

Scaling laws for amplitude surrogates

Joaquín Iturriza Ramirez

In collaboration with Anja Butter, Bertrand Laforge, Víctor Bresó Pla and Henning Bahl

IRN Terascale @ Montpellier 2025



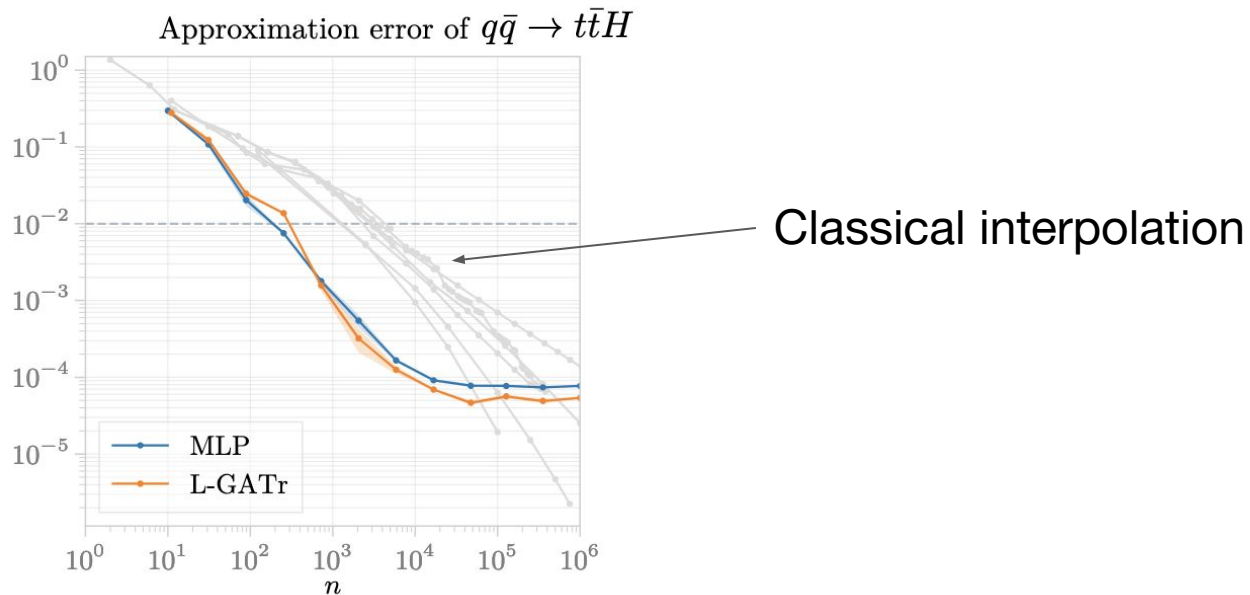
Co-funded by
the European Union

Motivation

Scattering amplitudes are at the heart of MC event generators

Calculating higher order scattering amplitudes is **expensive**

ML surrogate models are excellent interpolates



How accurate can we become and what is the price?

Motivation

Similar scaling behaviours have been observed in many applications of deep learning

Large Language Modeling [1]

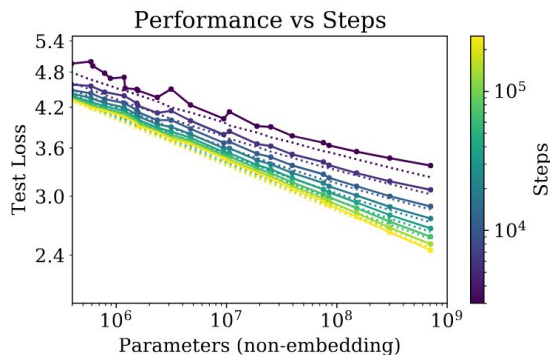
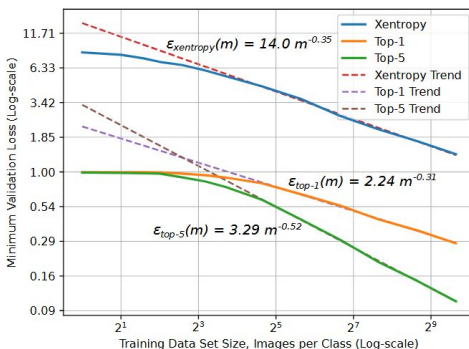
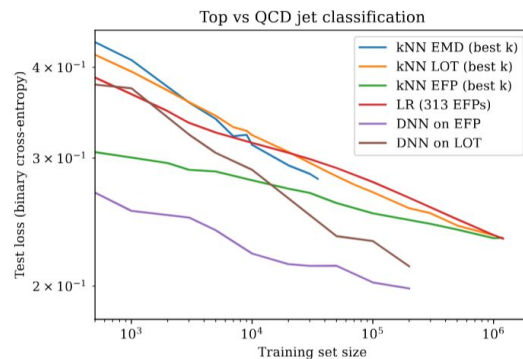


Image Classification [2]



Jet Classification [3]



Performance improves predictably as a power law with **training dataset size D** , **computing resources C** and **number of parameters N**

$$L(X, Y, Z) = (X_c/X)^{\alpha_X} + K(Y, Z)$$

Exponent of power law could be related to **intrinsic dimensionality** [4]

Can we determine scaling laws for scattering amplitude surrogates?

Compare

- ▶ Many different processes:

$$q\bar{q} \rightarrow t\bar{t}H, q\bar{q} \rightarrow Z + ng, q\bar{q} \rightarrow WZ + ng, q\bar{q} \rightarrow WWZ, gg \rightarrow \gamma\gamma + ng$$

- ▶ Scalings in D , C and N
- ▶ Different loss functions: **MSE** vs **Heteroscedastic Loss** for uncertainties estimation
- ▶ Different architectures: **MLP** vs **LloCa-Transformer**

If scaling laws are universal \longrightarrow We predict desired accuracy for given resources

We know the **degrees of freedom**: Test relation between scaling and intrinsic dimensionality

4-momentum of
particles in the process



Amplitude

All the data is generated with MadGraph, all amplitudes calculated at **lowest order**

Main case study $q\bar{q} \rightarrow t\bar{t}H$

Jet-associated Z production:

$$q\bar{q} \rightarrow Z + ng$$

Jet-associated W Z production:

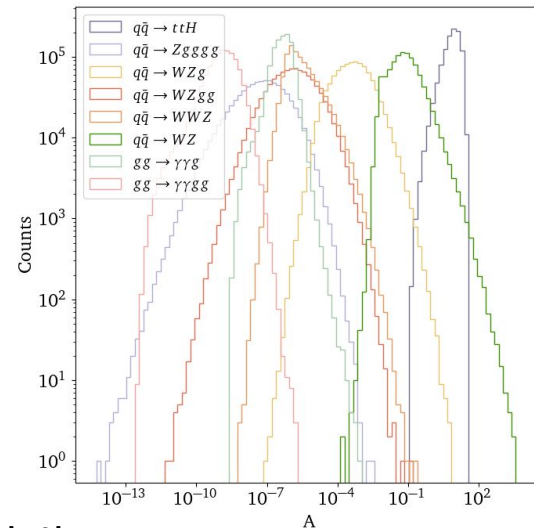
$$q\bar{q} \rightarrow WZ, q\bar{q} \rightarrow WZg, q\bar{q} \rightarrow WZgg$$

W W Z production:

$$q\bar{q} \rightarrow WWZ$$

Jet-associated di-photon production:

$$gg \rightarrow \gamma\gamma g, gg \rightarrow \gamma\gamma gg$$



Wide **variety** of processes covering many different characteristics

Neural Networks

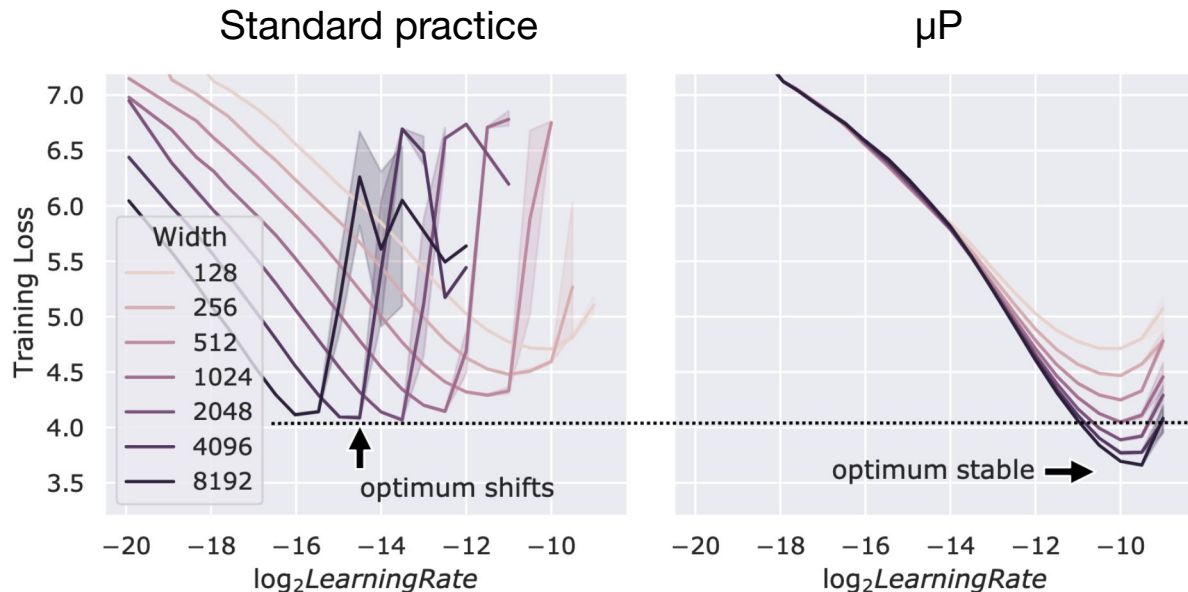
MLP set up:

