neural Network</u> (NN) - sequence of linear and non-linear functions

*input ("features")* $x$



*hidden layers*
$\theta_{NN}$

*output*

$y(x)$

$$\Sigma_i = \vec{w}_i \cdot \vec{x}_{input} + b_i$$
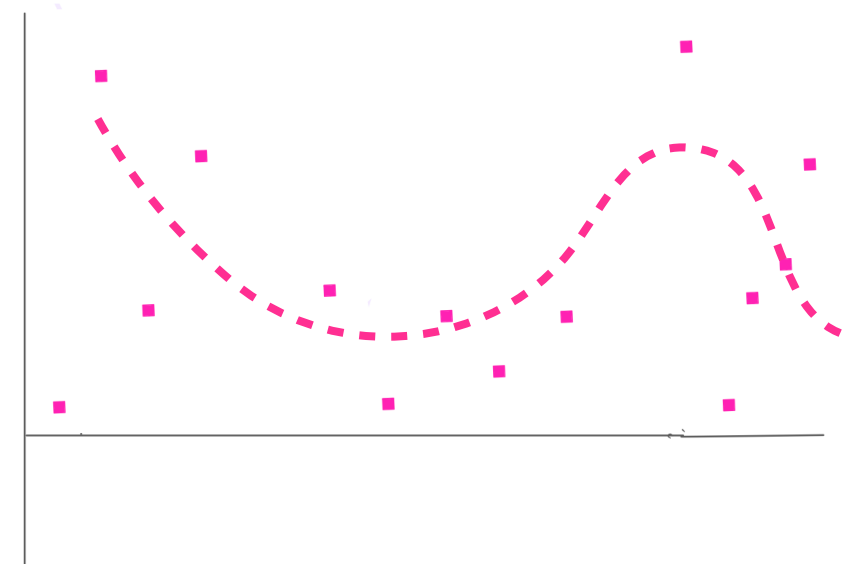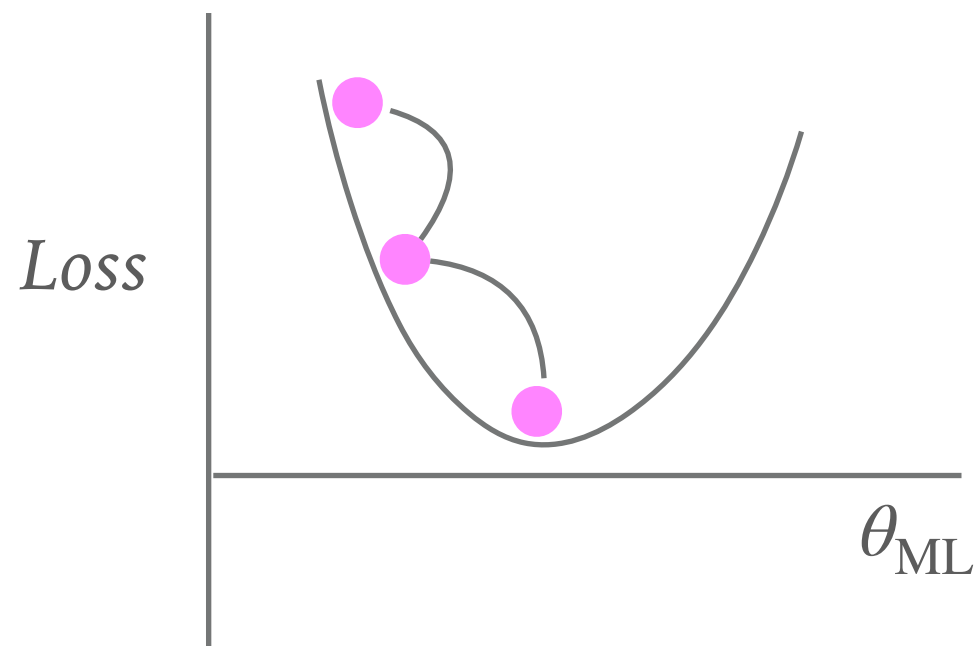
$\sigma$: non-linear activation

Flexible

➤ A family of functions - **expressive**, **universal approximators**

➤ <u>Transformer</u> ~ sequence of NN + attention (non-linear)



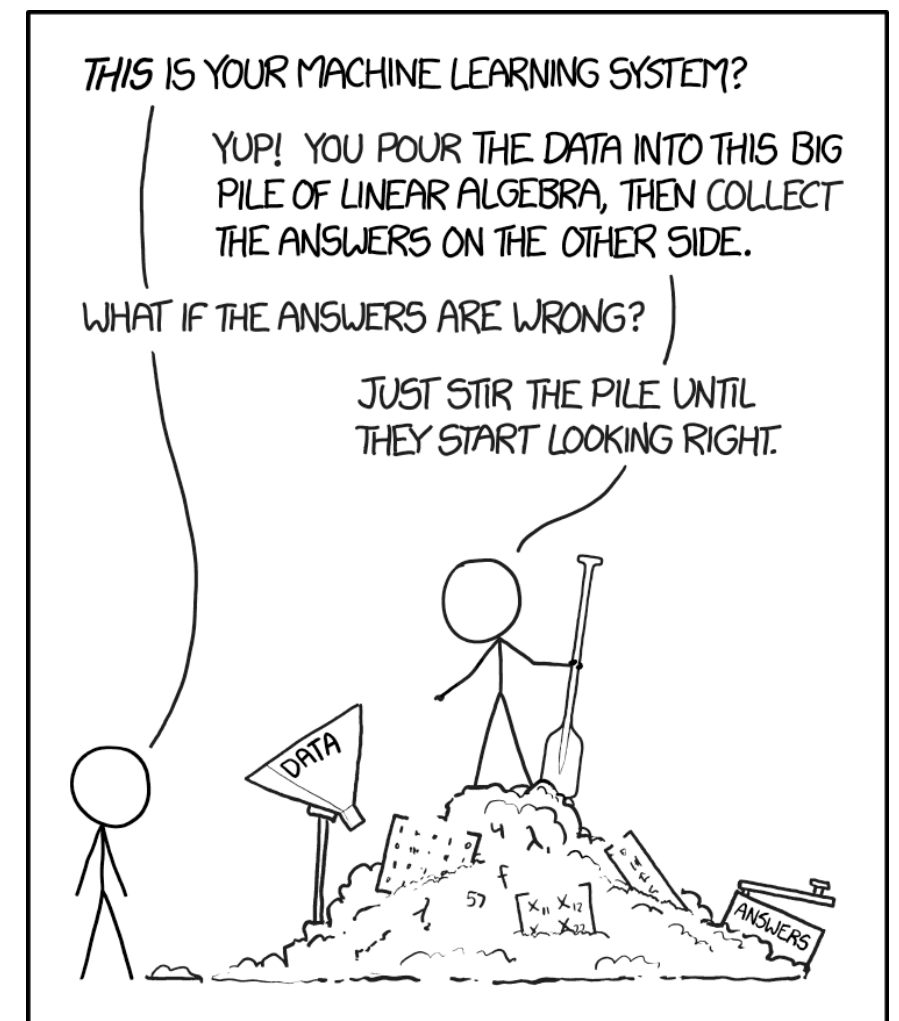*A. Vaswani et al, [1706.03762]*

**Flexible**

➤ A family of functions - **expressive, universal approximators**

➤ <u>Learning</u> - **fit to data.**

  ➤ <u>Training</u> - parameters of function found by minimizing "**loss**" calculated on given dataset.



$Loss$

$\theta_{\mathrm{ML}}$

**Great with a lot of data**

➤ A family of functions - **expressive, universal approximators**

   ➤ <u>Architecture</u> - the specific family of functions (NN, CNN, GNN, transformer, etc.)

➤ <u>**Learning**</u> - **fit to data.**

   ➤ <u>Training</u> - parameters of function found by minimizing "**loss**" calculated on given dataset.
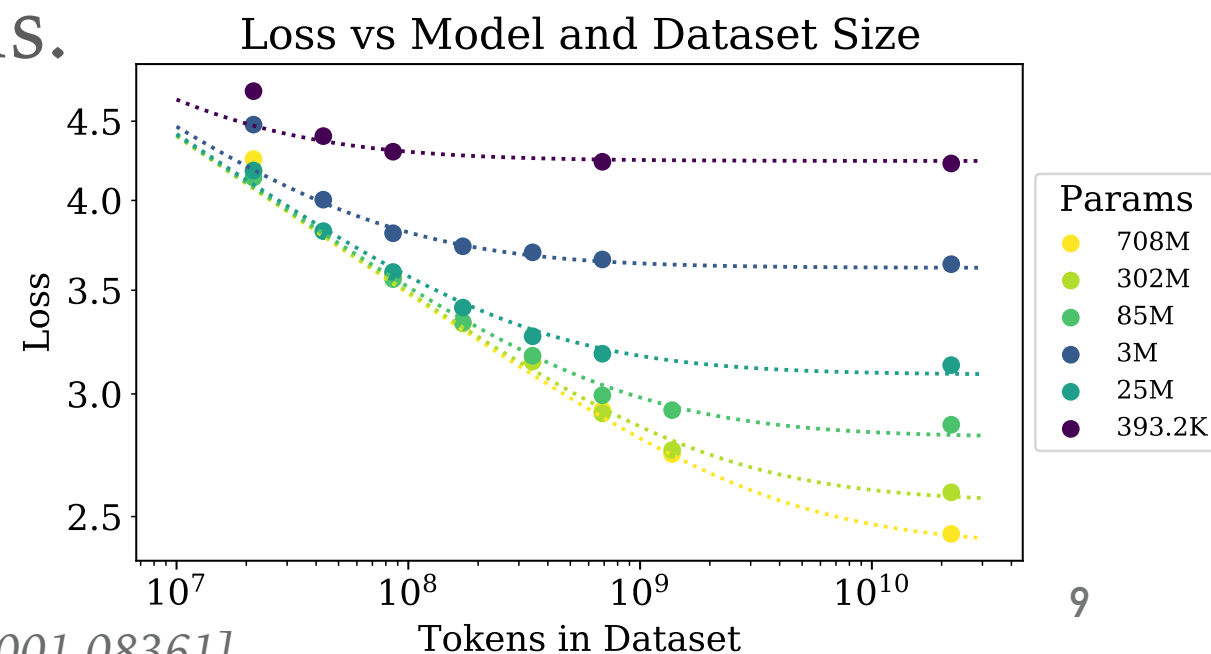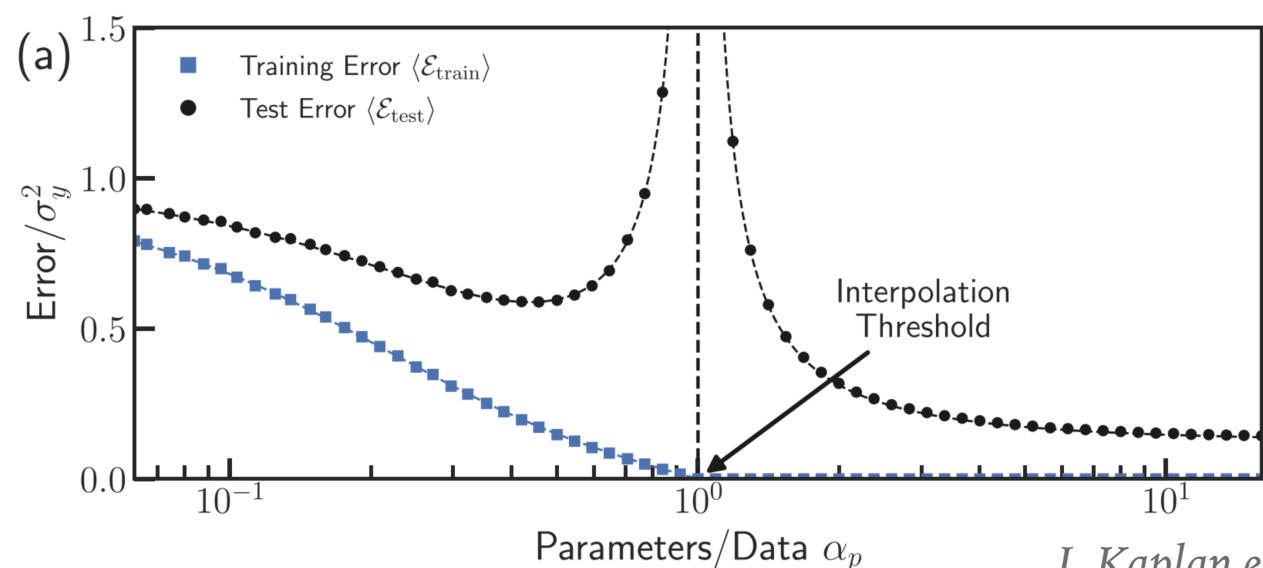
THIS IS YOUR MACHINE LEARNING SYSTEM?

