

What do Transformers learn when trained via Masked Language Modelling?

An analysis in the framework of the Generalized Potts model

Federica Gerace

Joint work with:



Riccardo Rende



Alessandro Laio

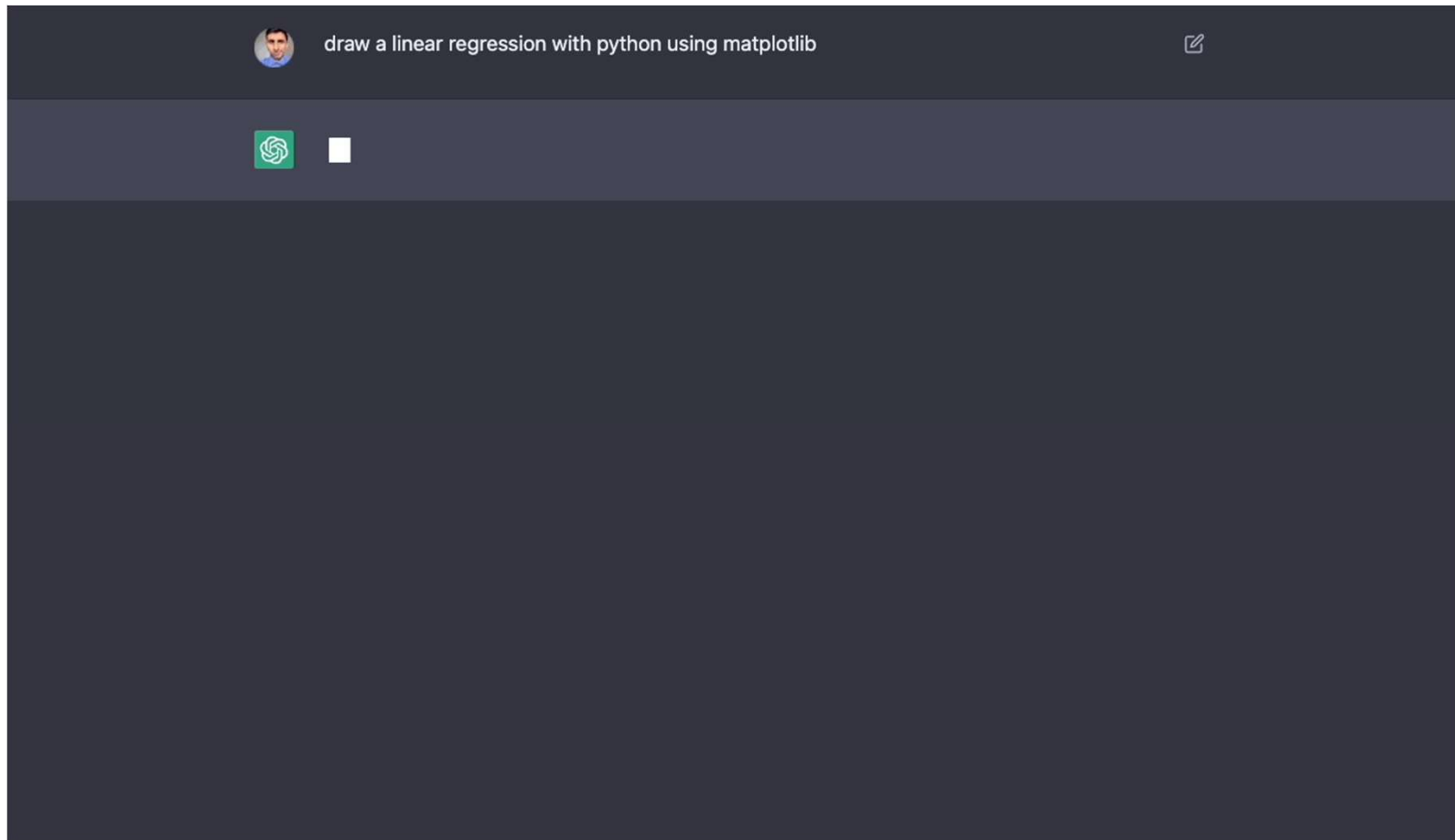


Sebastian Goldt



Bridging scales: At the crossroads among renormalisation group, multi-scale modelling, and deep learning - 16/04/2024

A famous Transformer example: Chat-GPT!



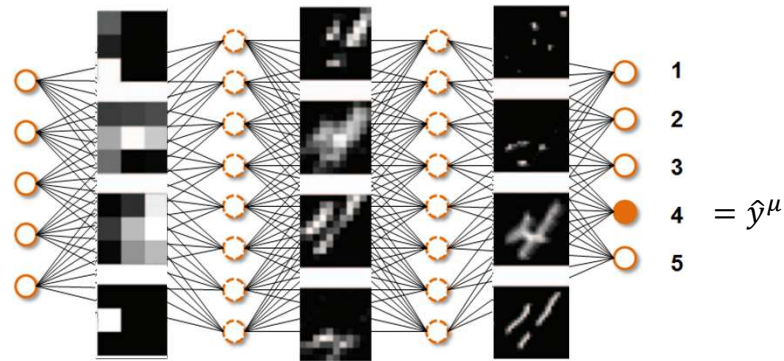
Artificial Neural Networks in short

- **Training Phase:**

Training Set

$$\mathcal{D} = \{\mathbf{x}^\mu, y^\mu\}_{\mu=1}^n$$

$$\mathbf{x}^\mu = \begin{matrix} \text{[Handwritten digit 4]} \\ y^\mu = 4 \end{matrix}$$



$$\mathbf{w}^{t+1,l} \leftarrow \mathbf{w}^{t,l} + \eta \frac{\partial \mathcal{L}}{\partial \mathbf{w}^{t,l}}$$

Loss Function

$$\mathcal{L} = \sum_{\mu=1}^n \ell(y^\mu, \hat{y}^\mu(\mathbf{w}; \mathbf{x}^\mu))$$

- **Testing Phase:**

$$\epsilon_g = \mathbb{E}_{\{\mathbf{x}_{new}, y_{new}\}} \left[(y^{new} - \hat{y}^{new}(\hat{\mathbf{w}}; \mathbf{x}^{new}))^2 \right]$$

$$\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmin}} \mathcal{L}(\mathbf{w}|\mathcal{D})$$

Transformers are taking the show

Language Modelling on WikiText-103

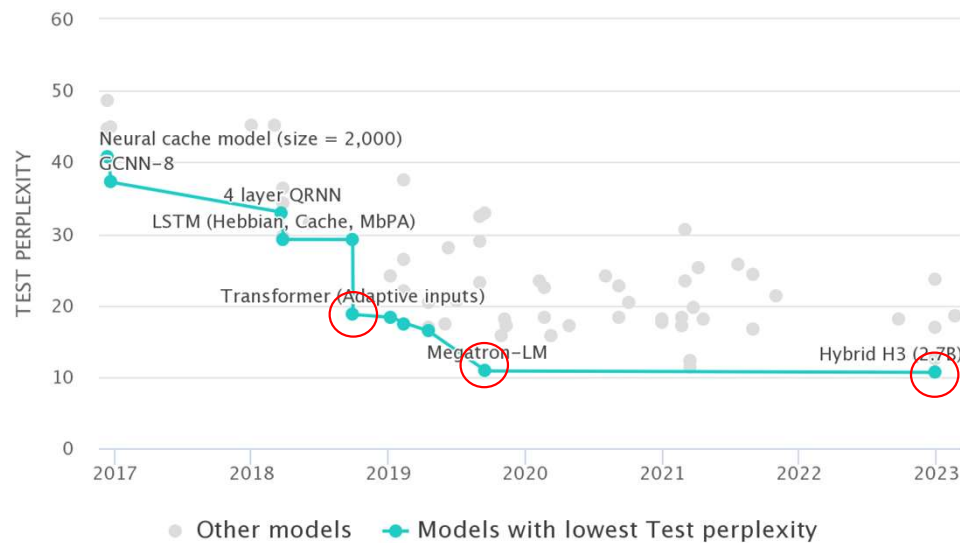
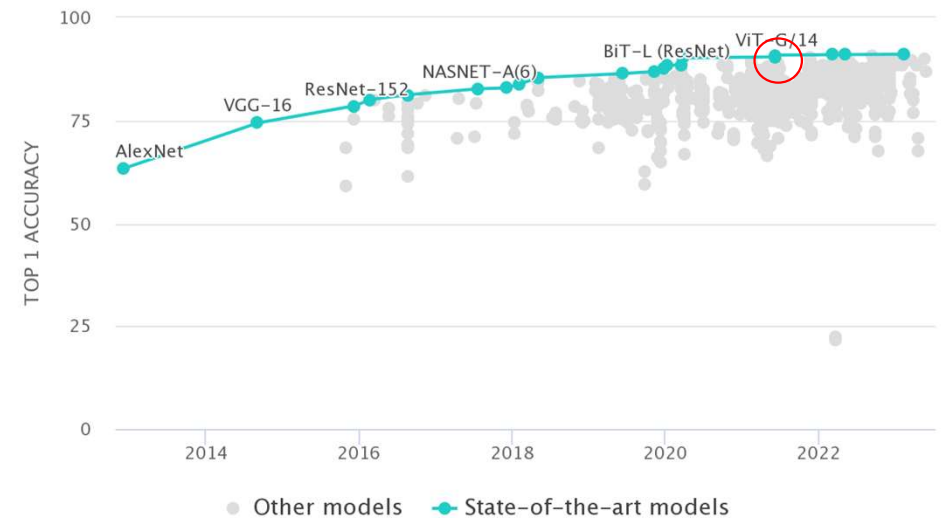