- ▶ Preprocessing:
 - Inputs: **Invariants** derived from 4-vectors
 - Amplitudes: **Log-standardized** $\hat{a} = \frac{\log(a) - \overline{\log(a)}}{\sigma_{\log(a)}}$ (except $q\bar{q} \rightarrow t\bar{t}H$)
- ▶ Unless explicitly changed:
 - 500 Hidden neurons, 4 hidden layers, **$\sim 10^6$ parameters**
 - **10^6 Training dataset size**
 - **10^4 Epochs**
- ▶ **GELU** non-linearities
- ▶ Batchsize = 256
- ▶ **Cosine Annealing** scheduler
- ▶ **Hyperparameter Transfer**

Hyperparameter Transfer^[1]

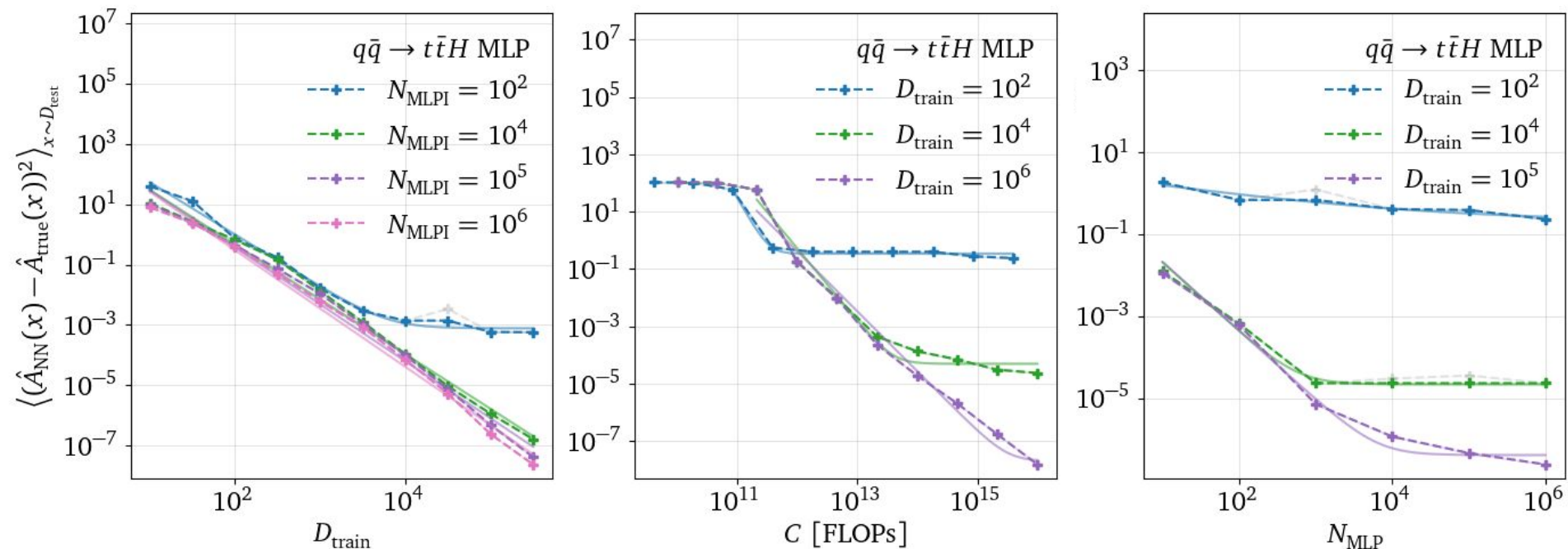
Optimal HP remain consistent
across widths

- ▶ Better initialization
- ▶ Better optimizer
- ▶ Optimizer treats layers differently
- ▶ Results are comparable or better than standard practice



Allows to optimize HP on a smaller (cheaper) model **once** and use them for different network sizes

Results: MLP, $q\bar{q} \rightarrow t\bar{t}H$ MSE loss



Very clean power laws, consistent slopes

Losses

MSE: $\mathcal{L}_{\text{MSE}} = \left\langle (A_{\text{true}}(x) - A_{\text{NN}}(x))^2 \right\rangle_{x \sim D_{\text{train}}}$

Heteroscedastic Loss:

Assume amplitude regression follows a normal distribution $p(A|x) = \mathcal{N}(A|\bar{A}(x), \sigma^2(x))$

Minimize negative log-likelihood:

$$\mathcal{L} = -\left\langle \log p(A|x) \right\rangle_{x \sim D_{\text{train}}} \longrightarrow \mathcal{L}_{\text{het}} = \left\langle \frac{(A_{\text{true}}(x) - \bar{A}(x))^2}{2\sigma^2(x)} + \log \sigma(x) \right\rangle_{x \sim D_{\text{train}}}$$

The NN predicts 2 outputs: $\bar{A}(x)$ and $\sigma(x)$

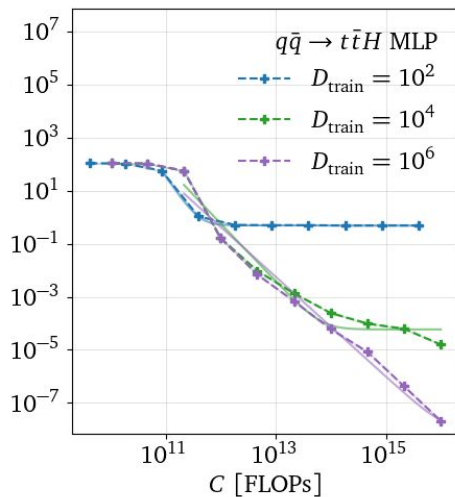
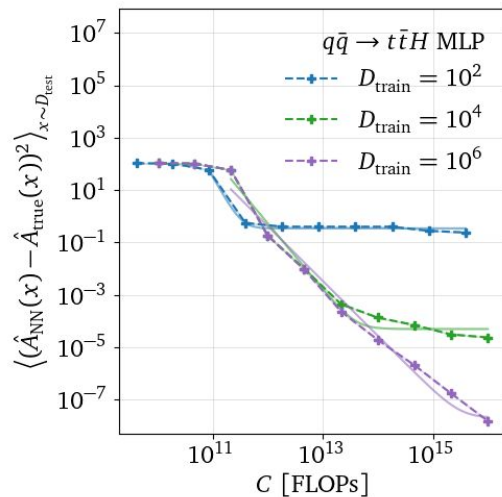
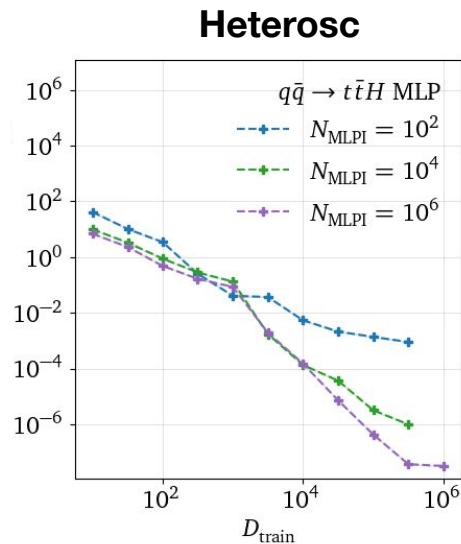
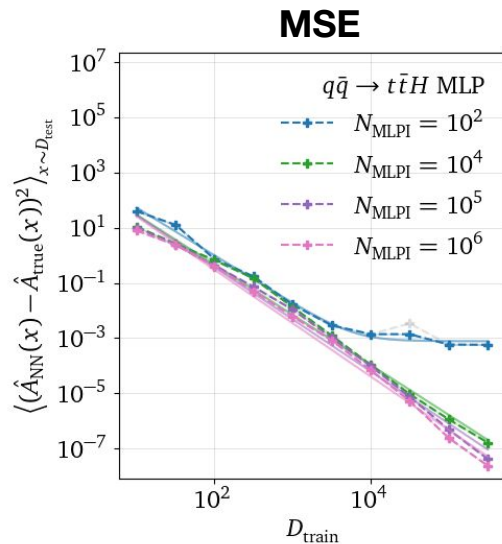
If it's well calibrated $t(x) = \frac{A_{\text{NN}}(x) - A_{\text{true}}(x)}{\sigma(x)}$ should follow a $\mathcal{N}(0, 1)$