YUP! YOU POUR THE DATA INTO THIS BIG PILE OF LINEAR ALGEBRA, THEN COLLECT THE ANSWERS ON THE OTHER SIDE.

WHAT IF THE ANSWERS ARE WRONG?

JUST STIR THE PILE UNTIL THEY START LOOKING RIGHT.

DATA

ANSWERS

**Great with a lot of data**

**Flexible**

*xkcd.com*

➤ <u>Modern ML is "more is more"</u> -

    ➤ More data, more parameters, more compute.

    ➤ Better capabilities, but also better generalization.

➤ <u>Modern ML is less specialized</u> -

    ➤ Transformers perform well on a wide range of tasks.

    ➤ Shift from carefully designing models for specific tasks to fine tuning foundational models.

*J. Kaplan et al, [2001.08361]*
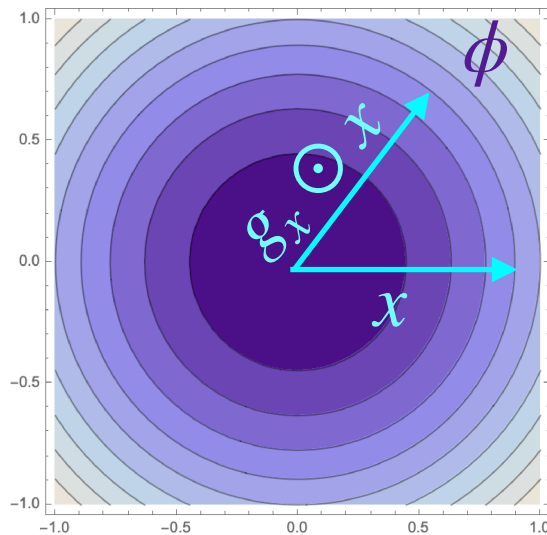
# MORE DATA/PARAMETERS VS. MORE STRUCTURE

➤ On the other hand, more information is also more - especially for **scientific applications**

➤ **Data**

    ➤ <u>Noisy</u> data can "trick" over-parameterized models

    ➤ Might require <u>more precision</u> than language or images

➤ **Theory**

    ➤ Often the <u>underlying truth is "simple"</u> - Ockham's razor

    ➤ We have **<u>guiding theoretical principles</u>** that can be easily phrased as clear mathematical/logical statements

**Physical~structure**

➤ Symmetries as theoretical input - physical information about the system we are trying to describe



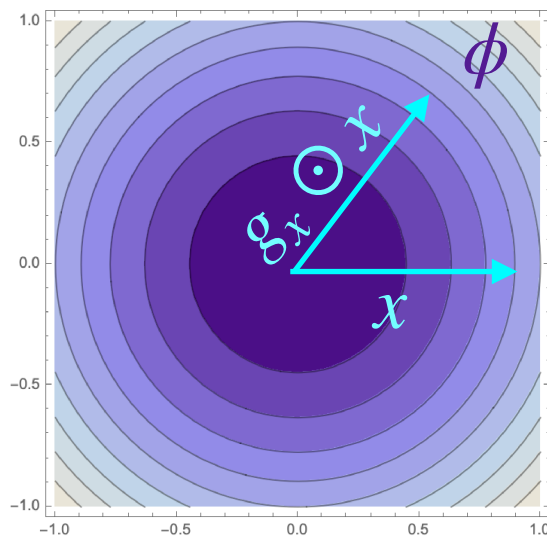$$\phi \left( g_x \odot x \right) = g_\phi \odot \phi \left( x \right) \qquad g \in G_{\text{symm}}$$

*property of input data*          *desired property of output*

➤ Ubiquitous in particle physics -  flavor, P/CP, rotations, translations…

➤ Often approximate - either theoretically or practically
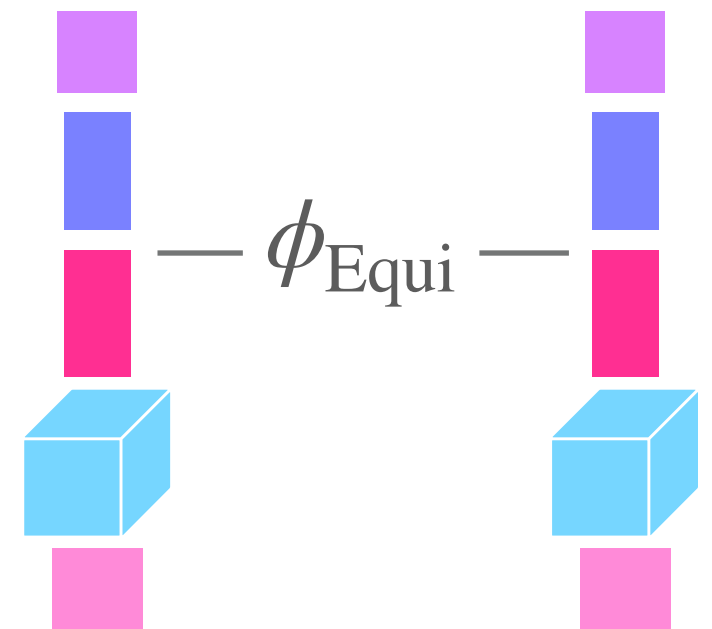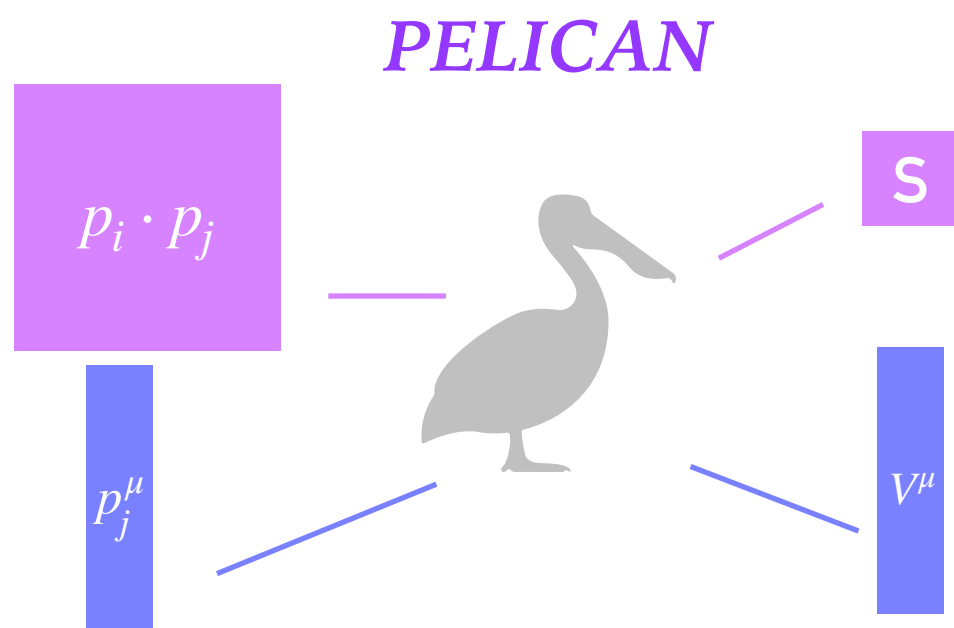
➤ <u>Symmetric architecture</u> - model can only output functions that transform in the correct way by construction.



$$\phi_{\text{ML}}\left(g_x \odot x\right) = g_\phi \odot \phi_{\text{ML}}\left(x\right)$$

*property of input data*  *desired property of output*

$— \phi_{\text{Equi}} —$

➤ Invariant - $g_\phi = 1$

➤ "Equivariant" - $g_\phi = g_x = g$ (covariant)

➤ <u>Symmetric architecture</u> - model can only output functions that transform in the correct way by construction.

