# Exploring Transformer-Based Models for Cascading Extremes

### Clemente Ferrer

clemente.ferrer@usm.cl

Department of Mathematics Universidad Técnica Federico Santa María

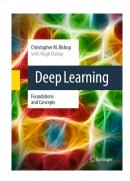
Joint work with Miguel de Carvalho & Ronny Vallejos

Workshop on Generative AI for Extreme Events

Edinburgh, UK June 13, 2025

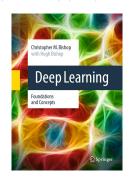
### During this talk I will cover the following topics:

- ► Why use Transformers?
- A Literature Review on Transformers.
- Transformer in Statistics.
- ► Large Language Models and Generative Pre-trained Transformers (GPT).
- Ideas on Transformers for Chained Extremes



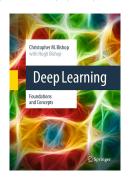
### During this talk I will cover the following topics:

- ▶ Why use Transformers?
- A Literature Review on Transformers.
- Transformer in Statistics.
- ► Large Language Models and Generative Pre-trained Transformers (GPT).
- Ideas on Transformers for Chained Extremes



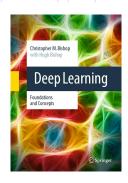
### During this talk I will cover the following topics:

- ▶ Why use Transformers?
- A Literature Review on Transformers.
- Transformer in Statistics.
- ► Large Language Models and Generative Pre-trained Transformers (GPT).
- Ideas on Transformers for Chained Extremes



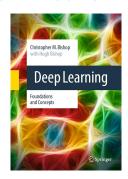
### During this talk I will cover the following topics:

- ▶ Why use Transformers?
- A Literature Review on Transformers.
- Transformer in Statistics.
- ► Large Language Models and Generative Pre-trained Transformers (GPT).
- Ideas on Transformers for Chained Extremes



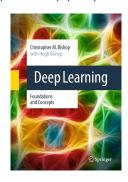
### During this talk I will cover the following topics:

- ▶ Why use Transformers?
- A Literature Review on Transformers.
- Transformer in Statistics.
- ► Large Language Models and Generative Pre-trained Transformers (GPT).
- Ideas on Transformers for Chained Extremes



During this talk I will cover the following topics:

- Why use Transformers?
- A Literature Review on Transformers.
- Transformer in Statistics.
- Large Language Models and Generative Pre-trained Transformers (GPT).
- Ideas on Transformers for Chained Extremes.



## Why use Transformers?

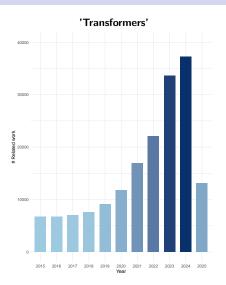


Figure: Number of works (articles, preprints, book chapter) referring to 'Transformers' since 2015 (Source: OpenAlex.org).

### ${\it Suppose}\,\,N\,\,{\it tokens}$

$$\mathbf{x}_1, \dots, \mathbf{x}_N \quad \leftarrow \quad \mathsf{Input data}.$$

where  $\mathbf{x}_n = (x_{n1}, \dots, x_{nD})^\mathrm{T} \in \mathbb{R}^D$ ,  $n = 1, \dots, N$ . The elements of the tokens are called features.

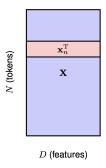


Figure: The structure of the data matrix  $\mathbf{X}$ , of dimension  $N \times D$ , in which row r represents the transposed data vector  $\mathbf{x}_n^T$ .

### ${\bf Suppose}\ N\ {\bf tokens}$

$$\mathbf{x}_1, \dots, \mathbf{x}_N \quad \leftarrow \quad \mathsf{Input data}.$$

where  $\mathbf{x}_n=(x_{n1},\dots,x_{nD})^\mathrm{T}\in\mathbb{R}^D$ ,  $n=1,\dots,N.$  The elements of the tokens are called features.

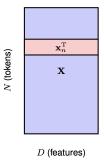


Figure: The structure of the data matrix  $\mathbf{X}$ , of dimension  $N \times D$ , in which row n represents the transposed data vector  $\mathbf{x}_n^T$ .

#### Goal

Define a transformation from tokens  $\mathbf{x}_1, \dots, \mathbf{x}_N$  to output embeddings  $\mathbf{y}_1, \dots, \mathbf{y}_N$ , where each  $\mathbf{y}_n$  is a weighted combination of ALL  $\mathbf{x}_m$ , with weights reflecting the relevance of  $\mathbf{x}_m$  to  $\mathbf{y}_n$ .