Image Classification on ImageNet



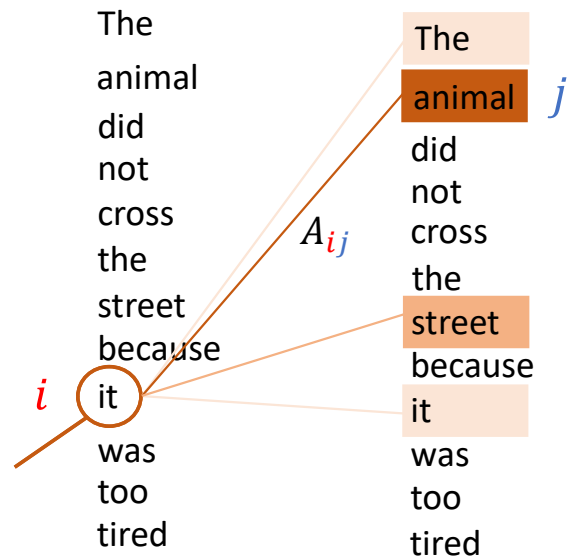
Transformers learn the context



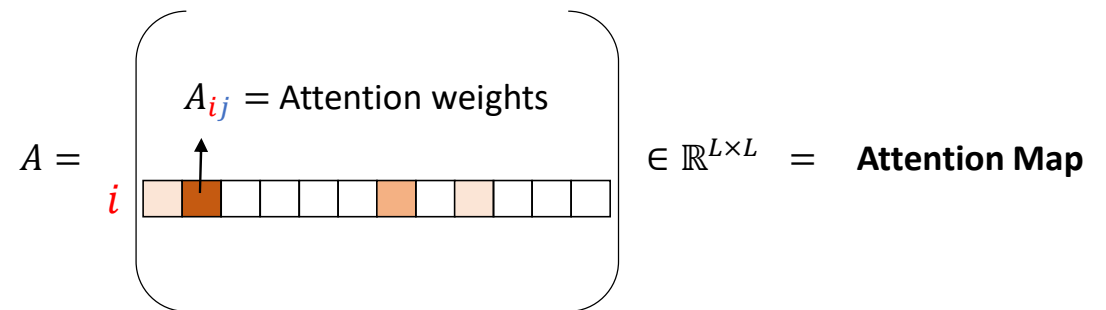
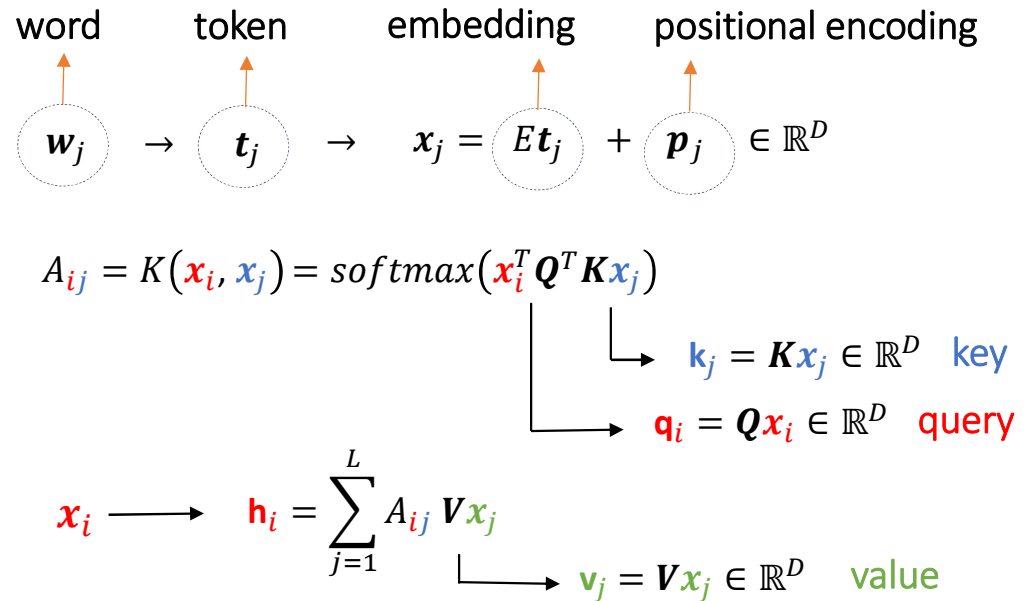
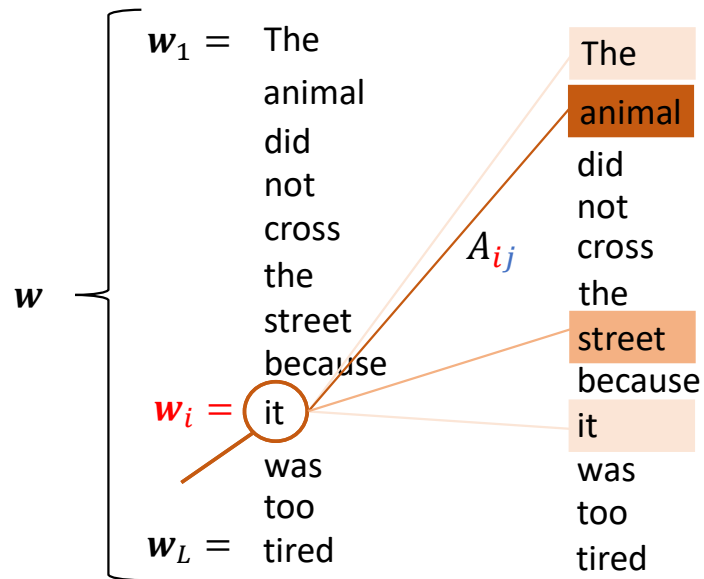
The animal did not cross the street because it was too tired.

The animal did not cross the street because it was too wide.

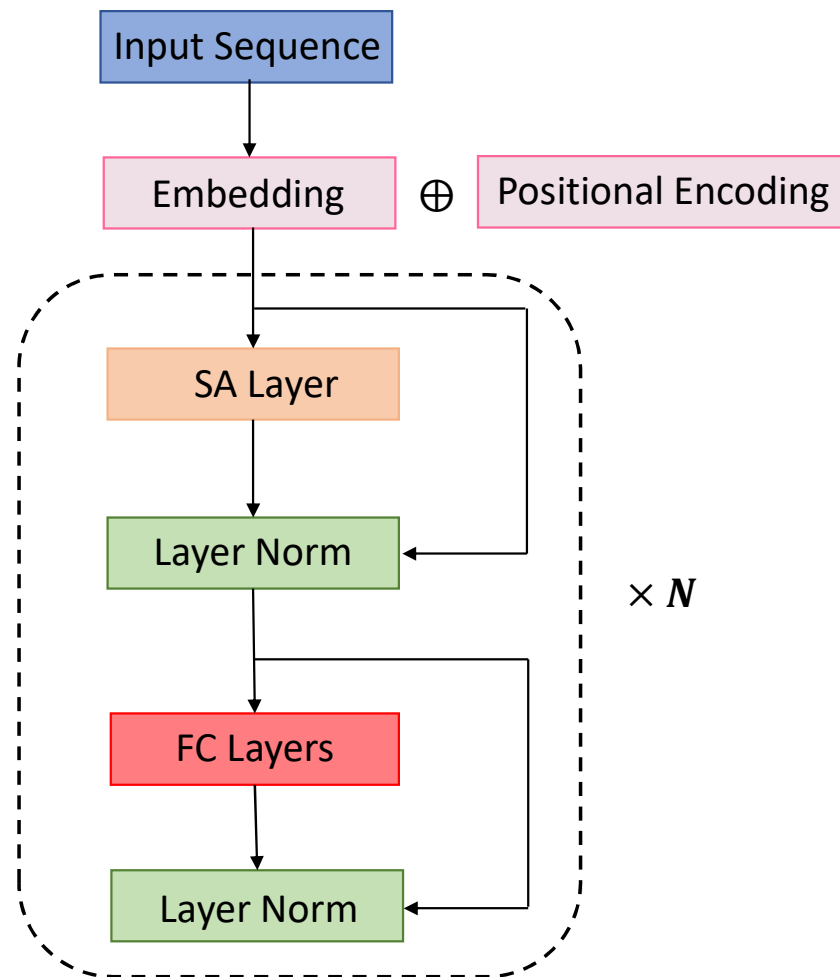
Self-Attention Mechanism



Self-Attention Mechanism



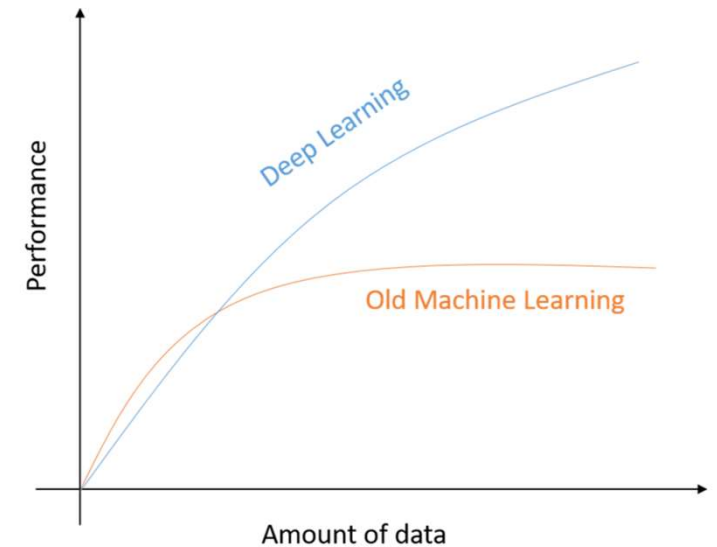
The building block of Transformers



Deep learning is data-hungry

“The analogy to deep learning is that the rocket engine is the deep learning models and the fuel is the huge amounts of data we can feed to these algorithms.”