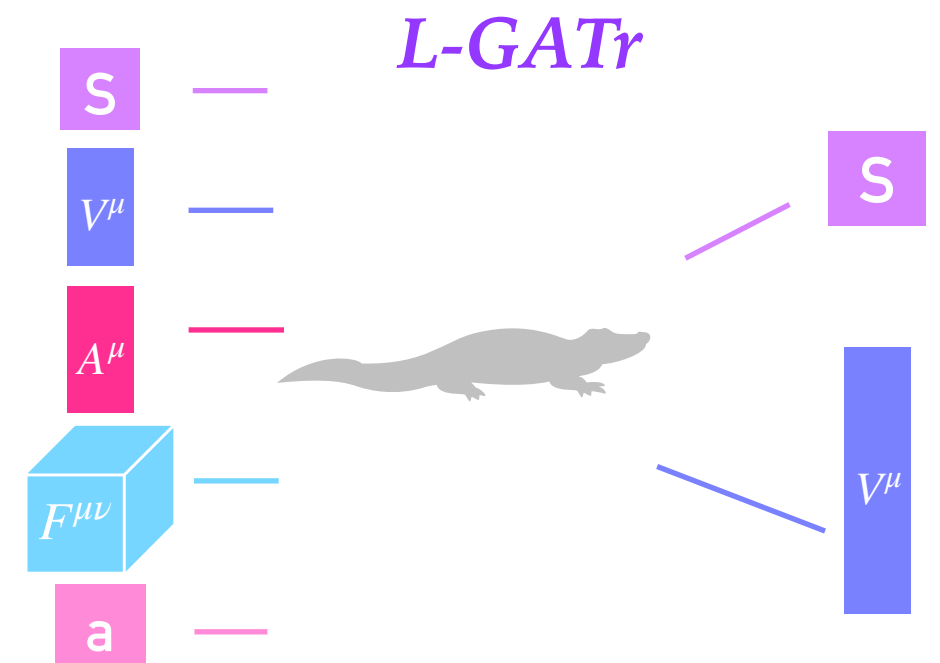$$\phi_{\text{ML}}\left(g_x \odot x\right) = g_\phi \odot \phi_{\text{ML}}\left(x\right)$$

➤ **<u>Lorentz invariance</u>** - theoretically exact, space-time symmetry, continuous and non-compact.     $g = \Lambda\left(\vec{\beta}, \vec{\theta}\right)$

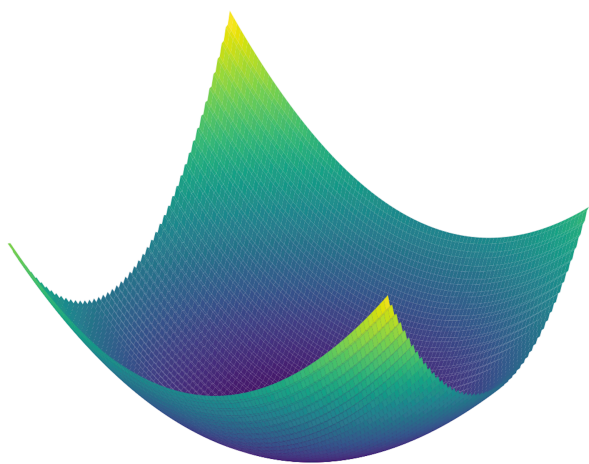➤ Systematically build representations



*PELICAN*

*L-GATr*

*A. Bogatskiy, T. Hoffman, D. W. Miller, J. T. Offermann, X. Liu, [2307.16506]*

*J. Spinner, V. Bresó, P. De Hann, T. Plehn, J. Thaler, J. Brehmer, [2405.14806]*
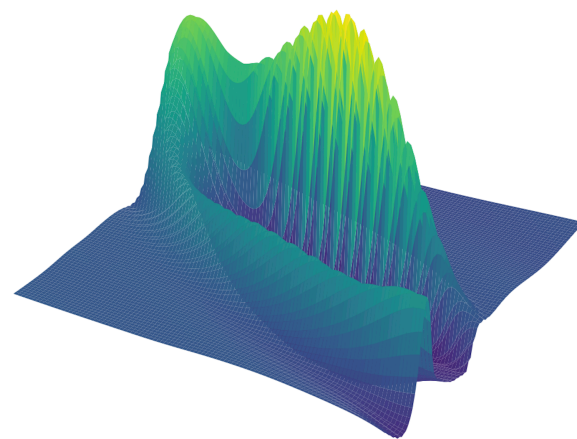
# CHALLENGES OF EQUIVARIANT MODELS

➤ Equivariant models have shown to improve performance on particle physics tasks.

➤ Expressivity could be challenging due to limited "building blocks".

➤ Can be more compute intensive - overhead evaluation time and more FLOPs per parameter.
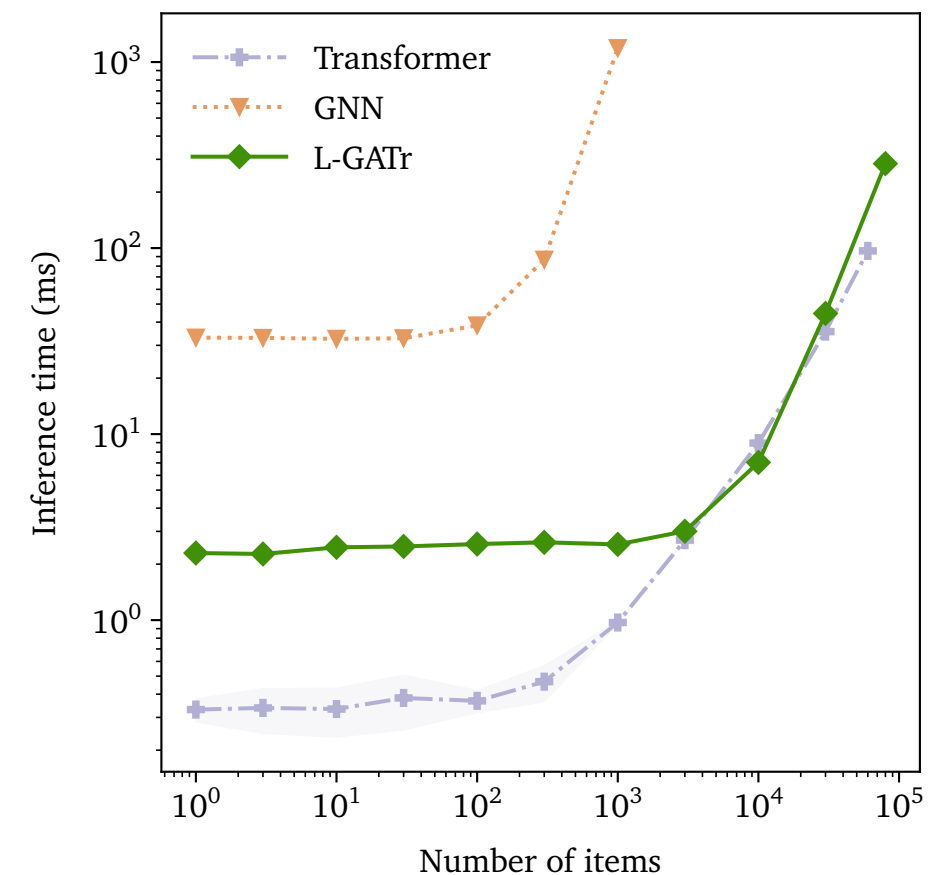
➤ Trainability - less smooth loss surface.



*Transformer*

*GATr*

A. Elhag, T. Rusch, F. Di Giovanni and M. Bronstein, [2410.17878]



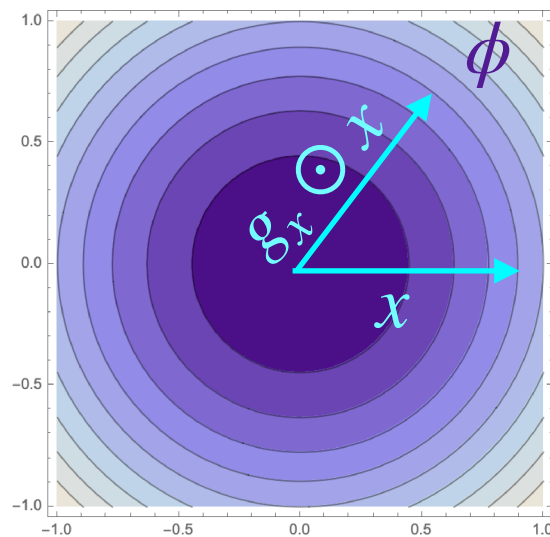J. Spinner, V. Bresó, P. De Hann, T. Plehn, J. Thaler, J. Brehmer, [2405.14806]

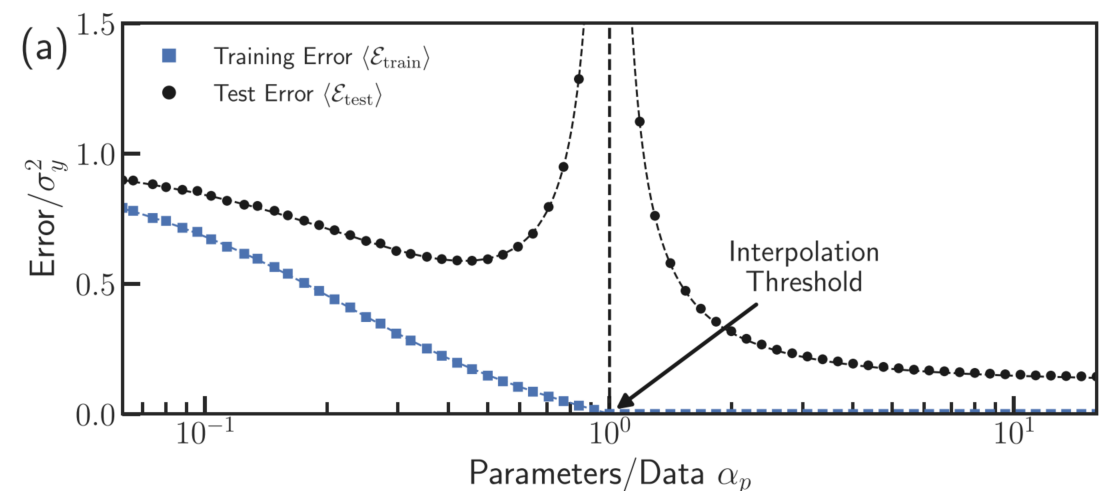# APPROXIMATE SYMMETRIES

➤ Often physical symmetries are only <u>approximate</u>.

➤ Although Lorentz invariance is exact, it is effectively broken if one only transforms the final state momenta

  ➤ <u>Beam</u> - introduces a preferred direction

  ➤ <u>Detector</u> - different energy efficiencies and spatial coverage/ sensitivity.

  ➤ <u>Clustering</u> - algorithm takes into account euclidean distances.