Results: MLPI, $q\bar{q} \rightarrow t\bar{t}H$, Heteroscedastic loss

Slightly worse performance for scaling
in dataset size

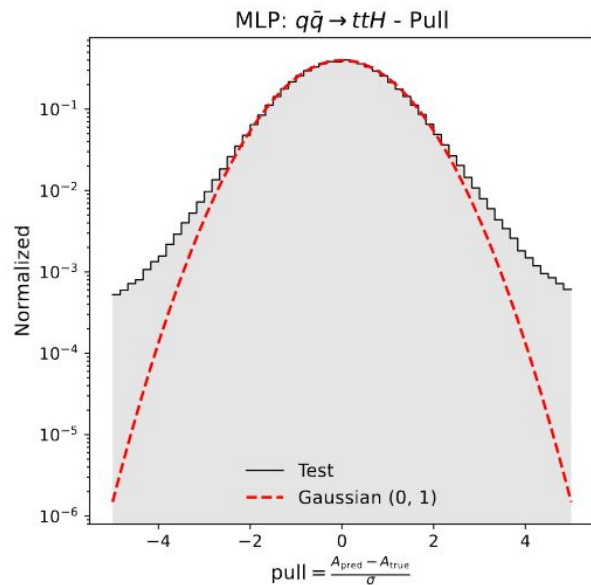
Very similar performance for scaling in
computing resources