Andrew Ng



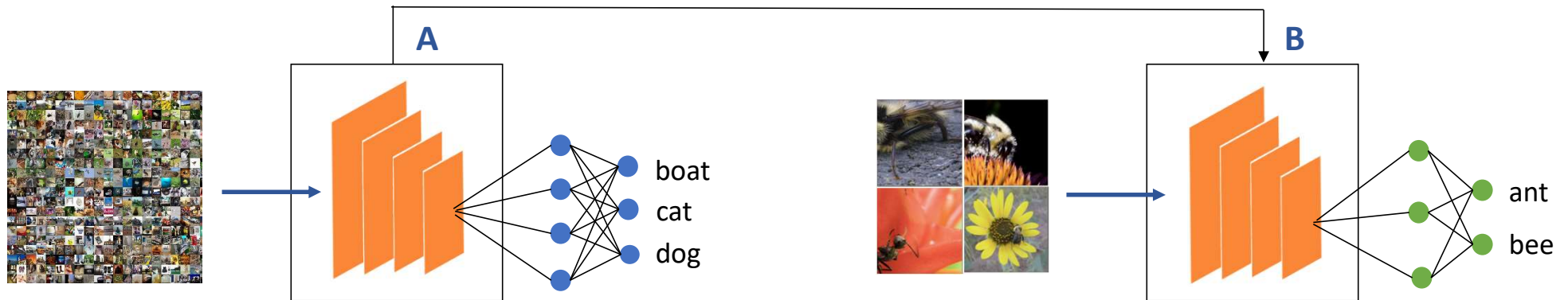
A possible solution: Transfer Learning

DATA-ABUNDANT
SOURCE TASK



DATA-SCARCE
TARGET TASK

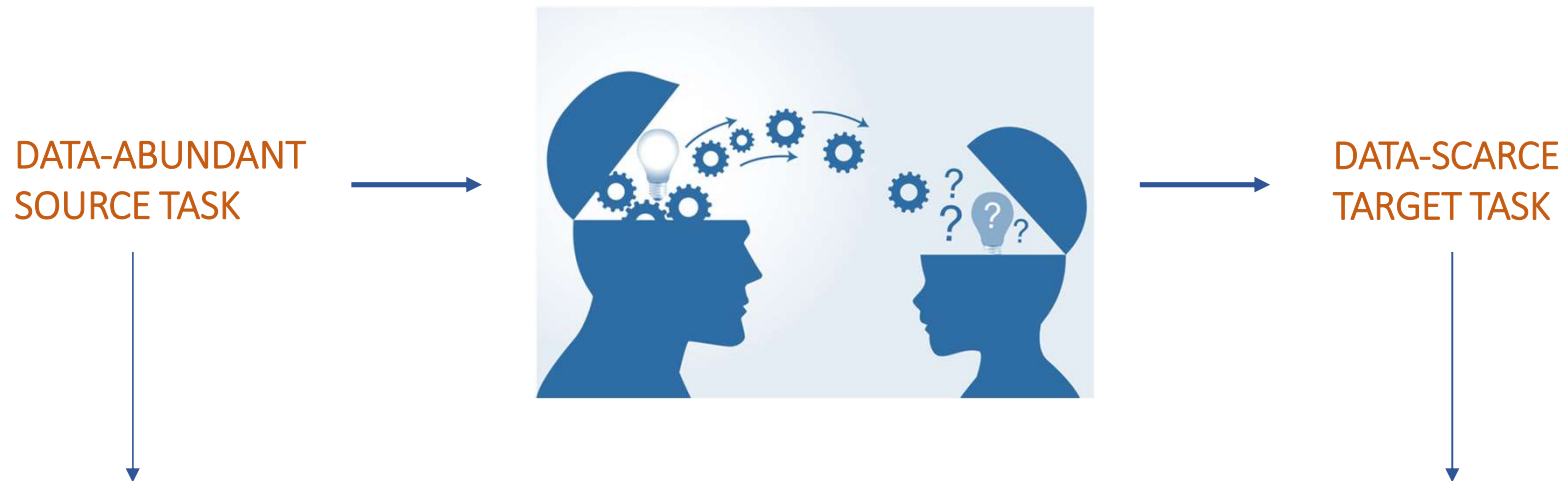
Feature map transfer



Gerace, Sarao Mannelli, Saglietti, Saxe, Zdeborová, *Machine Learning: Science and Technology*, 2022

Gerace, Doimo, Sarao Mannelli, Saglietti, Laio, *arXiv:2303.01429*, 2023

A possible solution: Transfer Learning



- Masked Language Modelling (MLM)

Input: $\mathbf{x}^\mu =$ We MASKED milk chocolate $\mu = 1, \dots, M$

Label: $\mathbf{y}^\mu =$ eat

Goal: MASKED = ?

- Text Translation;
- Text Generation (Chat-GPT);
- Sentiment Analysis...

Some open questions

-  What self-attention learns with Masked Language Modelling?
-  How many samples are required to achieve good generalization performances?

A key ingredient: **The Generalized Potts Model!**

The Generalized Potts Model

$$\mathbf{w} = \{w_1, \dots, w_i, \dots, w_L\}$$

word:

$$w_i = \text{street}$$



$$\mathbf{s} = \{s_1, \dots, s_i, \dots, s_L\}$$

one-hot encoded Potts spin:

$$s_i = [s_{i1}, \dots, s_{ic}, \dots, s_{iC}] \in \{0,1\}^C$$

C = size of the vocabulary

Each sequence is sampled from the Gibbs Measure of the Generalized Potts Model:

$$\mathcal{P}(\mathbf{s}) \propto \exp(-\beta \mathcal{H}(\mathbf{s}))$$

$$\mathcal{H}(\mathbf{s}) = \frac{1}{2} \sum_{i,j} J_{ij} \mathbf{s}_i^T U \mathbf{s}_j$$

$U \in \mathbb{R}^{C \times C}$ = Interaction among colors.

$J \in \{0,1\}^{L \times L}$ = Interaction among sites;

MLM with the Generalized Potts Model

Given a training set made of M *masked sequences* $\mathbf{s}_{\setminus i}^{\mu}$ with the corresponding *masked spin value* s_i^{μ} , e.g.:

$$\mathcal{D} = \left\{ \mathbf{s}_{\setminus i}^{\mu}, s_i^{\mu} \right\}_{\mu=1}^M$$

with each sequence \mathbf{s}^{μ} sampled from the *Generalized Potts Model*, e.g.: $\mathbf{s}^{\mu} \sim \mathcal{P}(\mathbf{s}^{\mu}; J, U)$.

The *goal* is to achieve the lowest possible generalization loss, e.g.:

$$\epsilon_g = -\mathbb{E}_{\mathbf{s} \sim \mathcal{P}(\mathbf{s}^{new}; J, U)} \left[\frac{1}{L} \sum_{i=1}^L \sum_{c=1}^C s_{ic}^{new} \log \left(\hat{\mathcal{P}}_{ic}(\mathbf{s}^{new}; A, V) \right) \right]$$

Exact masked spin value

True Gibbs Measure Transformer Gibbs Measure prediction

Question 1



What self-attention learns with Masked Language Modelling?

Vanilla Transformer on Generalized Potts

Task:

$$\mathcal{D} = \left\{ \mathbf{s}_{\setminus i}^{\mu}, \mathbf{s}_i^{\mu} \right\}_{\mu=1}^M$$