➤ We want flexible and easy to train models, that are aware of symmetries but can choose how to use that information.

➤ Instead of imposing symmetries, specify a <u>preference</u> towards respecting them.





**Physical**   **Scalable**   **Flexible**

# SYMMLOSS

➤ A symmetry-encouraging term added to the loss

$$\mathscr{L} = \mathscr{L}_{\mathrm{task}} + \lambda_{\mathrm{symm}}\mathscr{L}_{\mathrm{symm}}$$

$$\mathscr{L}_{\mathrm{symm}} = \| \phi_{\mathrm{ML}}\left(g_x \odot x\right) - g_\phi \odot \phi_{\mathrm{ML}}\left(x\right) \|^2$$
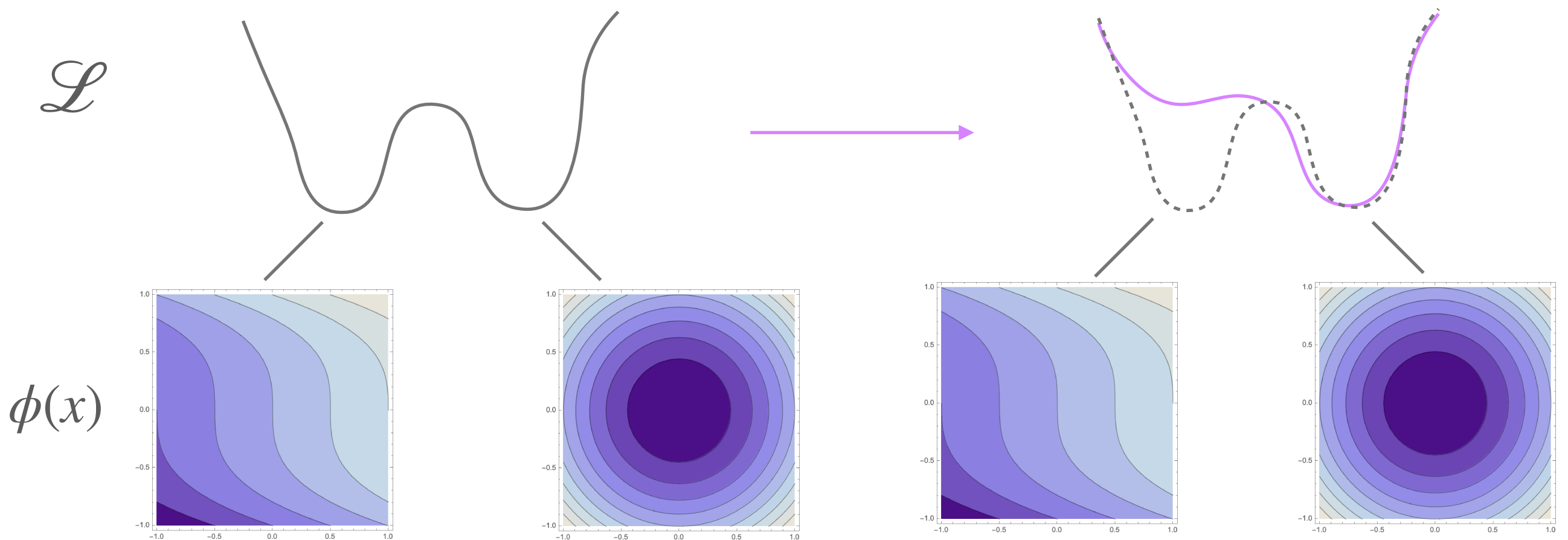
# SYMMLOSS

➤ A symmetry-encouraging term added to the loss

$$\mathscr{L} = \mathscr{L}_{\text{task}} + \lambda_{\text{symm}}\mathscr{L}_{\text{symm}}$$

$$\mathscr{L}_{\text{symm}} = \|\phi_{\text{ML}}\left(g_x \odot x\right) - g_\phi \odot \phi_{\text{ML}}\left(x\right)\|^2$$

➤ **Relax hard constraints** -

  ➤ Allow for <u>approximate symmetries</u> (and even no symmetries at all).

  ➤ Bias is <u>tunable</u> and controllable.

➤ **Flexible** - can be added to any model.

Scalable          Flexible

Physical

➤ A symmetry-encouraging term added to the loss

$$\mathcal{L} = \mathcal{L}_{\text{task}} + \lambda_{\text{symm}}\mathcal{L}_{\text{symm}}$$

$$\mathcal{L}_{\text{symm}} = \|\phi_{\text{ML}}\left(g_x \odot x\right) - g_\phi \odot \phi_{\text{ML}}\left(x\right)\|^2$$

➤ $\mathcal{L}_{\text{symm}} \to 0$ if $\phi$ is in the desired representation for <u>any group element $g$</u> and <u>any input $x$</u>.

➤ In practice:

  ➤ average over data

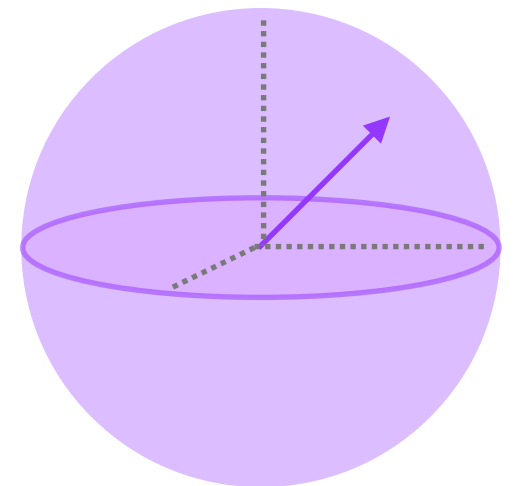  ➤ Group: *Gsymm - group sample*      *δsymm - infinitesimal*

# GSYMM

➤ Measure how different the output is on transformed inputs

*Gsymm:*  $$\mathcal{L}_G = \frac{1}{N} \sum_{i=1}^{N} \left\| \phi_{ML}\left(g_i^x \odot x_i\right) - g_i^\phi \odot \phi_{ML}\left(x_i\right) \right\|^2$$

*Sample $g_i \in G$*

➤ sample from the group.

➤ cheap to calculate.

➤ Measure how different the output is on transformed inputs

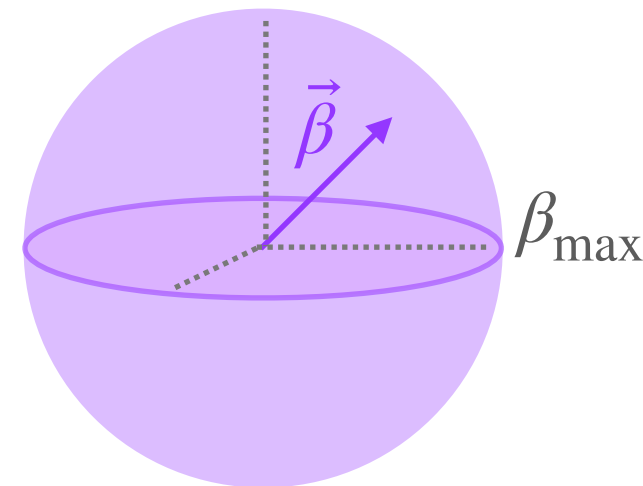$$Gsymm: \quad \mathscr{L}_G = \frac{1}{N} \sum_{i=1}^{N} \left\| \phi_{ML}\left( \Lambda\left(\vec{\beta}_i\right) \odot x_i \right) - g_i^\phi \odot \phi_{ML}\left(x_i\right) \right\|^2$$

*Sample $g_i \in G$*

*scalar:* $g_i^\phi = 1$

*4-vector:* $g_i^\phi = \Lambda\left(\vec{\beta}_i\right)$

➤ sample from the group.

➤ cheap to calculate.

➤ **Lorentz**: boost $\vec{\beta}$ uniformly sampled from a sphere of radius $\beta_{\max}$

# $\delta$SYMM

➤ Infinitesimal transformations by generator $L^a$:

$$\delta^a \phi(x) = \partial_x \phi(x) \, \delta^a \vec{x}$$

$\delta symm:$ $\qquad \mathscr{L}_\delta = \left\| \overset{\text{features}}{\underset{j=1}{\sum}} \left( \underset{\delta^a \phi}{L^a_\phi(\phi)} - L^a_x x_j \cdot \partial_{x_j} \phi \right) \right\|^2_{\text{gens, data}}$

$\qquad \qquad \qquad \qquad \qquad \qquad \delta^a \vec{x} = L^a_x \vec{x}$

$\qquad \qquad \qquad \qquad \qquad \qquad \qquad \quad L_x \text{ in the rep. of } x$

# $\delta$SYMM

➤ Infinitesimal transformations by generator $L^a$:

$$\delta^a \phi(x) = \partial_x \phi(x) \, \delta^a \vec{x}$$

$\delta symm:$

$$\mathcal{L}_\delta = \left\| \sum_{j=1}^{\text{features}} \left( L_\phi^a(\phi) - L_x^a x_j \cdot \partial_{x_j}\phi \right) \right\|_{\text{gens, data}}^2$$
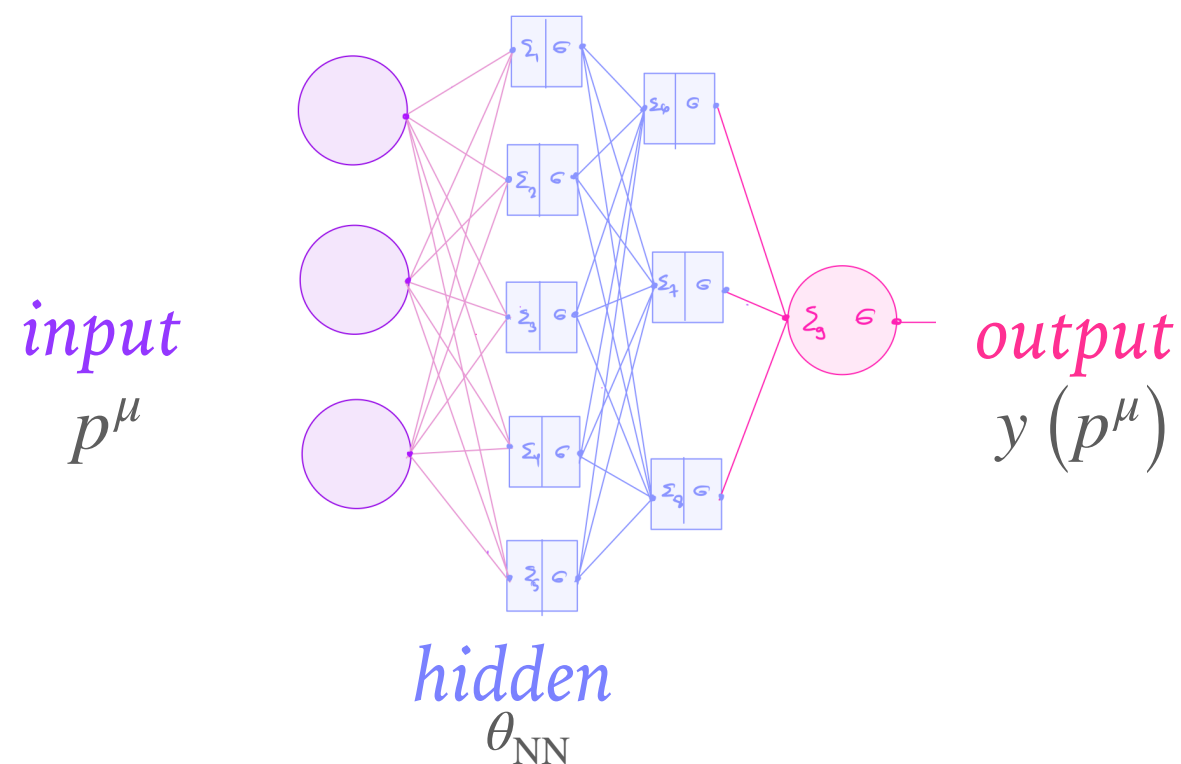
*$L_x$ in the rep. of x*

➤ Is already approximate.