**Important to note: Trained on
Heterosc loss, showing MSE**

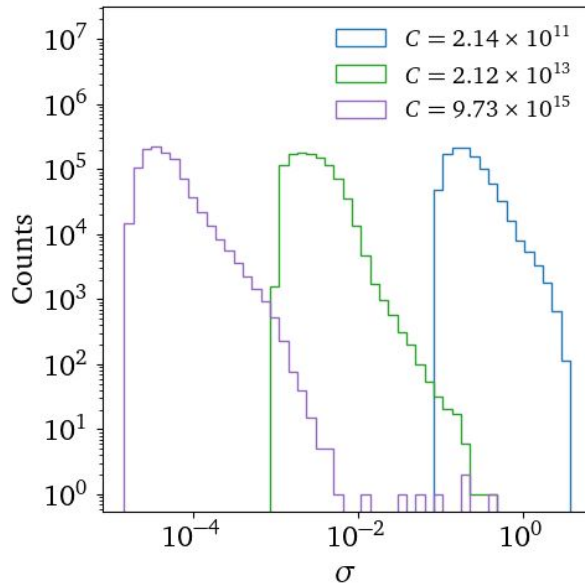


Uncertainties

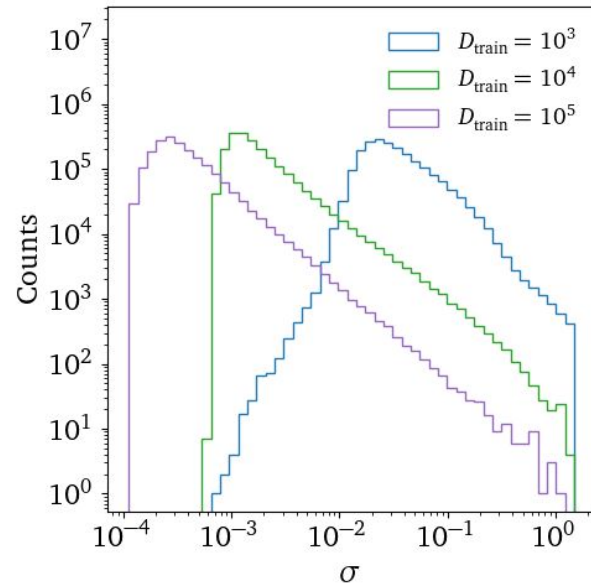
Calibration



$N=10^6$

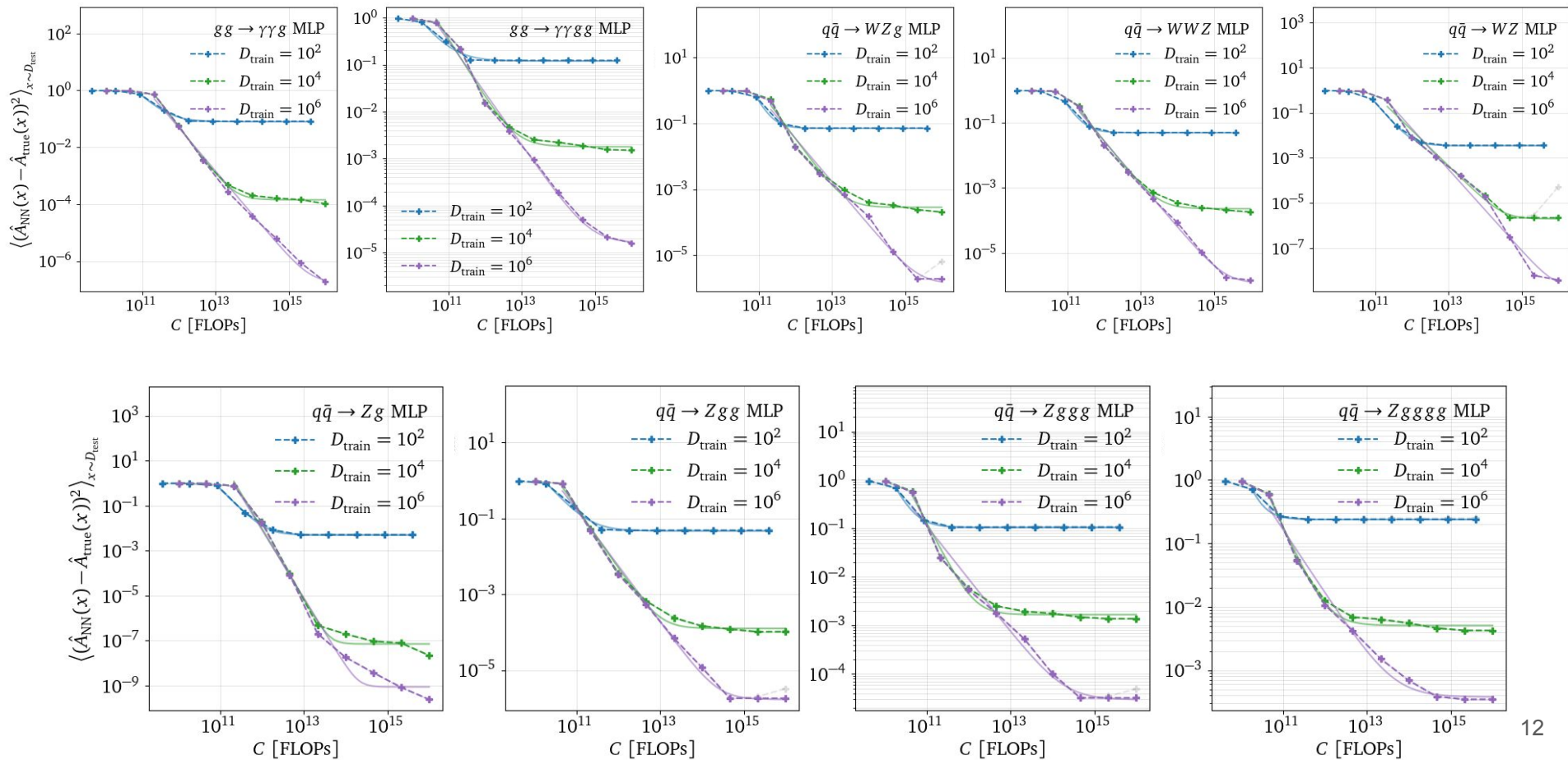


$N=10^4$

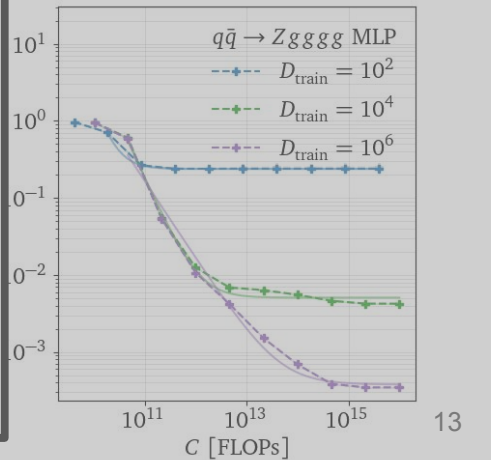
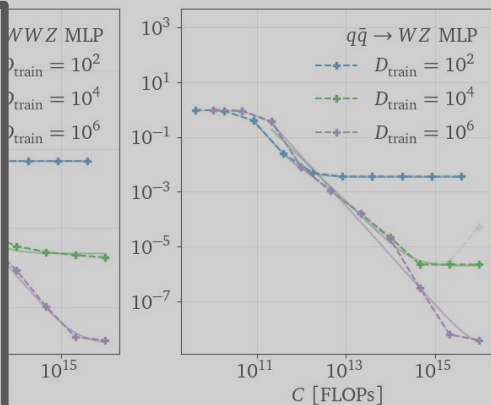
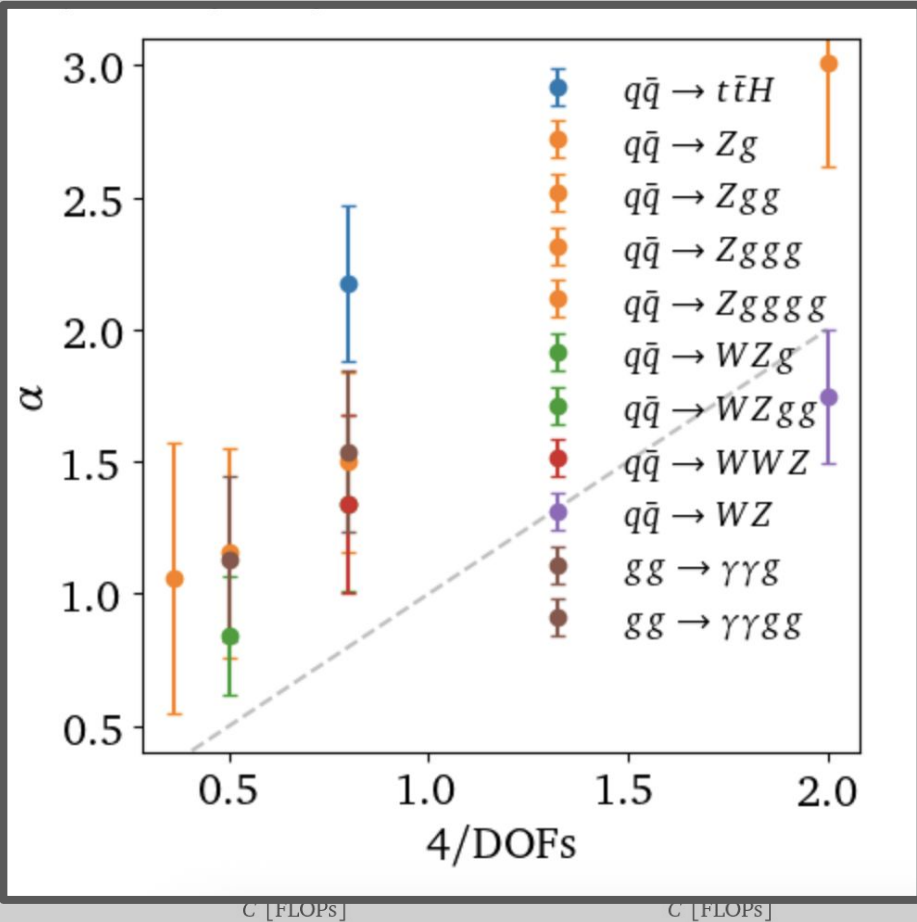
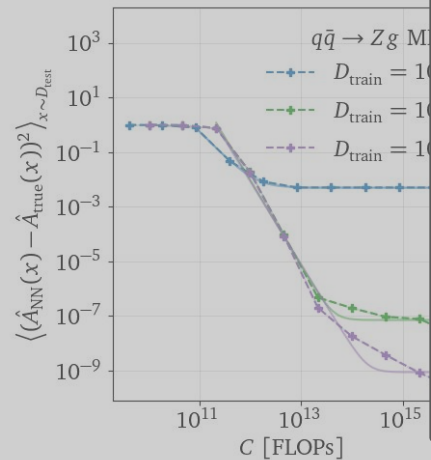
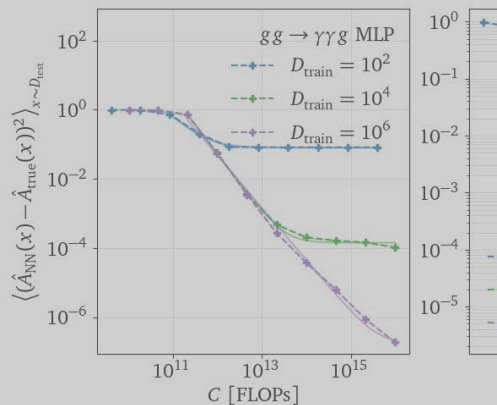