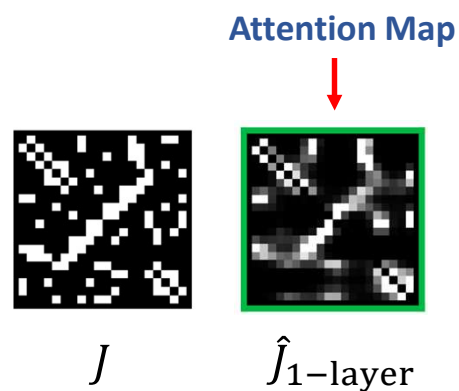
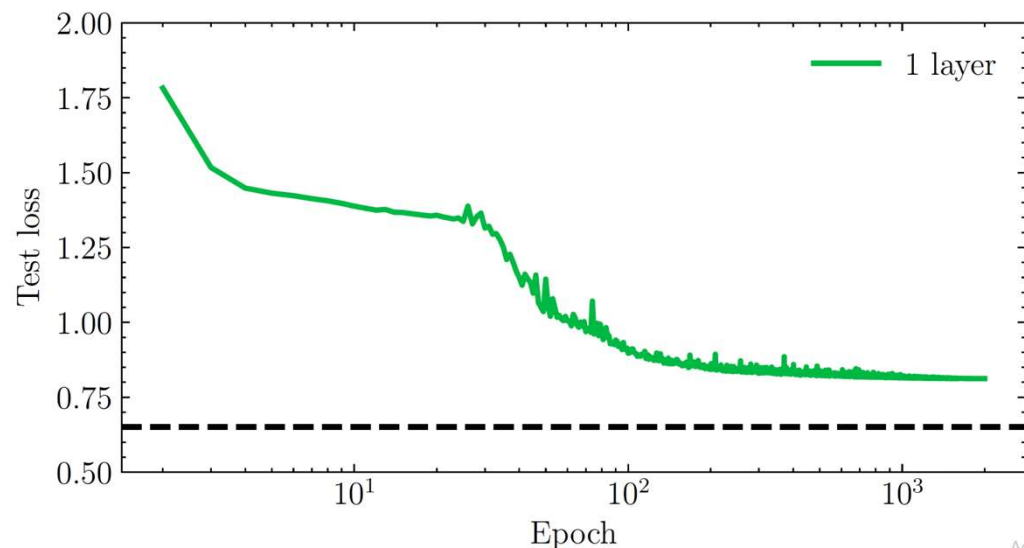
Test Loss:

$$\epsilon_g = -\mathbb{E}_{\mathbf{s} \sim \mathcal{P}(\mathbf{s}; J, U)} \left[\frac{1}{L} \sum_{i=1}^L \sum_{c=1}^C s_{ic} \log \left(\hat{\mathcal{P}}_{ic}(\mathbf{s}; A, V) \right) \right]$$

Exact masked spin value

Transformer Gibbs Measure estimate

True Gibbs Measure



Vanilla Transformer on Generalized Potts

Task:

$$\mathcal{D} = \left\{ \mathbf{s}_{\setminus i}^\mu, \mathbf{s}_i^\mu \right\}_{\mu=1}^M$$

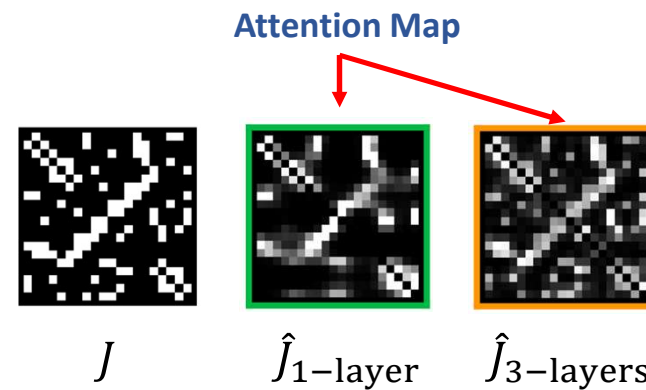
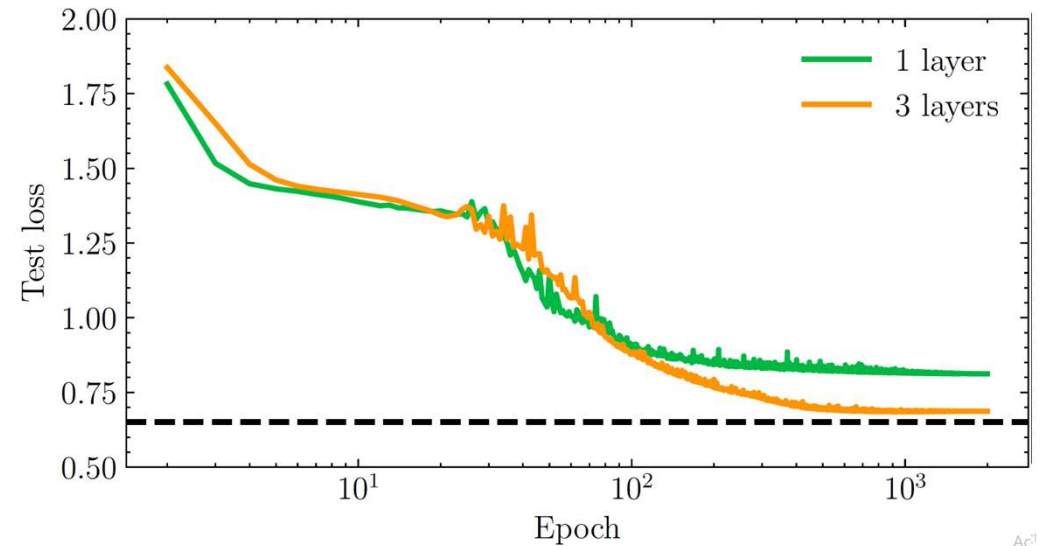
Test Loss:

$$\epsilon_g = -\mathbb{E}_{\mathbf{s} \sim \mathcal{P}(\mathbf{s}; J, U)} \left[\frac{1}{L} \sum_{i=1}^L \sum_{c=1}^C s_{ic} \log \left(\hat{\mathcal{P}}_{ic}(\mathbf{s}; A, V) \right) \right]$$

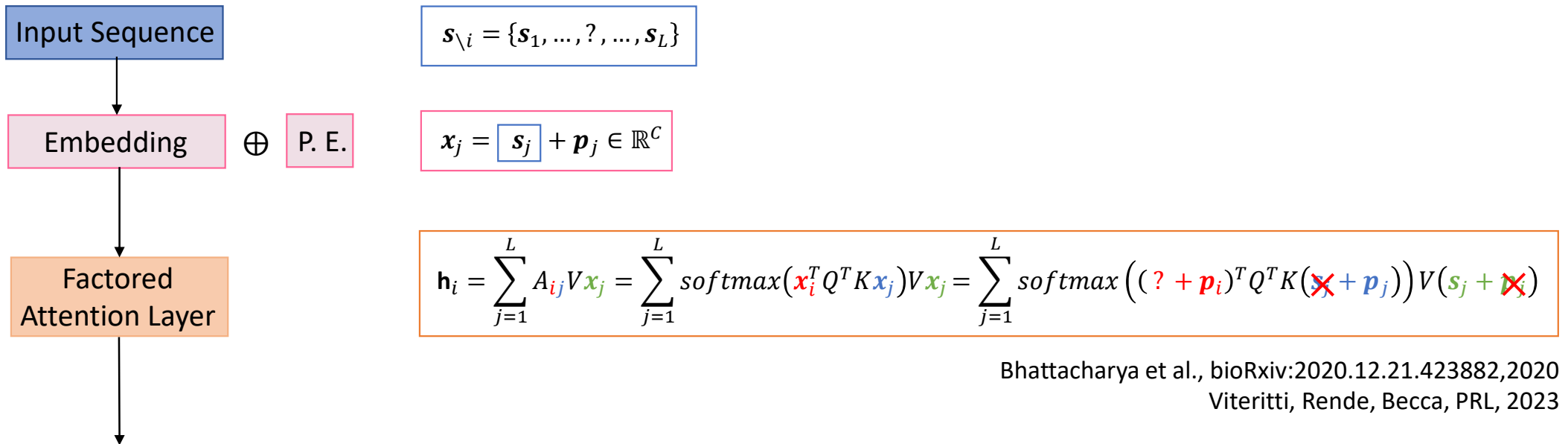
Exact masked spin value

Transformer Gibbs Measure estimate

True Gibbs Measure



Factored Self-Attention learns Generalized Potts



Transformer prediction

$$\hat{s}_{ic} \sim \text{softmax}(h_{ic}) \propto \exp\left(\sum_{j=1}^L A_{ij} (V s_j)_c\right)$$

$$A_{ij} \dashrightarrow \beta J_{ij}$$

$$V \dashrightarrow U$$

Exact Potts conditionals

$$s_{ic} \sim \exp\left(\beta \sum_{j=1}^L J_{ij} (U s_j)_c\right)$$

Rende, Gerace, Laio, Goldt, *Physical Review Research*, 2024

Vanilla Transformer on Generalized Potts

Task:

$$\mathcal{D} = \left\{ \mathbf{s}_{\setminus i}^{\mu}, \mathbf{s}_i^{\mu} \right\}_{\mu=1}^M$$