➤ No need to figure out sampling over group.

➤ On the other hand - computationally more expensive.

➤ Infinitesimal transformations by generator $L^a$:

$$\delta^a \phi(x) = \partial_x \phi(x) \, \delta^a \vec{x}$$

$\delta$*symm:* $\quad \mathscr{L}_\delta = \left\| \sum_{j=1}^{\text{features}} \left( L_\phi^a(\phi) - L_x^a x_j \cdot \partial_{x_j}\phi \right) \right\|_{\text{gens, data}}^2$

*scalar: 0*

*4-vector: $L\phi$*

*$L_x$ in the rep. of x*

➤ Is already approximate.

➤ No need to figure out sampling over group.

➤ On the other hand - computationally more expensive.

➤ **Lorentz**: 6 generators: $K_x$, $K_y$, $K_z$, $L_x$, $L_y$, $L_z$

# EXPERIMENTS & RESULTS

➤ Input: list of 4-momenta $p_i^\mu$

➤ NN with 3 hidden layers of width 300, GeLU activation.

➤ **Exact Symmetry**: $f_{\text{truth}}\left(p_i^\mu\right) = \text{poly}\left(\text{p}_i \cdot \text{p}_j\right)$

➤ MSE loss - $\mathscr{L}_{\text{MSE}} = \sum_i \left| f_{\text{truth}}\left(p_i^\mu\right) - y_{\text{pred}}\left(p_i^\mu\right) \right|^2 + \lambda_{\text{symm}} \mathscr{L}_{\text{symm}}$



*input*
$p^\mu$

*hidden*
$\theta_{\text{NN}}$

*output*
$y\left(p^\mu\right)$

# TOYS – EXACT SYMMETRY

$$\mathscr{L} = \mathscr{L}_{\text{MSE}} + \lambda_{\text{symm}}\mathscr{L}_G \qquad \mathscr{L}_G = \left\| \phi_{ML}\left(B_i\left(x_i\right)\right) - \phi_{ML}\left(x_i\right) \right\|^2$$

➤ Gsymm can achieve better performance than baseline on boosted inputs.

➤ Larger training $\beta_{\text{max}}$ - flatter as function of boost, but can under-perform for small boosts.

$$\mathcal{L}_\delta = \left\| \frac{\partial \phi_{ML}}{\partial p_\mu} \cdot \left( L_{\mu\nu} p_i^\nu \right) \right\|^2 \qquad \mathcal{L}_G = \left\| \phi_{ML} \left( B_i \left( x_i \right) \right) - \phi_{ML} \left( x_i \right) \right\|^2$$

➤ Even infinitesimal loss achieves better performance than baseline, and can extend to non-infinitesimal boosts!

➤ $\delta$symm better at smaller $\beta$.

➤ Big $\lambda$ doesn't hurt for small transformations.

$$\mathcal{L} = \mathcal{L}_{\mathrm{MSE}} + \lambda_{\mathrm{symm}} \mathcal{L}_{\mathrm{symm}}$$

➤ Input: list of 4-momenta $p_i^\mu$

➤ NN with 3 hidden layers of width 300, GeLU activation.

➤ **Exact Symmetry**: $f_{\text{truth}}\left(p_i^\mu\right) = \text{poly}\left(p_i \cdot p_j\right)$

➤ **Approximate Symmetry**: $f_{\text{truth}}\left(p_i^\mu\right) = \text{poly}\left(p_i \cdot p_j, p_i \cdot s\right)$

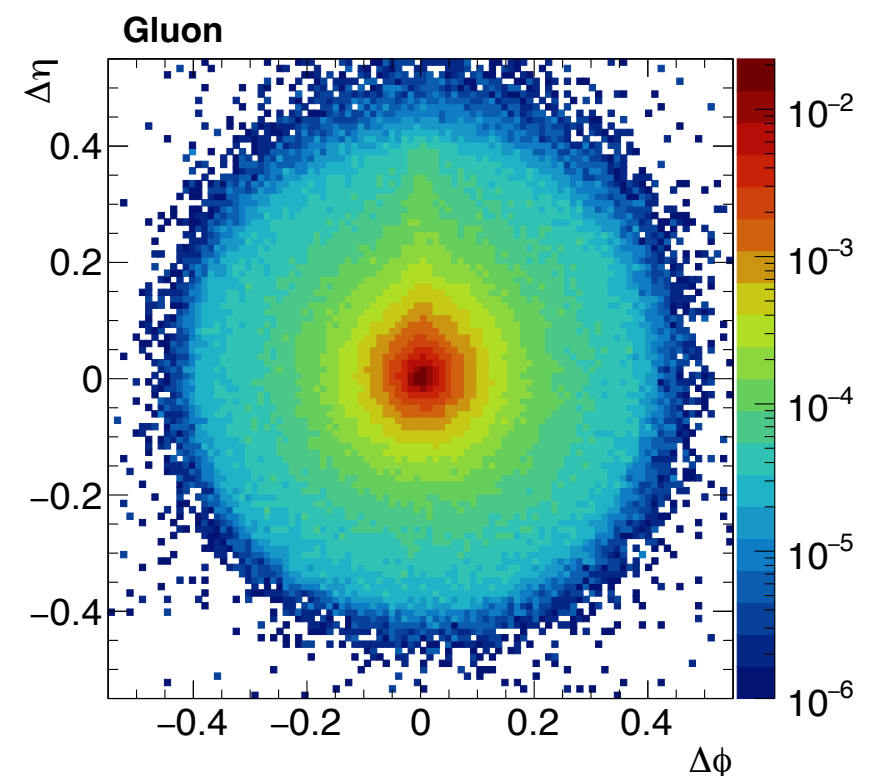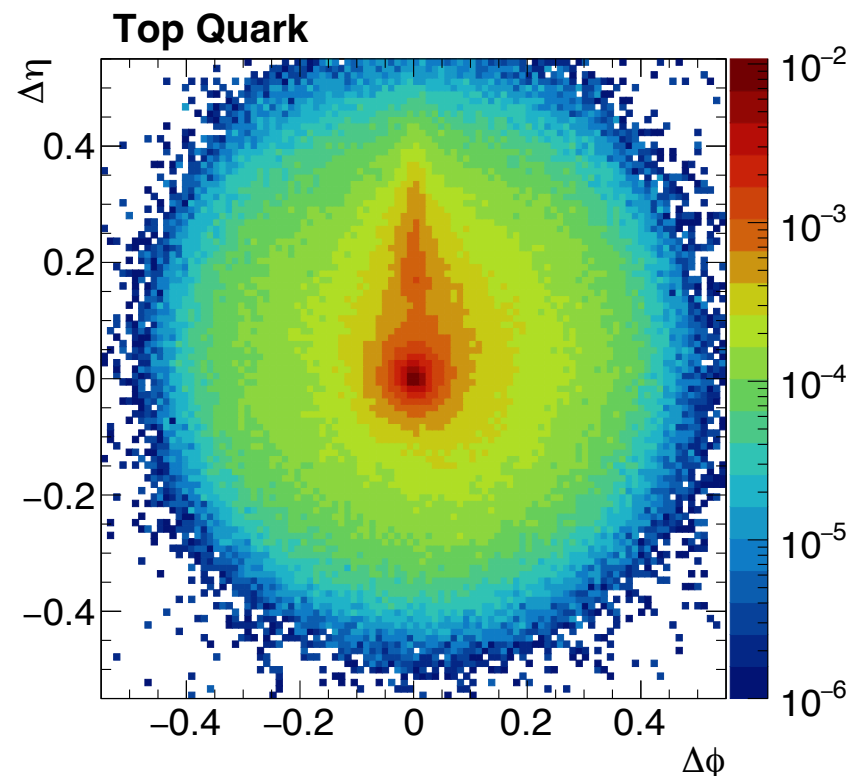  ➤ "Spurion" $s = \begin{pmatrix} 0 & 0 & 0 & 10^{-3} \end{pmatrix}$

$$\mathscr{L} = \mathscr{L}_{\text{MSE}} + \lambda_{\text{symm}}\mathscr{L}_{\text{symm}}$$

➤ Gain even when the symmetry is not exact.

$$\mathscr{L}_{\delta} = \left\| \frac{\partial \phi_{ML}}{\partial p_{\mu}} \cdot \left( L_{\mu\nu} p_i^{\nu} \right) \right\|^2 \qquad \mathscr{L}_G = \left\| \phi_{ML}\left( B_i\left( x_i \right) \right) - \phi_{ML}\left( x_i \right) \right\|^2$$

- ➤ Physics example - QCD vs. top-jets

  - ➤ Precision measurements

  - ➤ BSM studies

- ➤ Goal - learn $p(x|\text{top})$ vs. $p(x|\text{QCD}) \rightarrow$ classify jet.