Well calibrated uncertainties for large enough dataset sizes

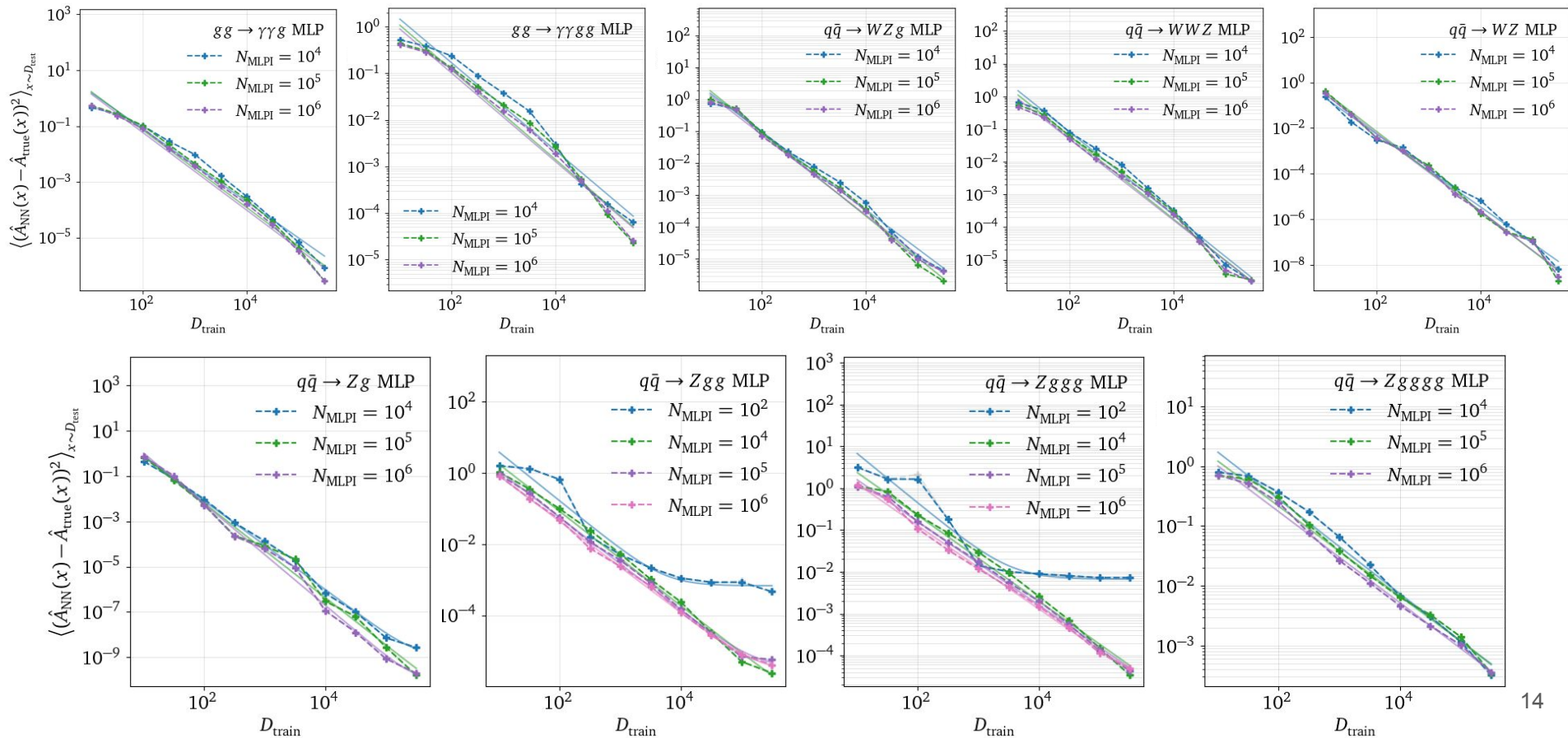
Results scaling on computing resources



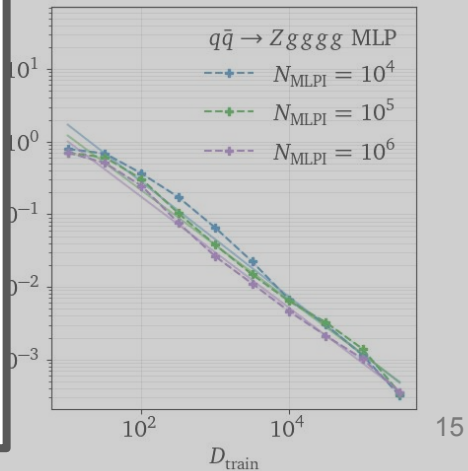
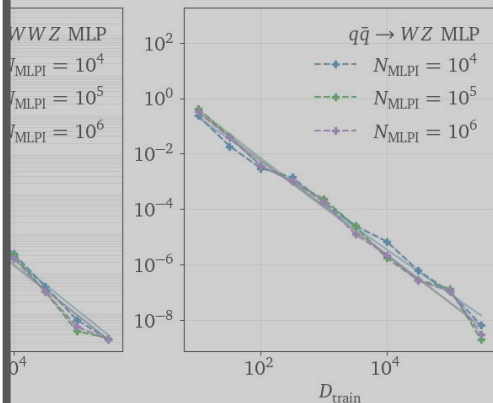
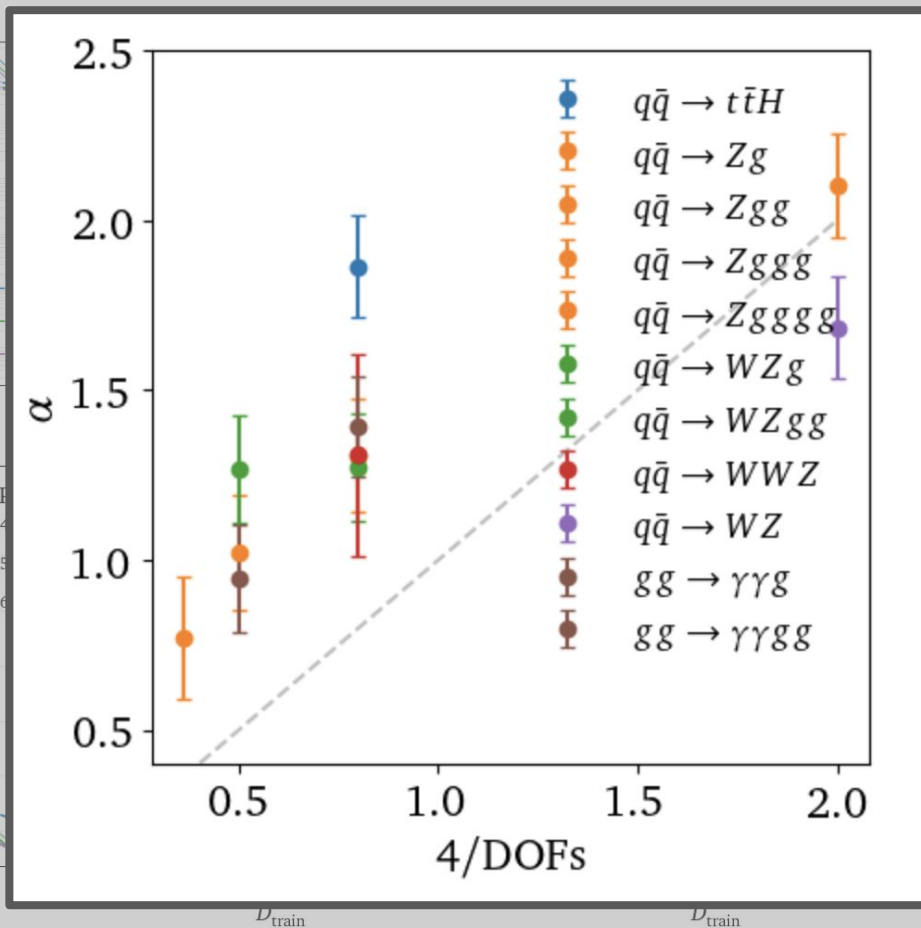
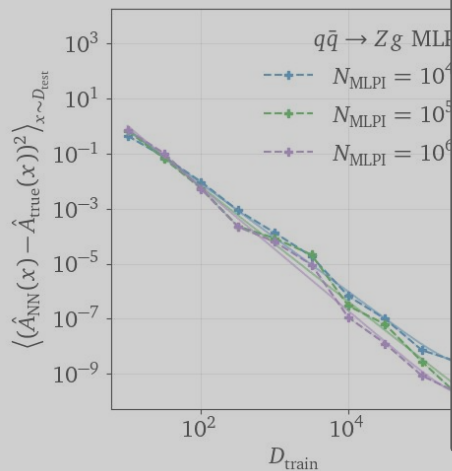
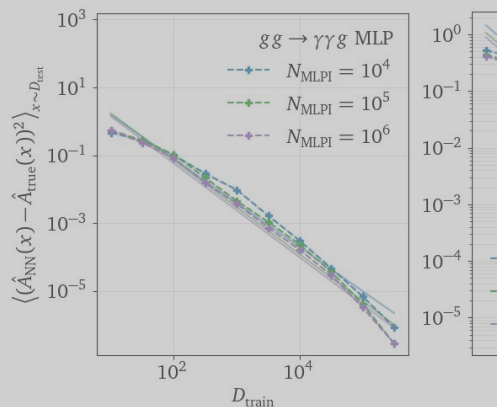
Results scaling on computing resources



Results scaling on dataset size

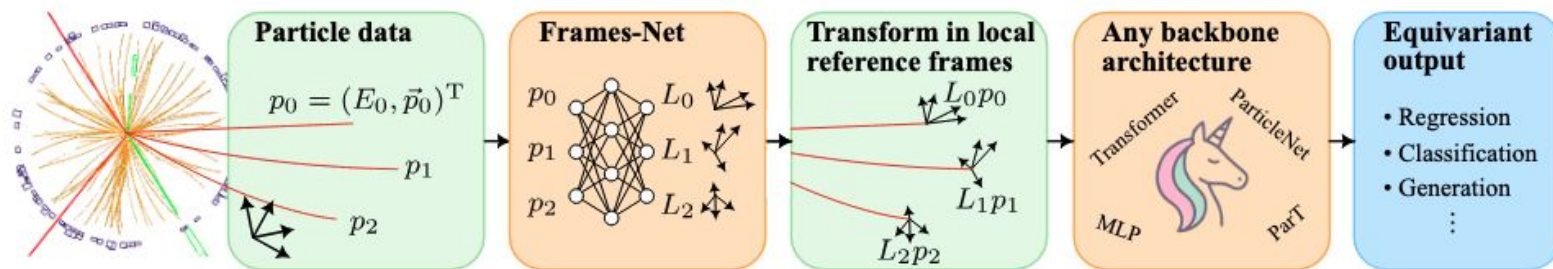


Results scaling on dataset size

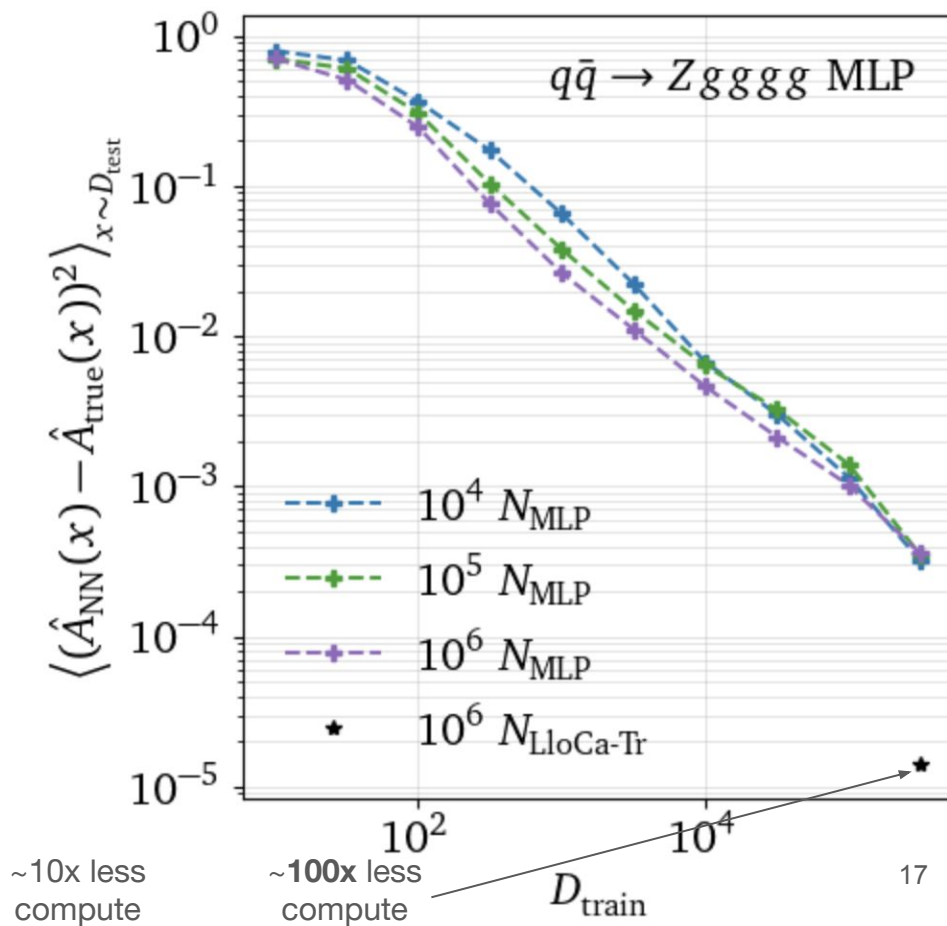
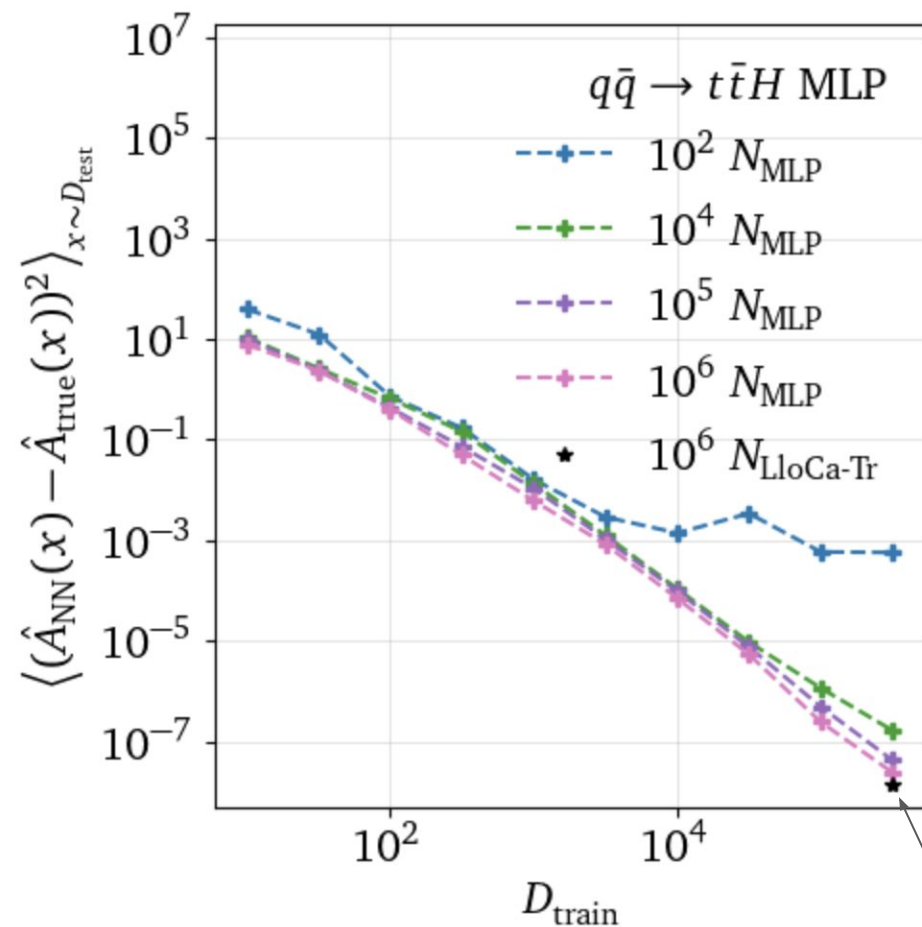