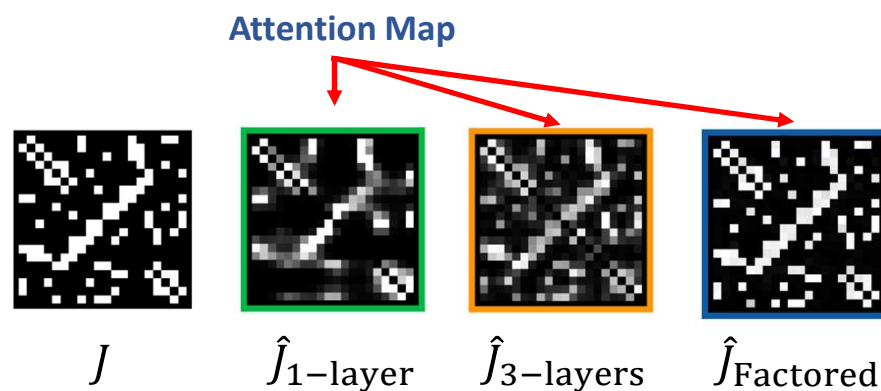
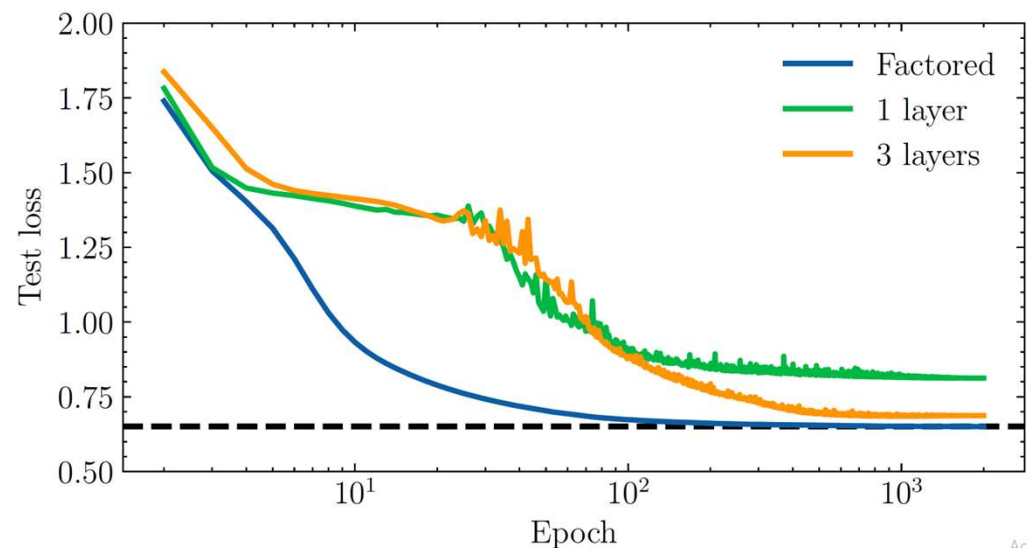
Test Loss:

$$\epsilon_g = -\mathbb{E}_{\mathbf{s} \sim \mathcal{P}(\mathbf{s}; J, U)} \left[\frac{1}{L} \sum_{i=1}^L \sum_{c=1}^C s_{ic} \log \left(\hat{\mathcal{P}}_{ic}(\mathbf{s}; A, V) \right) \right]$$

Exact masked spin value

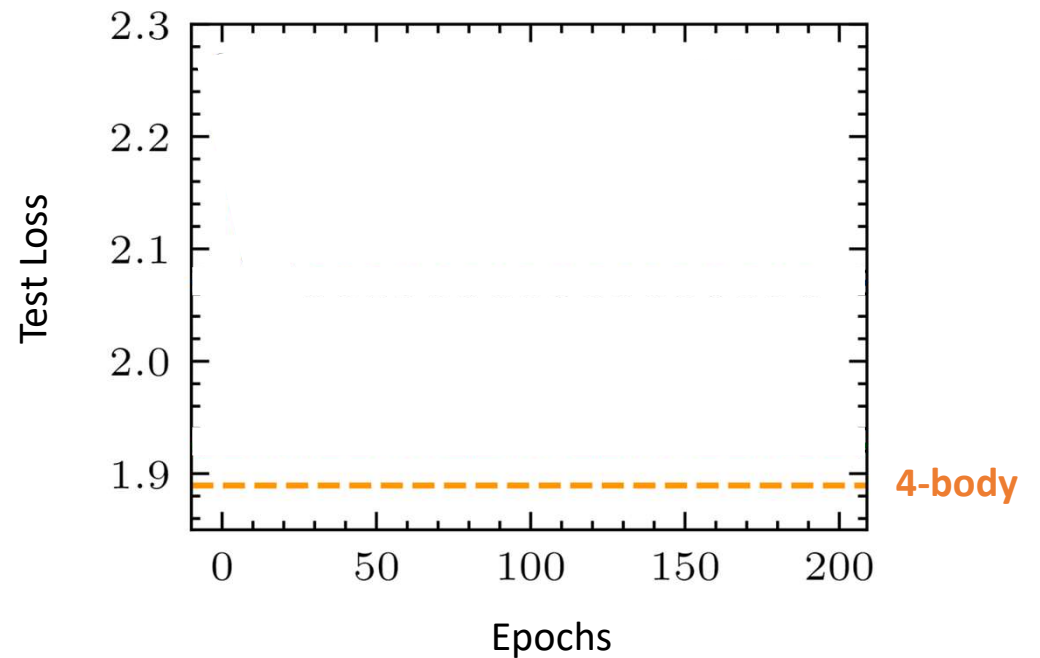
Transformer Gibbs Measure estimate

True Gibbs Measure

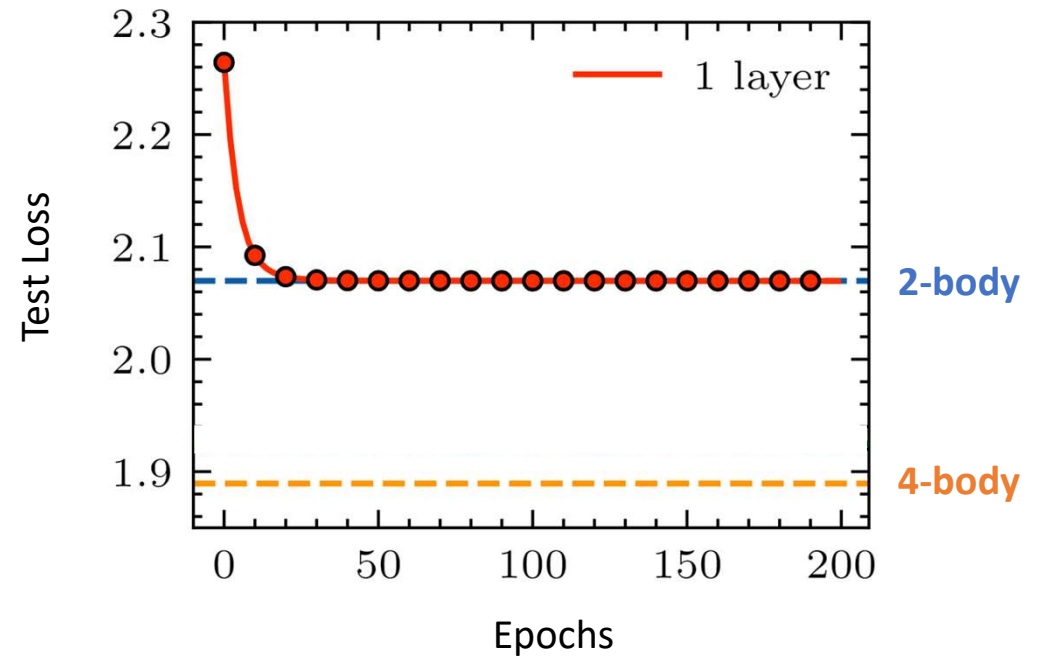
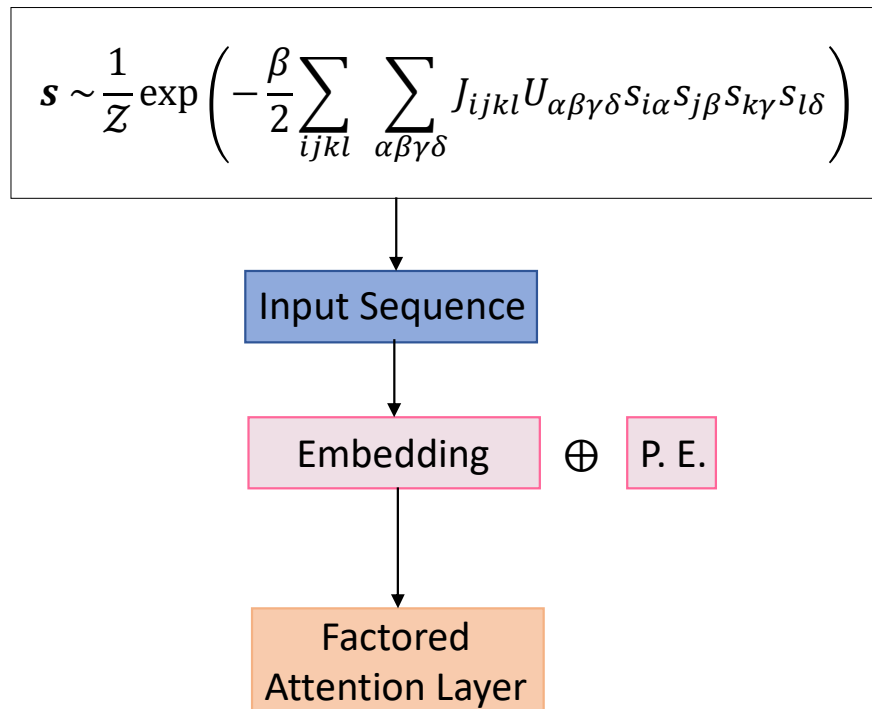


And if the interactions are of higher-order?