*V. Mikuni, F. Canelli, [2102.05073]*

➤ ATLAS top tagging dataset

*ATLAS collaboration (2022), https://opendata.cern.ch/record/15013*

➤ Most realistic dataset

    ➤ Full LHC Run-2 conditions (including pile-up)

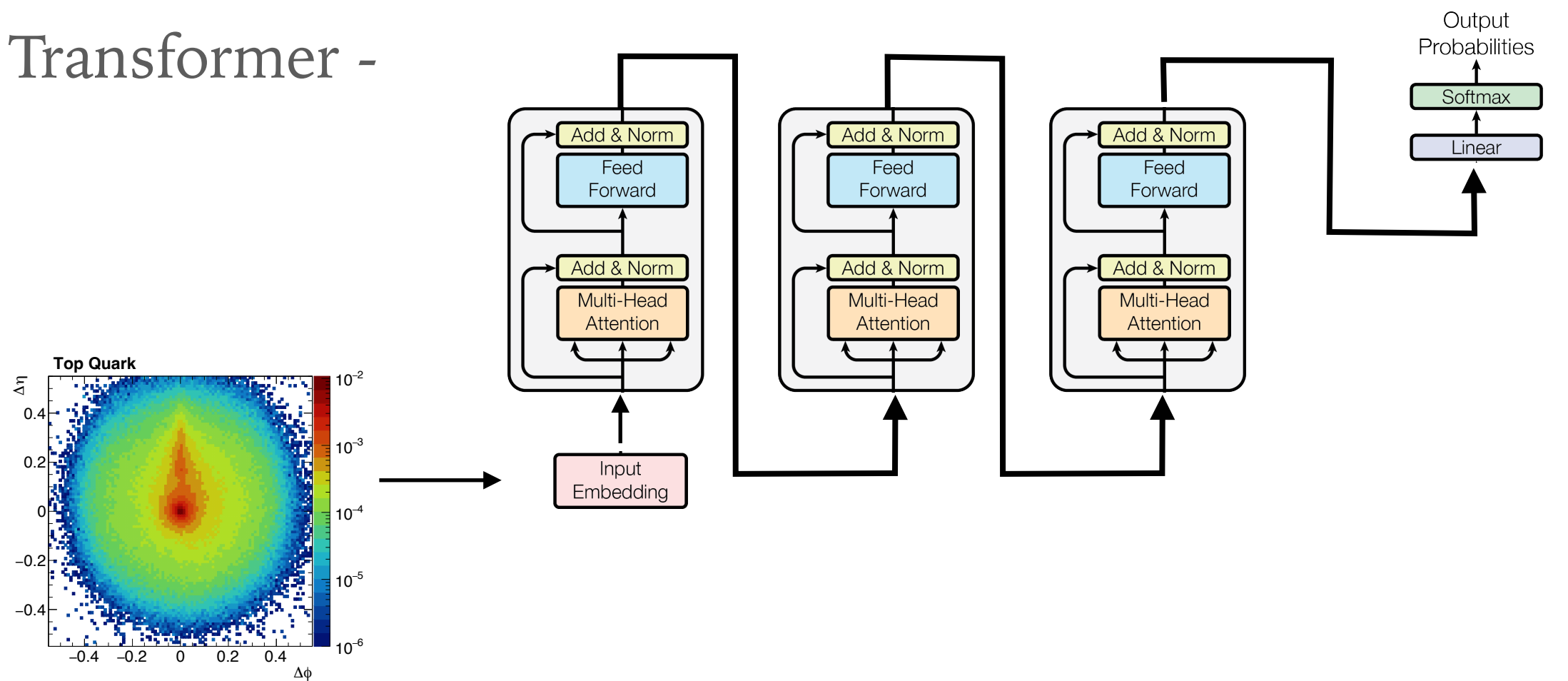    ➤ Full detector simulation

    ➤ Event reconstruction

➤ Input - jet constituents 4-momenta

$$\left\{ p_T^i, E^i, \frac{p_T^i}{p_T^{\text{jet}}}, \frac{E^i}{E^{\text{jet}}}, \Delta\phi^i, \Delta\eta^i, \Delta R^i = \sqrt{\Delta\eta^2 + \Delta\phi^2} \right\}$$
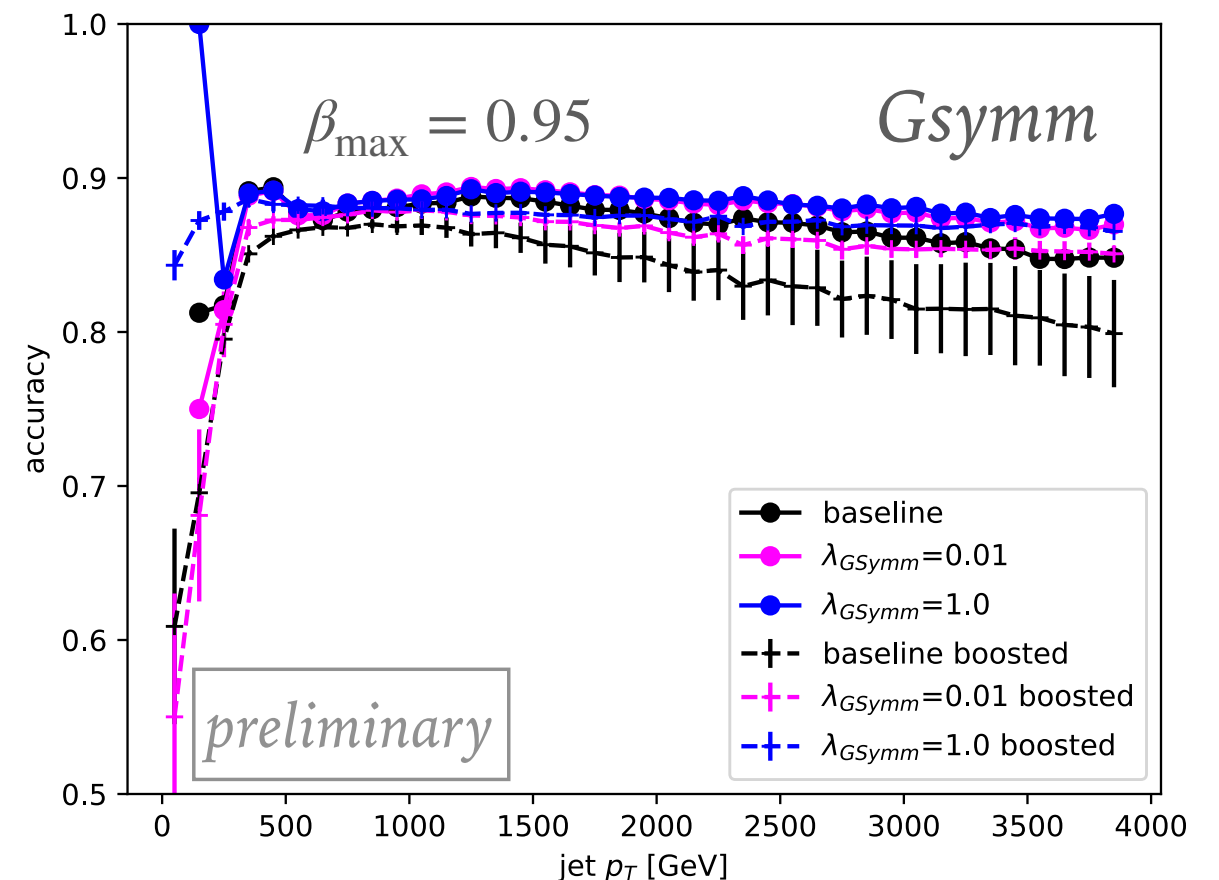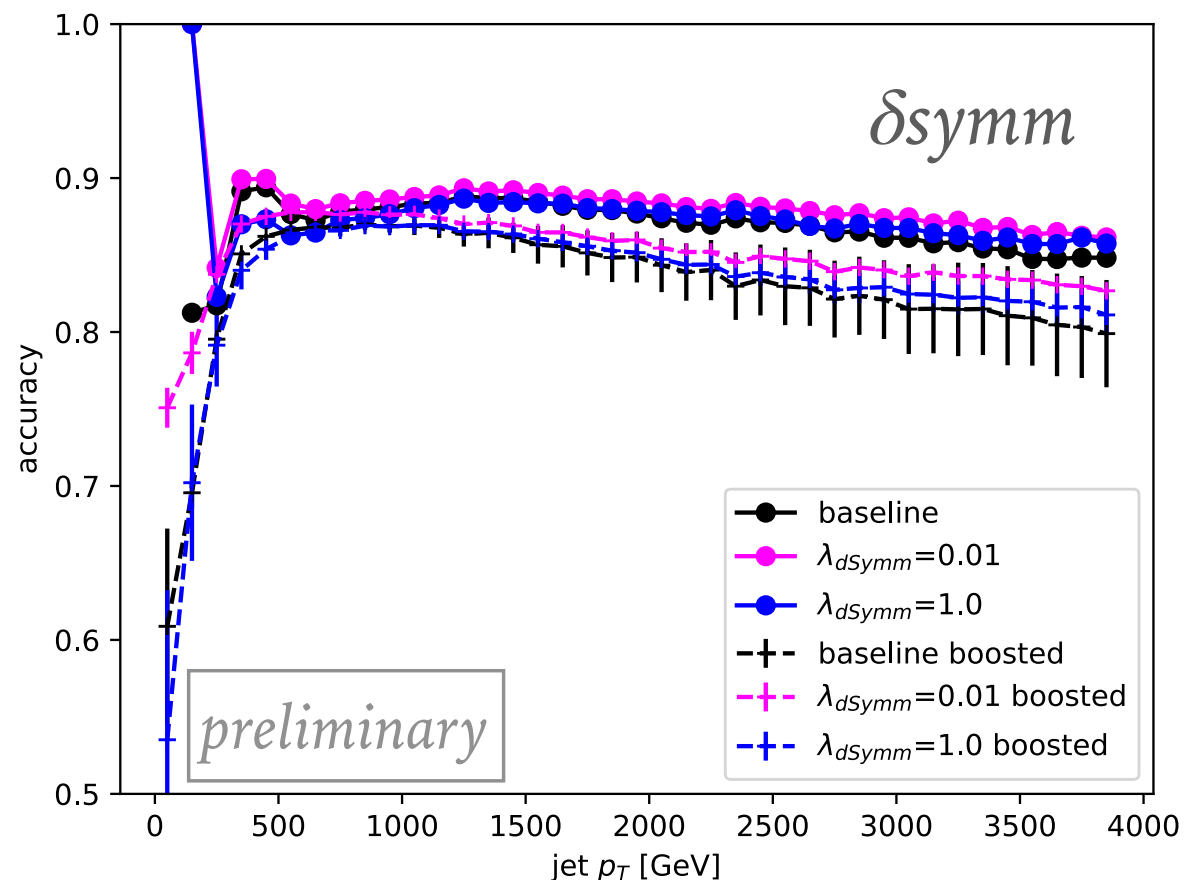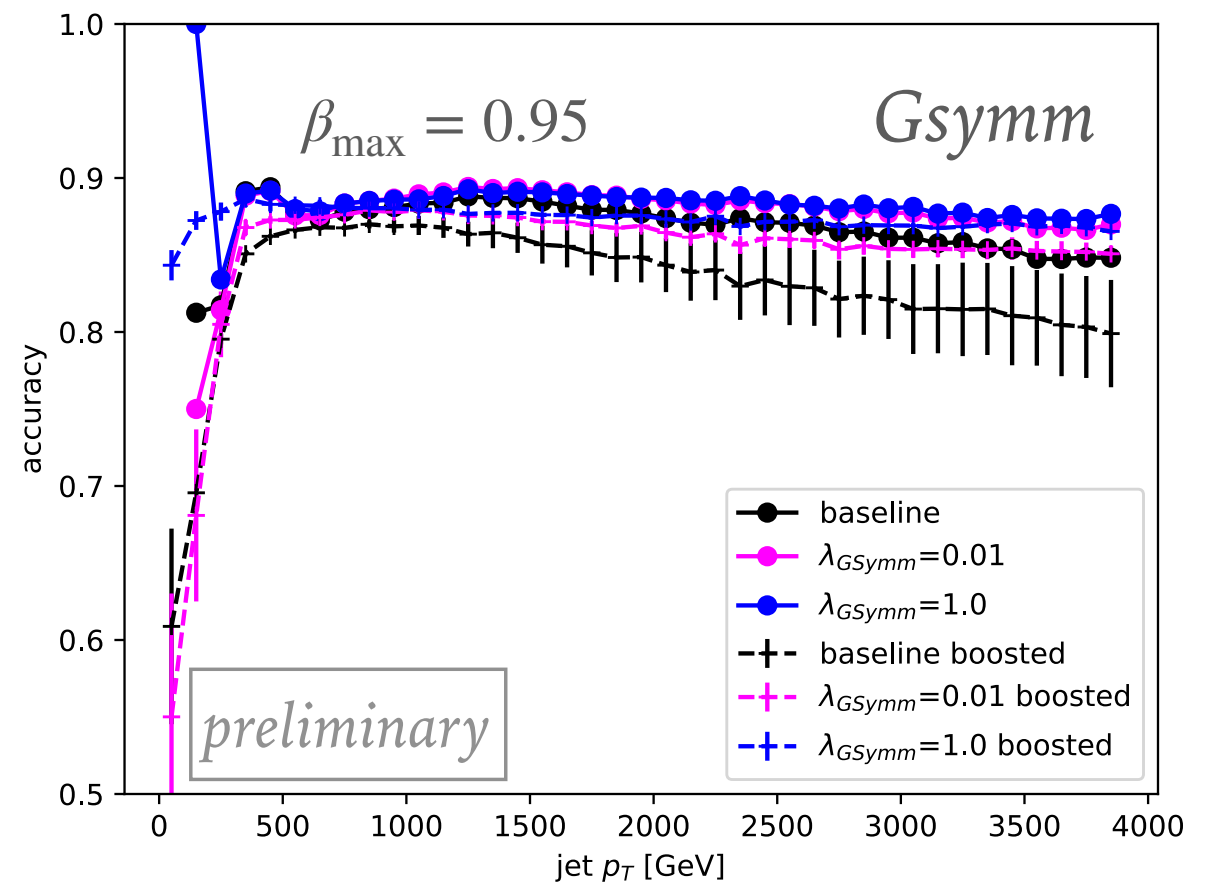