Work in progress: LloCa^[1]

Permutation invariance and equivariance yield better results



Work in progress: LloCa



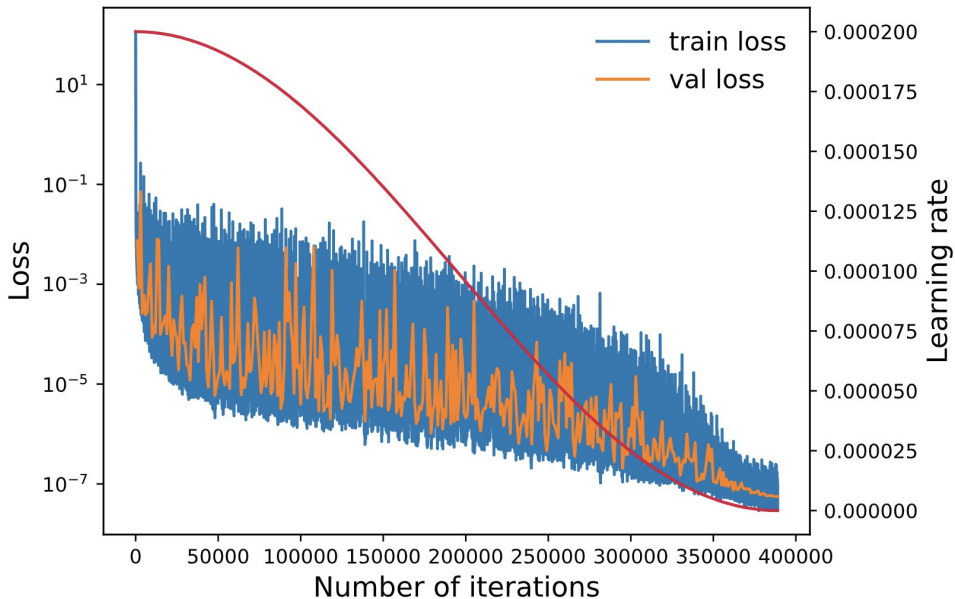
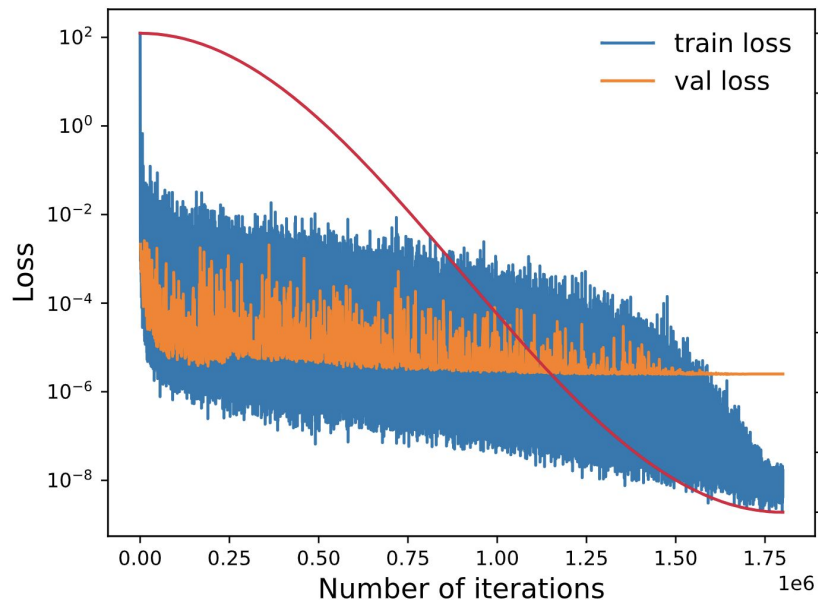
Conclusions

- ▶ Observed very clean and predictable scaling laws for amplitude surrogates
- ▶ Dataset size bottlenecks are the most important
- ▶ Scaling with MSE and Heterosc loss
- ▶ Well calibrated uncertainties across many orders of magnitude
- ▶ Observed relationship between scaling and DOFs
- ▶ Promising results with equivariant NNs

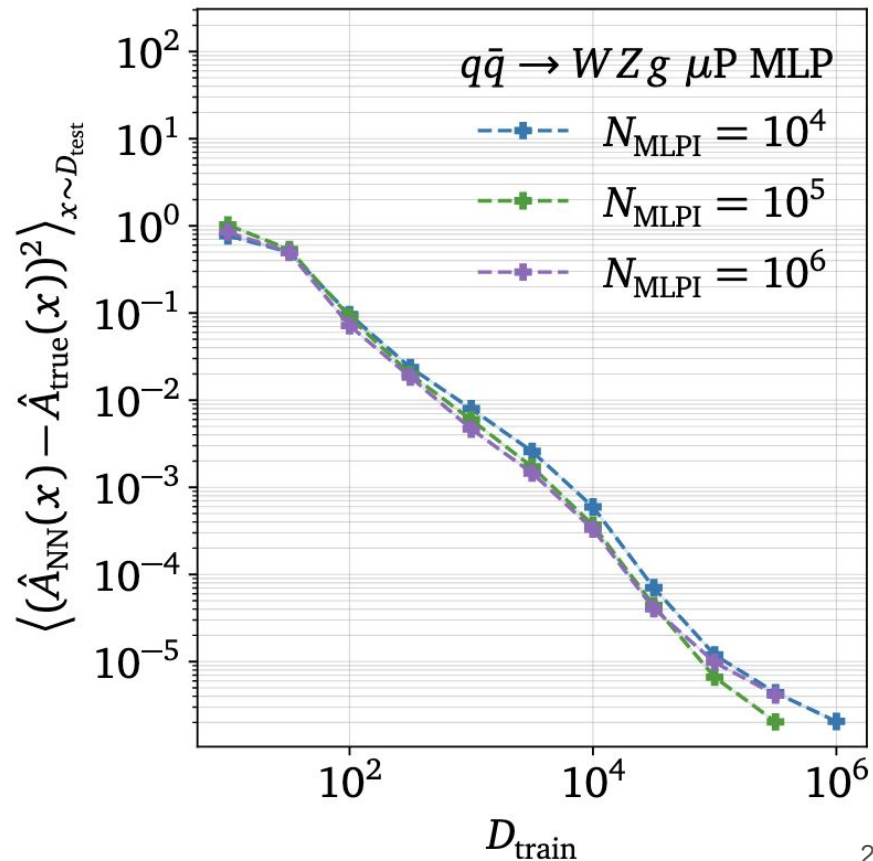
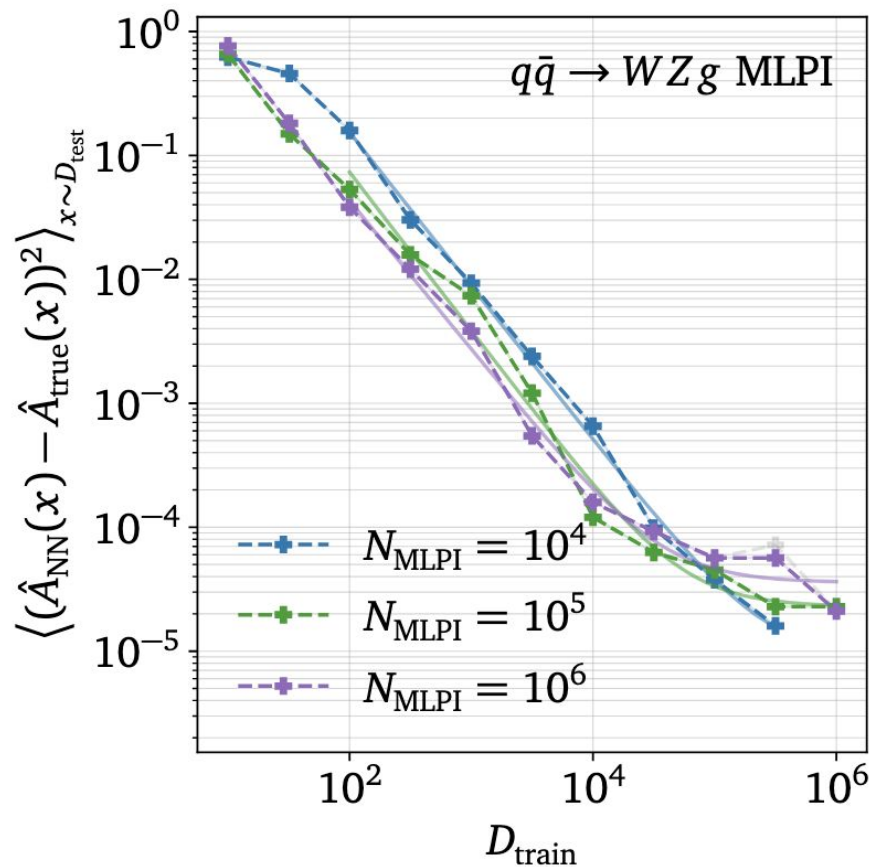
Thanks!