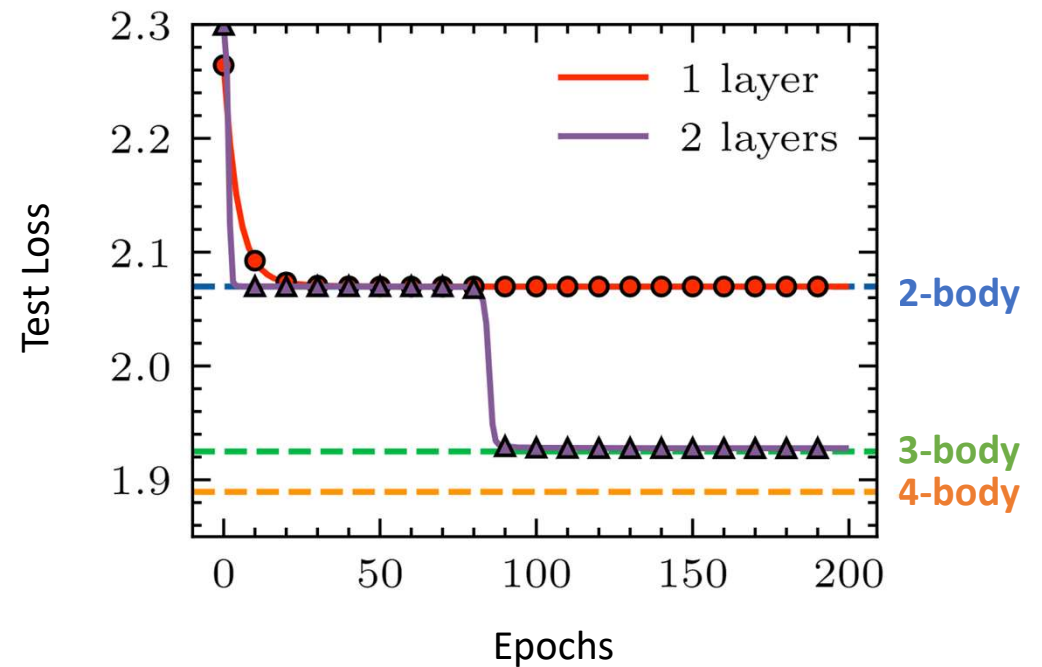
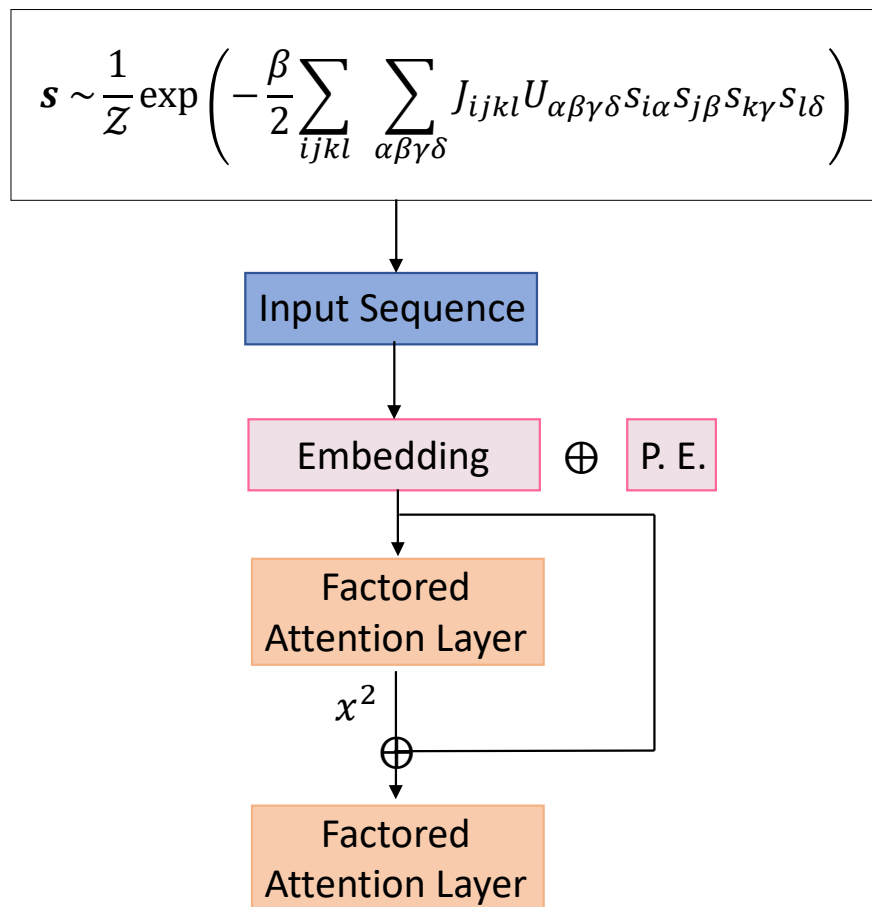
$$\mathbf{s} \sim \frac{1}{Z} \exp \left(-\frac{\beta}{2} \sum_{ijkl} \sum_{\alpha\beta\gamma\delta} J_{ijkl} U_{\alpha\beta\gamma\delta} s_{i\alpha} s_{j\beta} s_{k\gamma} s_{l\delta} \right)$$



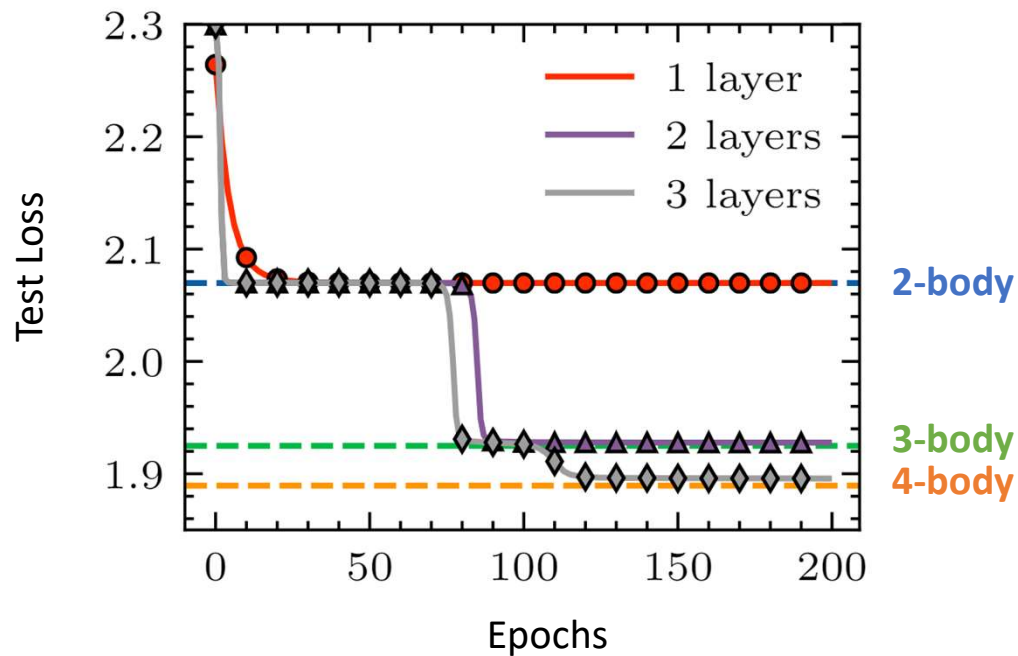
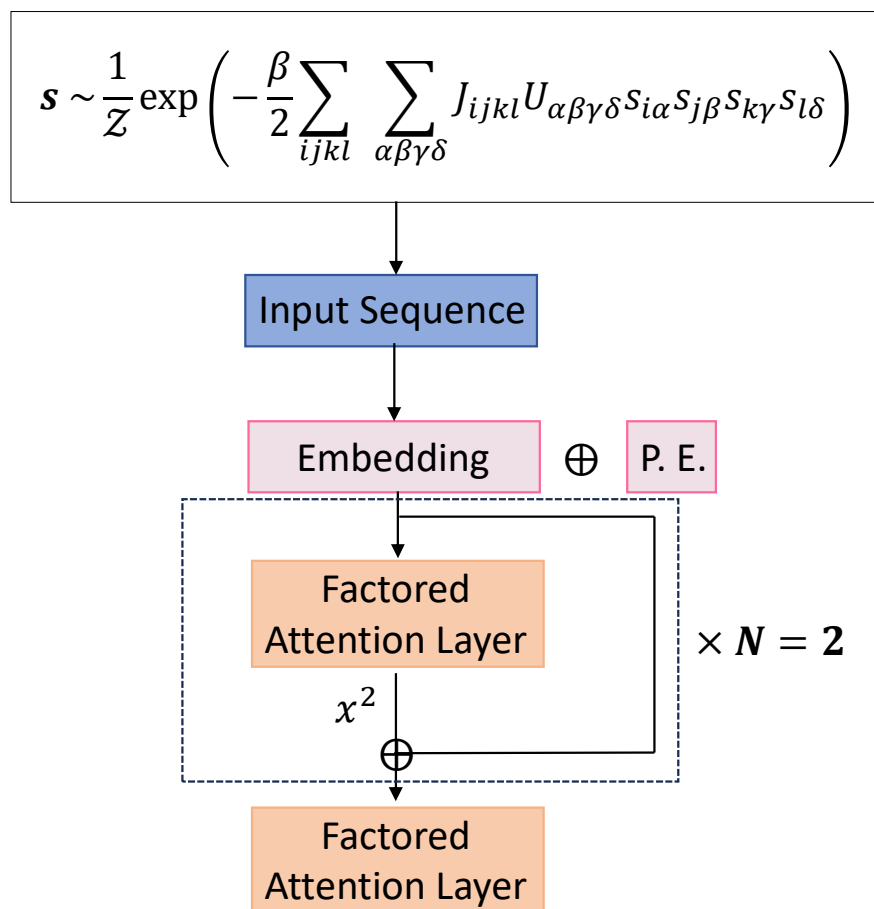
And if the interactions are of higher-order?



And if the interactions are of higher-order?