$$p\left(x | \text{top}\right)$$

➤ Transformer -
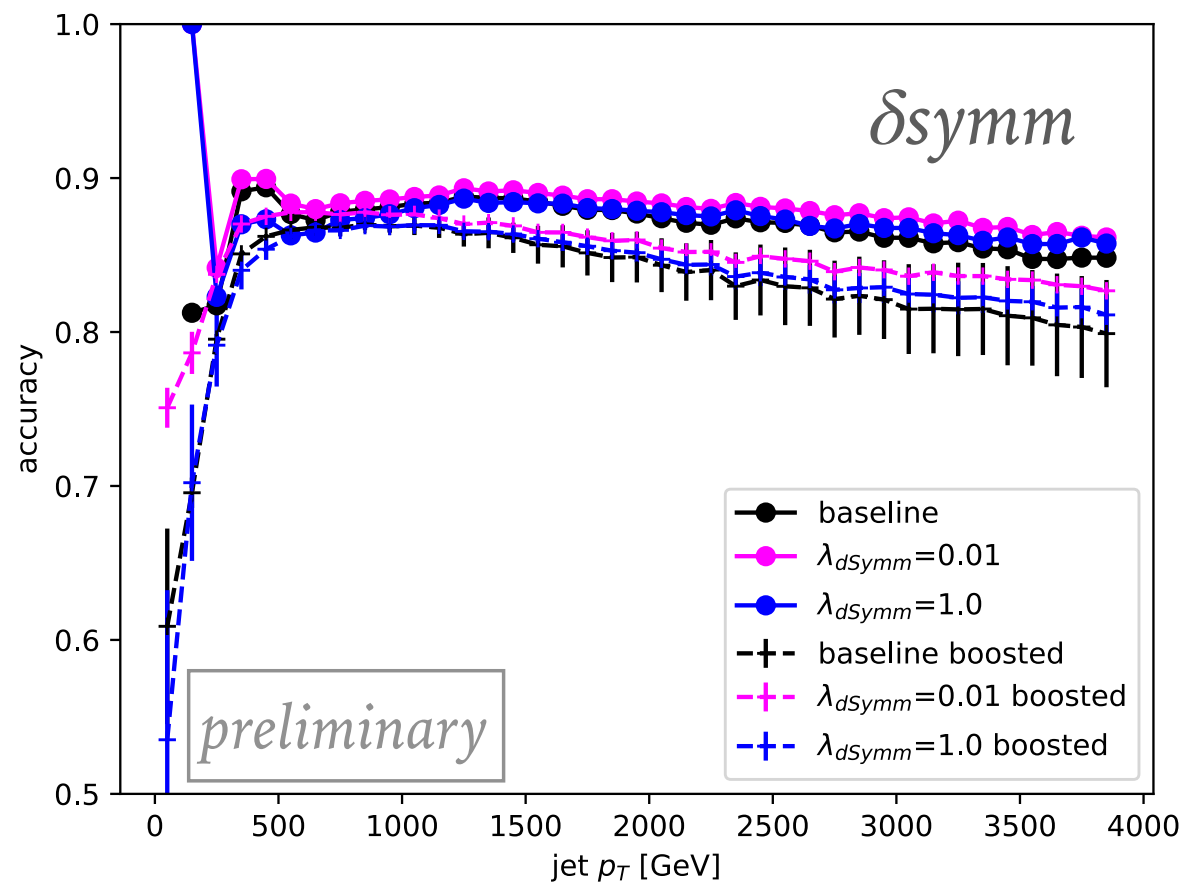
➤ Invariance check - boosting inputs and assigning them with the same truth values as original inputs.