Scheduler

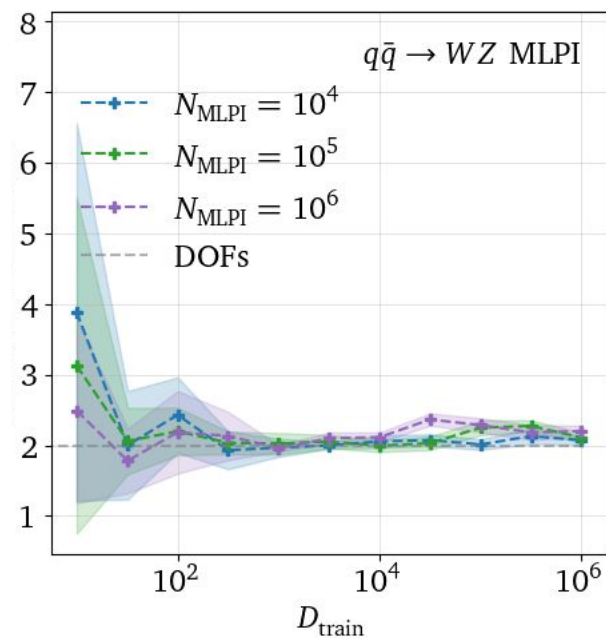
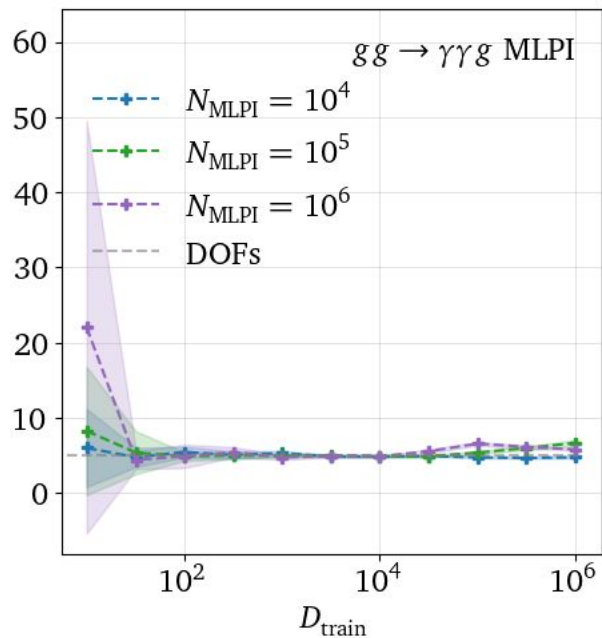
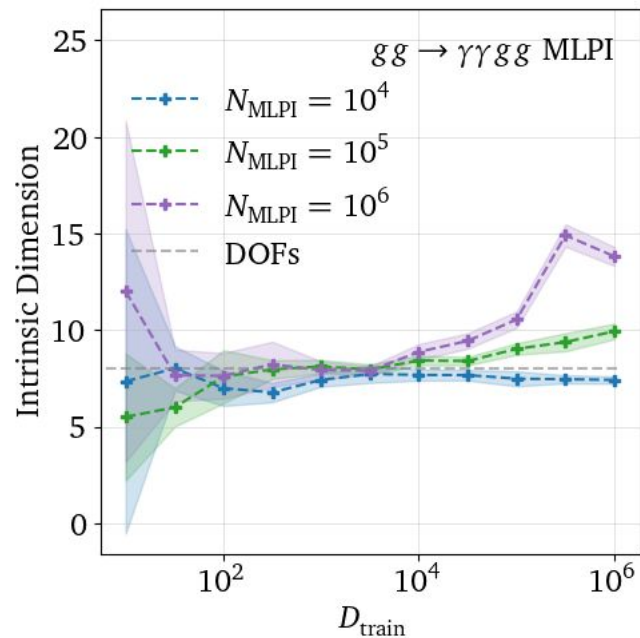
Scheduler has to go to 0 **just at the right time**



muP vs normal MLP

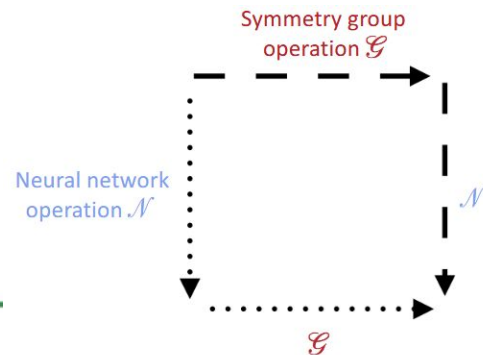


Internal Intrinsic Dimension



L-GATr = Equivariance + Transformer

$$\mathcal{G}(\mathcal{N}(x)) = \mathcal{N}(\mathcal{G}(x))$$

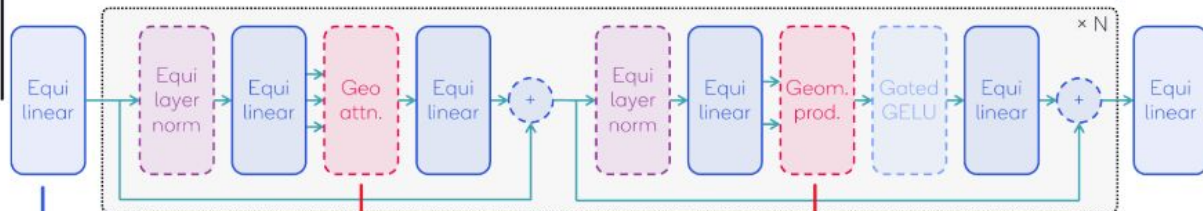


Input and output data

can have one or multiple token dimensions

Attention blocks

can be stacked to large depth, gradients are propagated efficiently



Linear layers

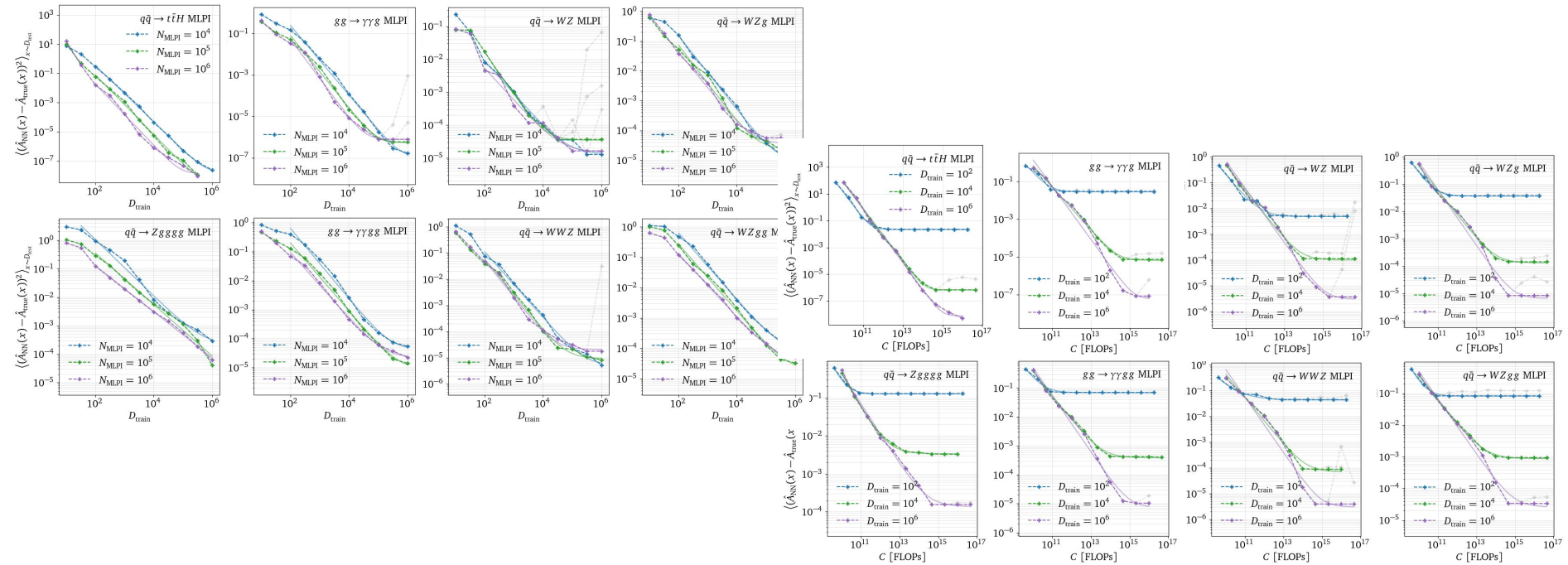
between GA representations with equivariance constraint