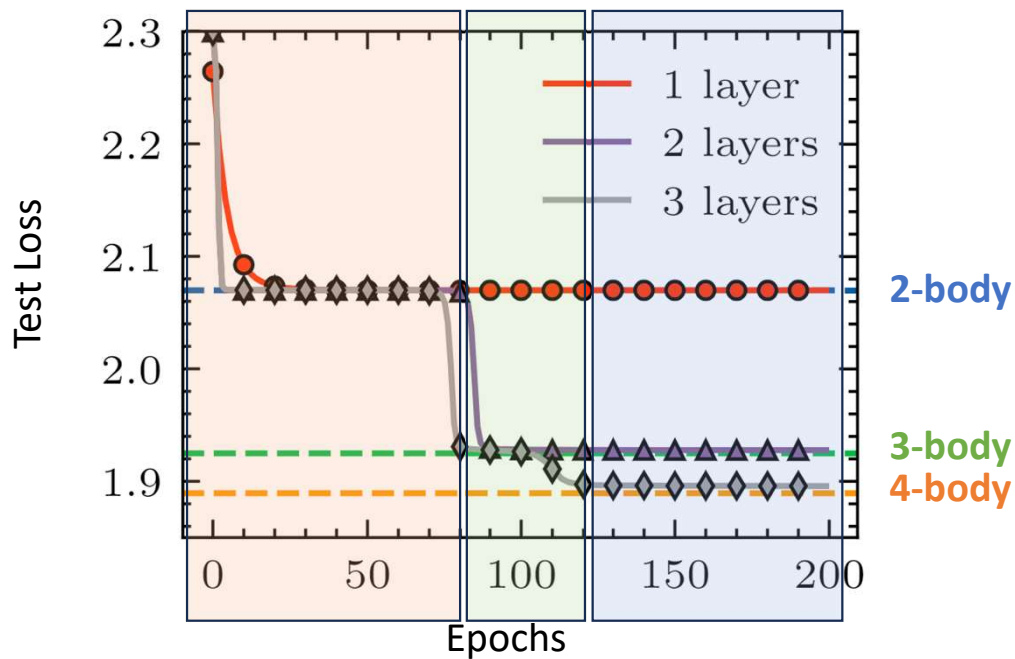


And if the interactions are of higher-order?

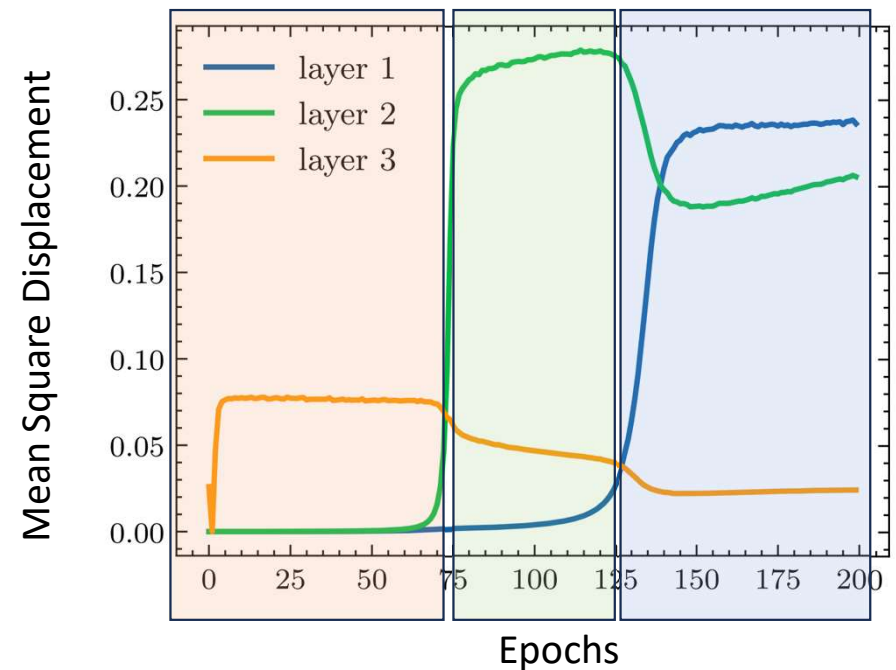


N layers learn interactions of order $N + 1$...

And if the interactions are of higher-order?

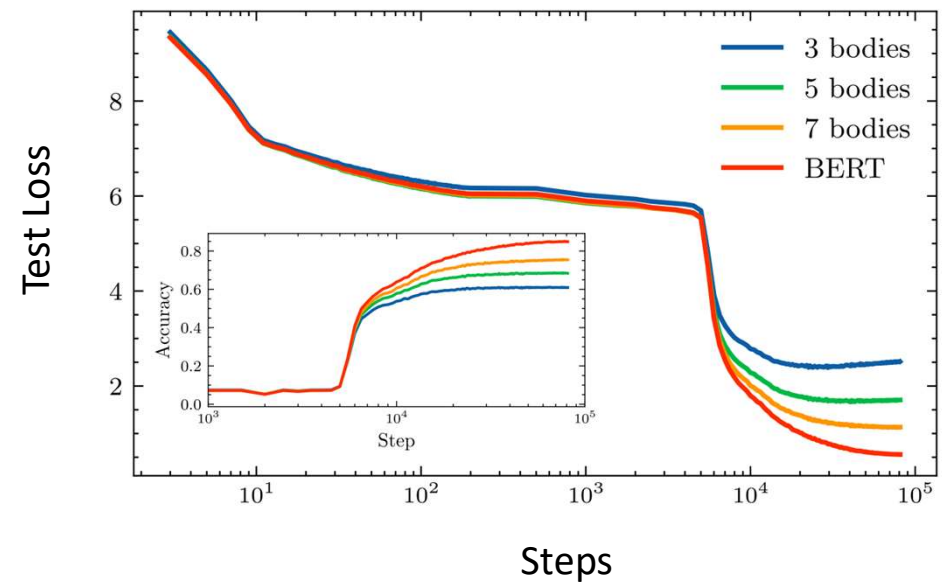
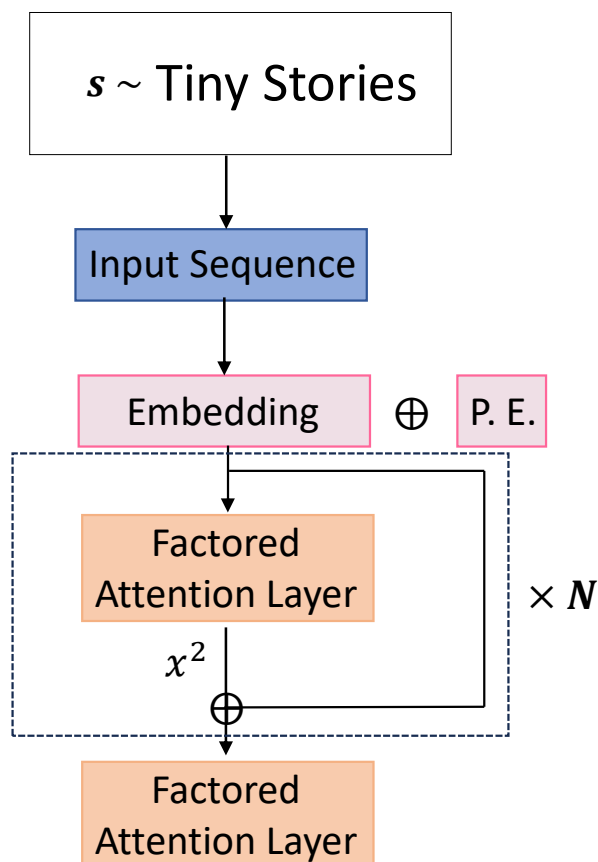


N layers learn interactions of order $N + 1$...



$$MSD_l(t) = \frac{1}{L} \|A_l(t) - A_l(0)\|^2 + \frac{1}{C} \|V_l(t) - V_l(0)\|^2$$

And if the interactions are of higher-order?



...and in real NLP datasets!

Question 2



How many samples are required to achieve good generalization performances?

On the way to fill the gap

Statistical Physics: $x \sim \mathcal{N}(0, \mathbb{I})$



Gardner, Derrida, *Journal of Physics A*, 1988...