➤ As expected - improved and flatter performance over boosted inputs.
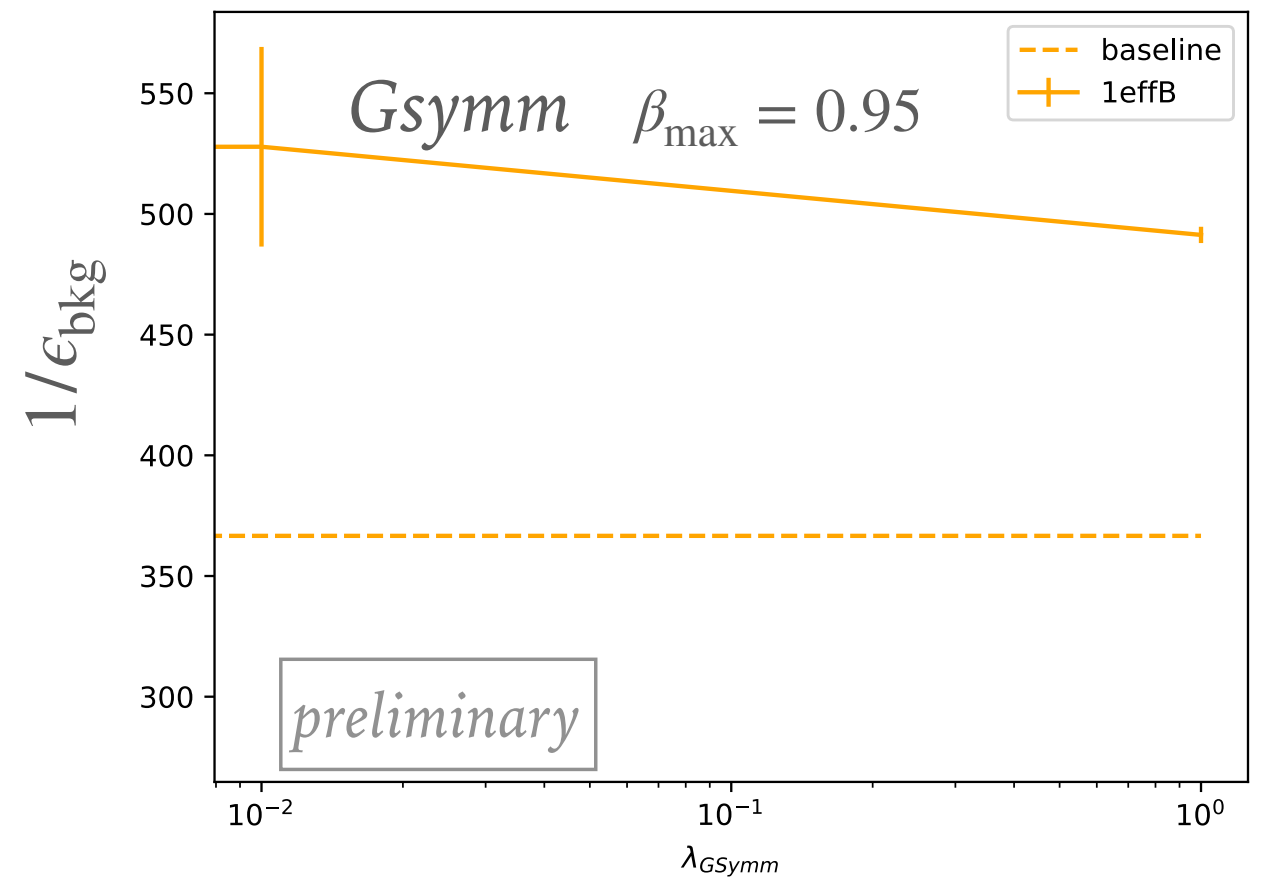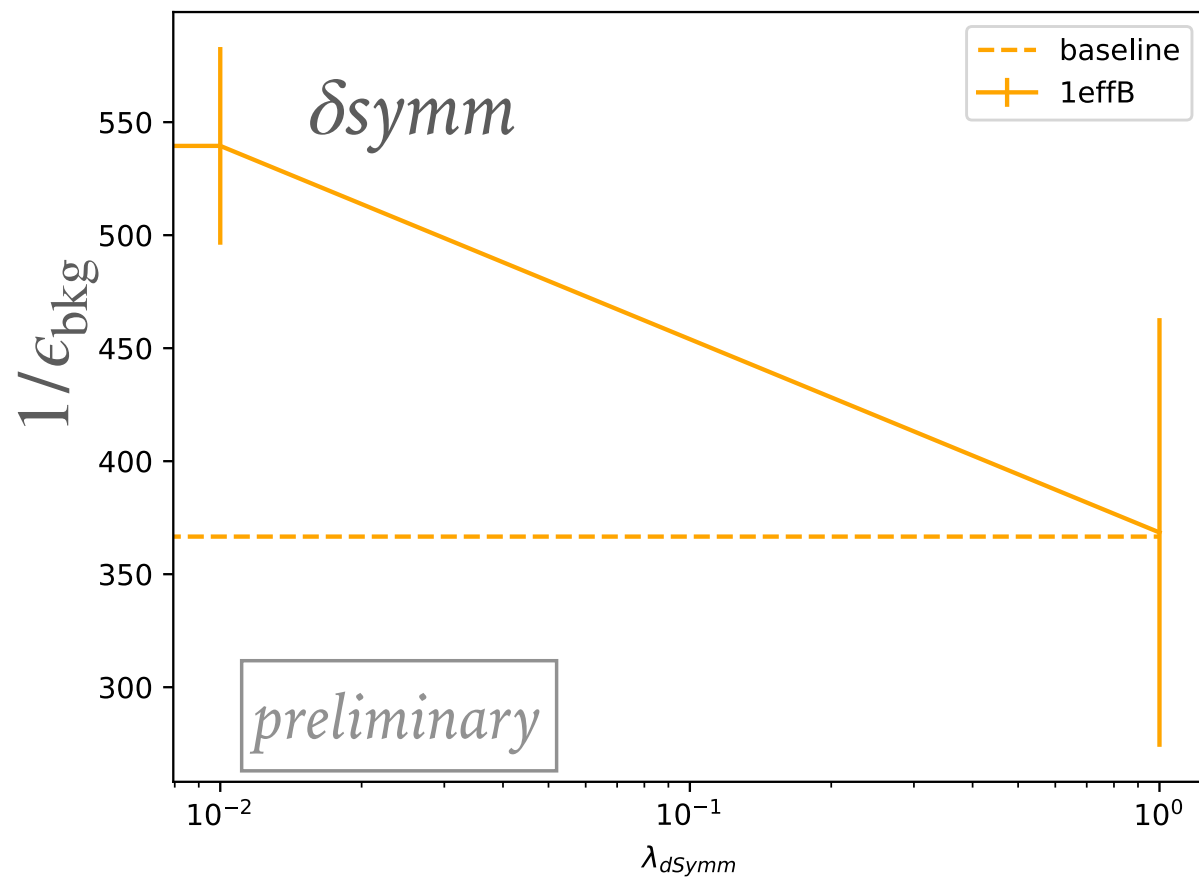
➤ For the real data - performance at least on-par with baseline.
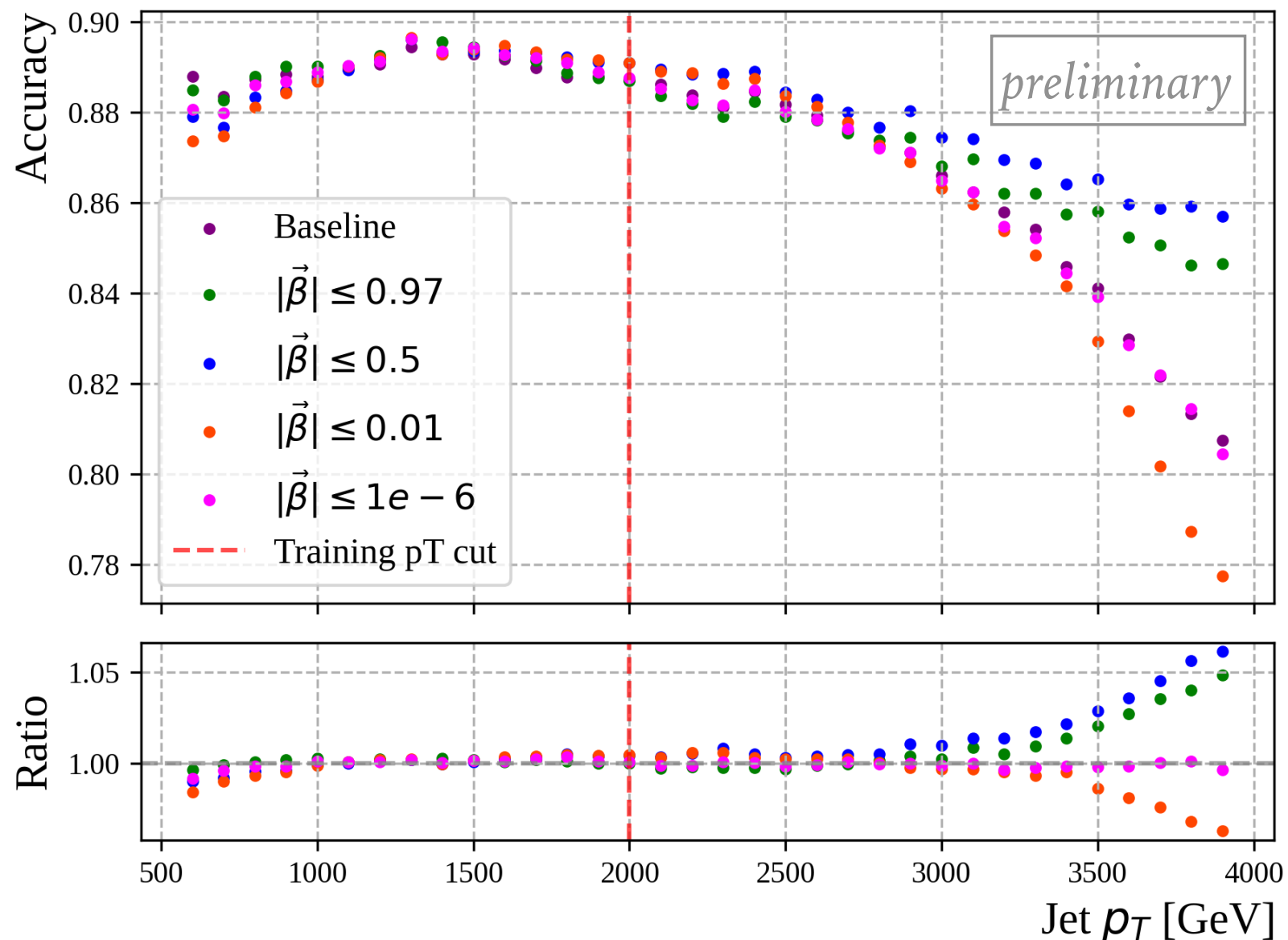
➤ For the real data - performance at least on-par with baseline.

➤ Improved background rejection at signal efficiency of 0.3.

➤ Extrapolation test: train only on $p_T \leq 2$ TeV
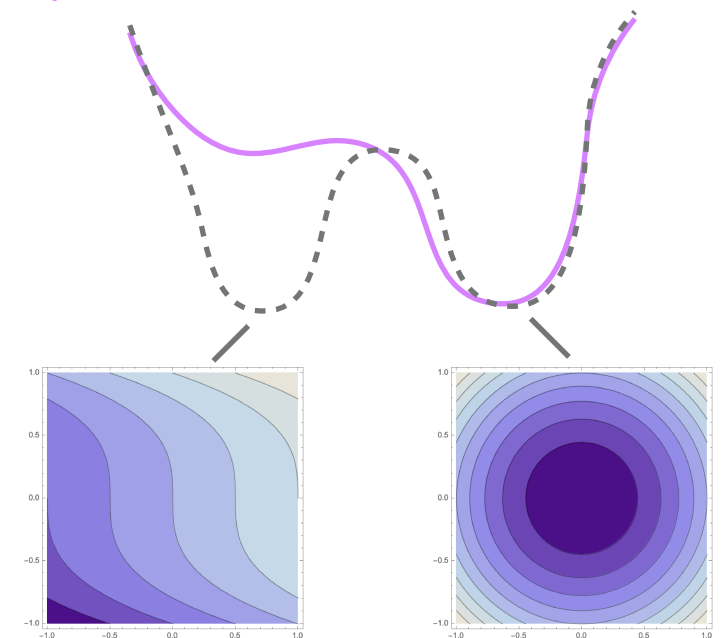
➤ Gsymm with $\beta_{\max} = 0.5$ extrapolates best to unseen $p_T$ !

➤ ML + physical knowledge help extract more information from data.

➤ Symmetries can be imposed on architecture-level, but can be challenging to build and train.

➤ **SymmLoss - bias the model towards respecting symmetries.**

$$\mathscr{L} = \mathscr{L}_{\text{task}} + \lambda_{\text{symm}}\mathscr{L}_{\text{symm}}$$

➤ <u>Flexible</u> - can be added to any model, easy to implement.

➤ <u>Multi-purpose</u> - accommodates <u>approximate symmetries</u> (and no symmetries).

➤ Bias is <u>tunable</u> and controllable.

➤ **Better results for symmetric problems, even if the symmetry is broken.**

# FUTURE WORK

➤ Working on full comparison to PELICAN & L-GATr

➤ Performance -

  ➤ combine with SOTA models

  ➤ Other ideas for broken symmetry losses

➤ Scaling behavior

# THANK YOU!