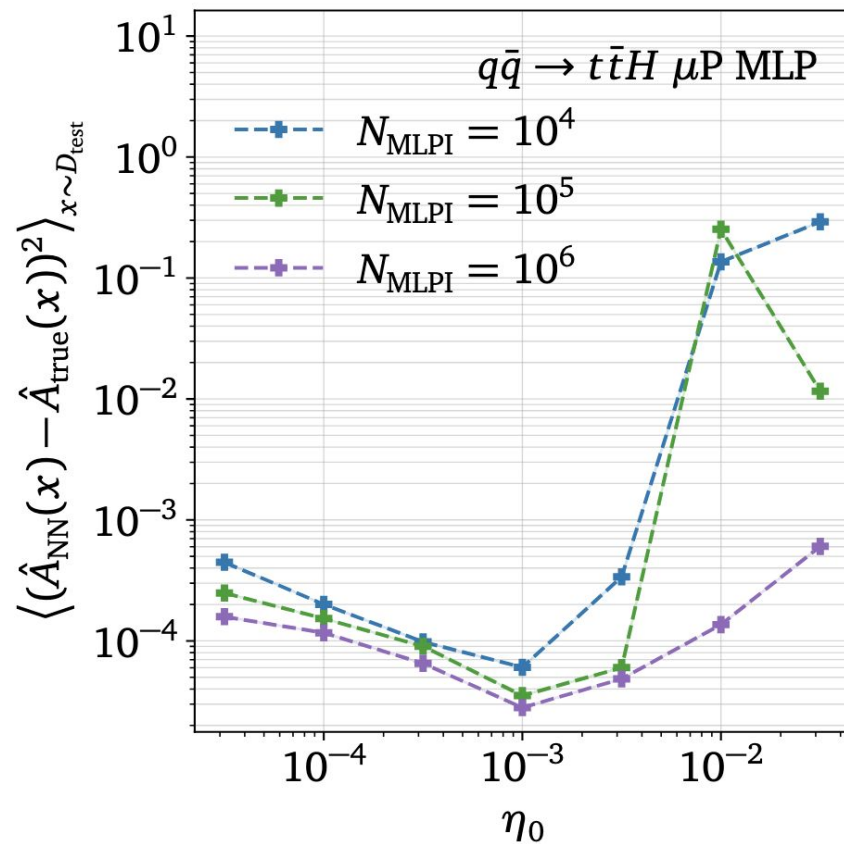
Geometric attention

generalizes scaled dot-product attention

Geometric product

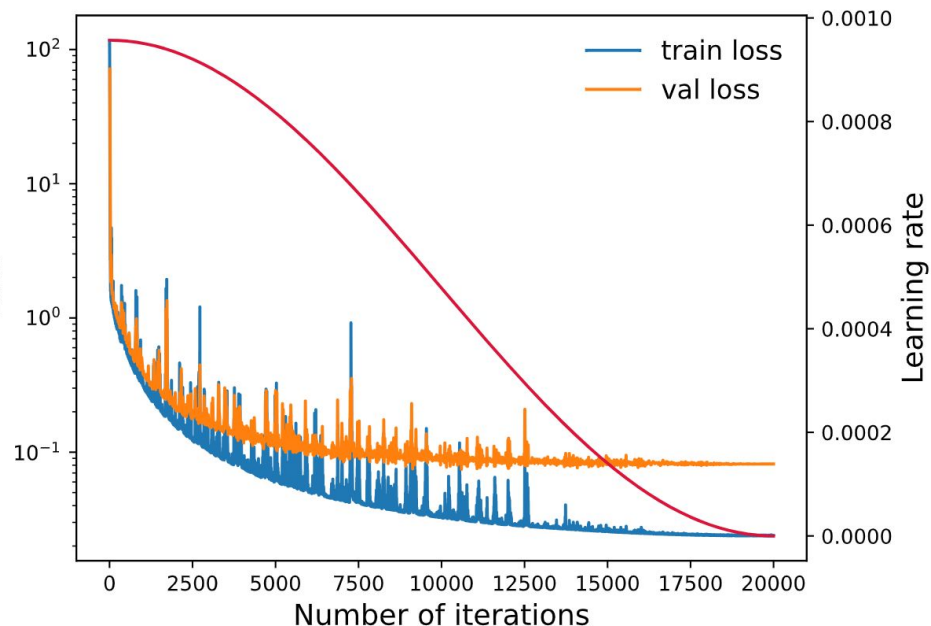
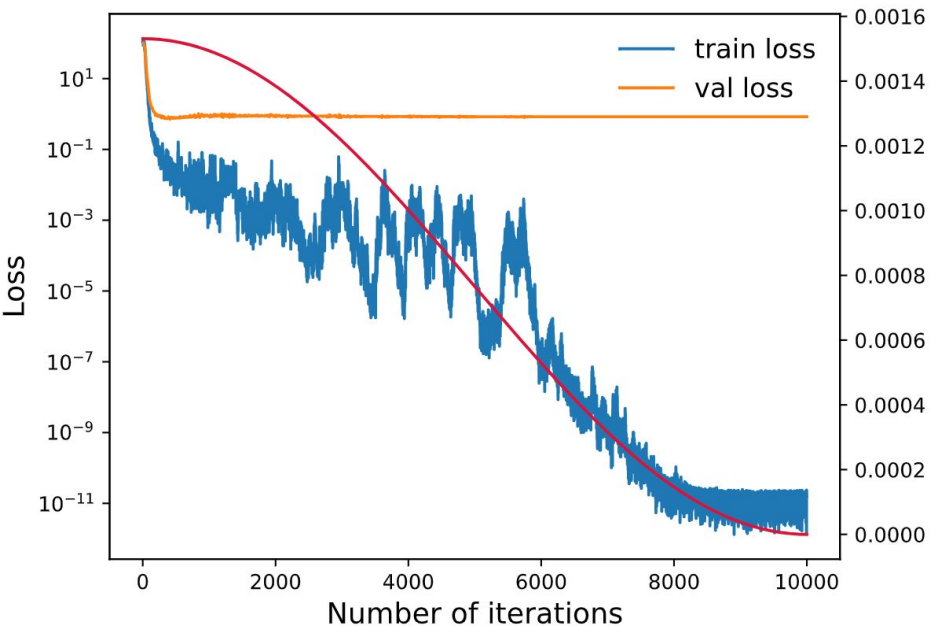
allow for construction of new geometric types



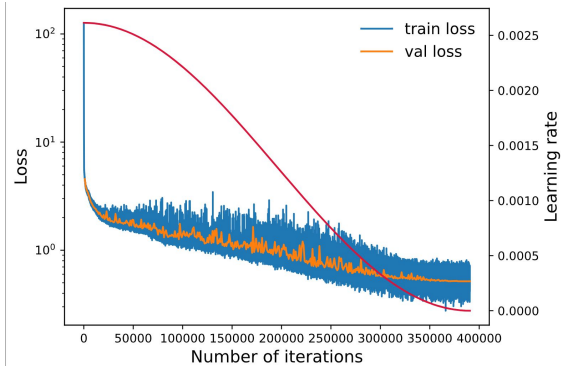


muP overfitting for small datasets

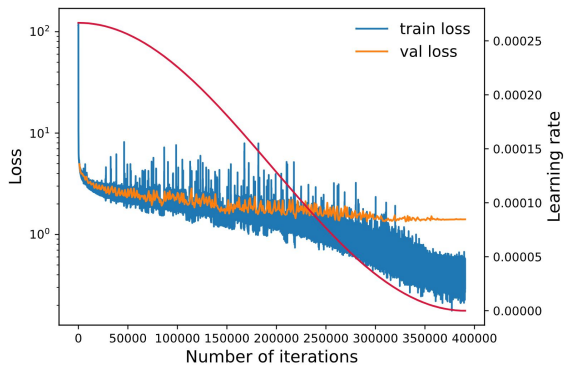
100 training points



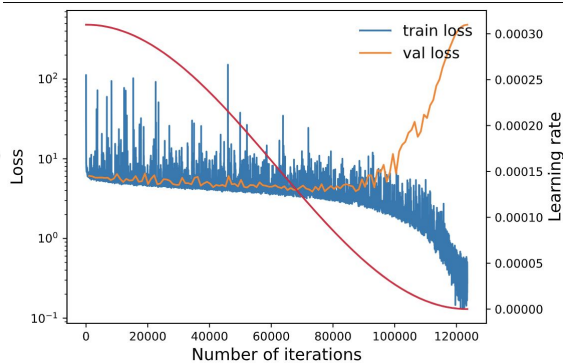
Overfitted Het loss



$\mu\text{P MLP: } q\bar{q} \rightarrow t\bar{t}H$ - Pull



$\mu\text{P MLP: } q\bar{q} \rightarrow t\bar{t}H$ - Pull



$\mu\text{P MLP: } q\bar{q} \rightarrow t\bar{t}H$ - Pull