Up to now: $x \sim \sum_{k=1}^K p_k \mathcal{N}(\mu_k, \Omega_k)$



Machine Learning: $x \sim ?$ 🤔



Gerace, Loureiro, Mézard, Krzakala, Zdeborová, *ICML*, 2020
Loureiro, Gerbelot, Cui, Goldt, Krzakala, Mézard, Zdeborová, *NeurIPS*, 2021
Loureiro, Sicuro, Gerbelot, Pacco, Krzakala, Zdeborová, *NeurIPS*, 2021
Gerace, Loureiro, Stephan, Krzakala, Zdeborová, *PRE*, 2023
Sarao Mannelli, Gerace, Rostamzadeh, Saglietti, *arXiv:2205.15935*, 2022

A Simplified Gaussian Data Model

$$\mathbf{s} = \{\mathbf{s}_1, \dots, \mathbf{s}_i, \dots, \mathbf{s}_L\}$$



Generalized Potts Model:

$$\mathcal{P}(\mathbf{s}) \propto \exp\left(\frac{\beta}{2} \sum_{i,j} J_{ij} \mathbf{s}_i^T U \mathbf{s}_j\right)$$



$$\mathbf{m} = \{m_1, \dots, m_i, \dots, m_L\}$$

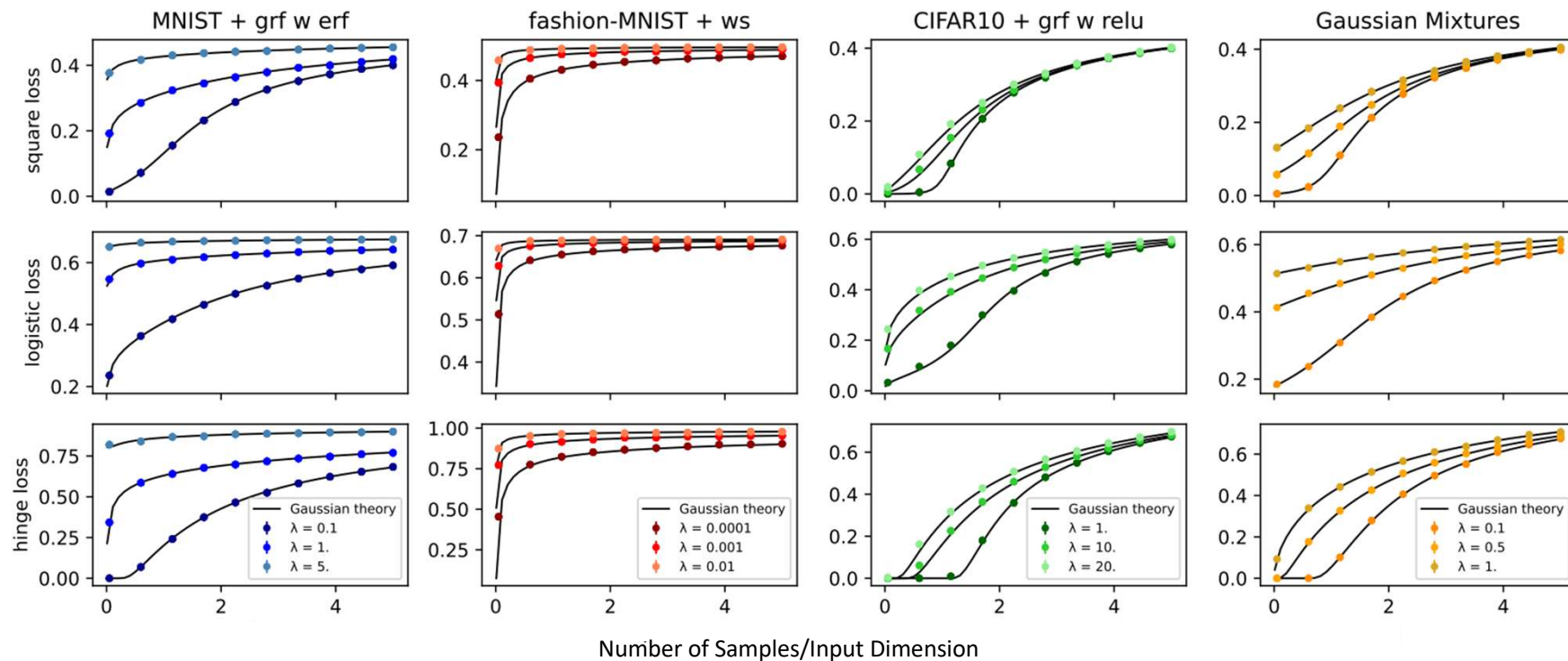


Gaussian Model:

$$\mathcal{N}(\mathbf{m}) \propto \exp\left(-\frac{1}{2} \sum_{i,j} \Sigma_{ij}^{-1} m_i m_j\right)$$

Gaussian data are representative of real data

Single-Layer Neural Networks with Random Labels in classification tasks:



★ Shallow Neural Networks only care of data covariance with random labels!

MLM from a Statistical Physics point of view

Masked Language Modelling:

Data:

$$\mathcal{D} = \left\{ \mathbf{m}_{\setminus i}^{\mu}, m_i^{\mu} \right\}_{\mu=1}^M$$

↓

$$m_i^{\mu} \sim \mathcal{N} \left(m_i^{\mu} | \mathbf{m}_{\setminus i}^{\mu}; \Sigma_i^{-1} \right)$$

Empirical Risk Minimization:

$$\hat{A}_i = \operatorname{argmin}_A \mathbb{E}_{\mathcal{D}} \left[\frac{1}{2} \left(m_i^{\mu} - \sum_{j \neq i=1}^L A_{ij} m_j^{\mu} \right)^2 \right]$$

Test on unseen data:

$$\epsilon_g = \mathbb{E}_{\{x^{new}, y^{new}\}} \left[\frac{1}{2} \sum_{i=1}^L \left(m_i^{new} - \sum_{j \neq i=1}^L \hat{A}_{ij} m_j^{new} \right)^2 \right]$$

Teacher-Student:

Data:

$$\mathcal{D} = \{ \mathbf{x}^{\mu}, y^{\mu} \}_{\mu=1}^M$$

↓

$$y^{\mu} \sim P(y^{\mu} | \mathbf{x}^{\mu}; \mathbf{T})$$

Empirical Risk Minimization:

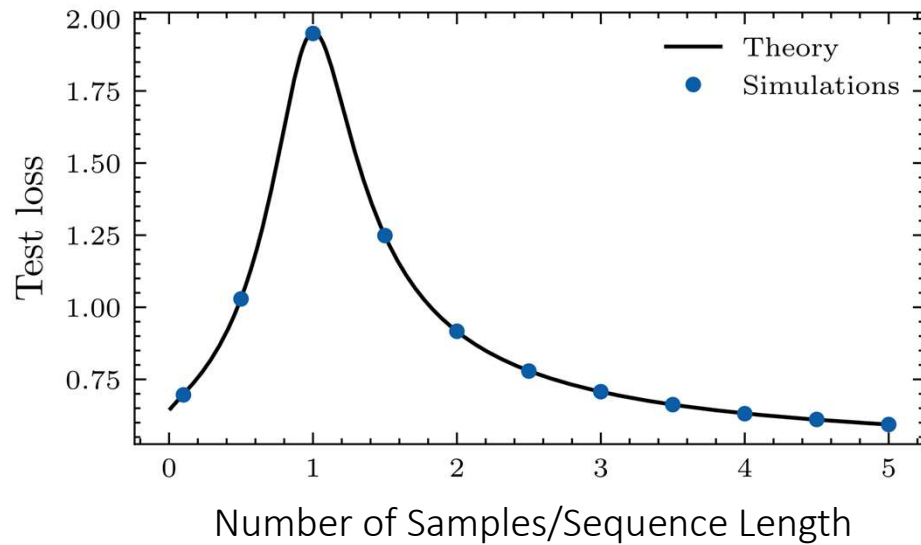
$$\hat{W} = \operatorname{argmin}_W \mathbb{E}_{\mathcal{D}} \left[\frac{1}{2} \left(y^{\mu} - \sum_{j=1}^L W_j x_j^{\mu} \right)^2 \right]$$

Test on unseen data:

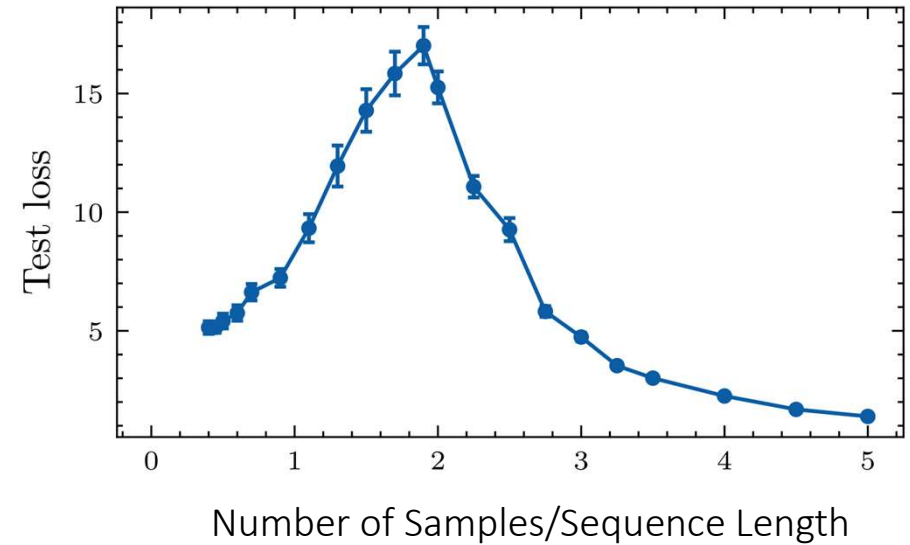
$$\epsilon_g = \mathbb{E}_{\{x^{new}, y^{new}\}} \left[\frac{1}{2} \sum_{i=1}^L \left(y^{new} - \sum_{j=1}^L \hat{W}_j x_j^{new} \right)^2 \right]$$

A New Generalization Behavior

Gaussian Model



Generalized Potts Model



★ Qualitative similar behavior!

To conclude few take home messages...

- Factored Attention learns the exact conditionals of the Generalized Potts Model;
- Deep Transformers sequentially learn high-order interactions in the input data;
- The interpolation peak appears in self-supervised learning too but it is now triggered by the noise inherent in the training data;
- There exist universality classes qualitatively describing the behavior of learning models on more complex data distribution.