Proposal

$$\mathbf{y}_n = \sum_{m=1}^N \alpha_{nm} \mathbf{x}_m,$$

where  $\alpha_{nm} \geq 0$  are attention weights such that  $\sum_{m=1}^{N} \alpha_{nm} = 1$ 

FIRST APPROACH

$$a_{nm} = \frac{\exp(\mathbf{x}_n^{\mathrm{T}}\mathbf{x}_m)}{\sum_{m'=1}^{N} \exp(\mathbf{x}_n^{\mathrm{T}}\mathbf{x}_{m'})} \quad \leftarrow \quad \mathsf{Softmax}$$

$$\mathbf{Y} = \operatorname{Softmax}(\mathbf{X}\mathbf{X}^{\mathrm{T}})\mathbf{X} \quad \leftarrow \quad \mathsf{Self-Attention}$$

#### GOAL

Define a transformation from tokens  $\mathbf{x}_1, \dots, \mathbf{x}_N$  to output embeddings  $\mathbf{y}_1, \dots, \mathbf{y}_N$ , where each  $\mathbf{y}_n$  is a weighted combination of ALL  $\mathbf{x}_m$ , with weights reflecting the relevance of  $\mathbf{x}_m$  to  $\mathbf{y}_n$ .

Proposal

$$\mathbf{y}_n = \sum_{m=1}^N \alpha_{nm} \mathbf{x}_m,$$

where  $\alpha_{nm} \geq 0$  are attention weights such that  $\sum_{m=1}^{N} \alpha_{nm} = 1$ .

FIRST APPROACH

$$a_{nm} = \frac{\exp(\mathbf{x}_n^{\mathrm{T}}\mathbf{x}_m)}{\sum_{m'=1}^{N} \exp(\mathbf{x}_n^{\mathrm{T}}\mathbf{x}_{m'})} \quad \leftarrow \quad \mathsf{Softmax}$$

$$\mathbf{Y} = \operatorname{Softmax}(\mathbf{X}\mathbf{X}^{\mathrm{T}})\mathbf{X} \quad \leftarrow \quad \mathsf{Self-Attention}$$

Goal

Define a transformation from tokens  $\mathbf{x}_1, \dots, \mathbf{x}_N$  to output embeddings  $\mathbf{y}_1, \dots, \mathbf{y}_N$ , where each  $\mathbf{y}_n$  is a weighted combination of ALL  $\mathbf{x}_m$ , with weights reflecting the relevance of  $\mathbf{x}_m$  to  $\mathbf{y}_n$ .

Proposal

$$\mathbf{y}_n = \sum_{m=1}^N \alpha_{nm} \mathbf{x}_m,$$

where  $\alpha_{nm} \geq 0$  are attention weights such that  $\sum_{m=1}^{N} \alpha_{nm} = 1$ .

FIRST APPROACH

$$a_{nm} = \frac{\exp(\mathbf{x}_n^{\mathrm{T}}\mathbf{x}_m)}{\sum_{m'=1}^{N} \exp(\mathbf{x}_n^{\mathrm{T}}\mathbf{x}_{m'})} \quad \leftarrow \quad \mathsf{Softmax}.$$

$$\mathbf{Y} = \operatorname{Softmax}(\mathbf{X}\mathbf{X}^{\mathrm{T}})\mathbf{X} \quad \leftarrow \quad \mathsf{Self-Attention}$$

Goal

Define a transformation from tokens  $\mathbf{x}_1, \dots, \mathbf{x}_N$  to output embeddings  $\mathbf{y}_1, \dots, \mathbf{y}_N$ , where each  $\mathbf{y}_n$  is a weighted combination of ALL  $\mathbf{x}_m$ , with weights reflecting the relevance of  $\mathbf{x}_m$  to  $\mathbf{y}_n$ .

Proposal

$$\mathbf{y}_n = \sum_{m=1}^N \alpha_{nm} \mathbf{x}_m,$$

where  $\alpha_{nm} \geq 0$  are attention weights such that  $\sum_{m=1}^{N} \alpha_{nm} = 1$ .

FIRST APPROACH

$$a_{nm} = \frac{\exp(\mathbf{x}_n^{\mathrm{T}} \mathbf{x}_m)}{\sum_{m'=1}^{N} \exp(\mathbf{x}_n^{\mathrm{T}} \mathbf{x}_{m'})} \quad \leftarrow \quad \mathsf{Softmax}.$$

$$\mathbf{Y} = \operatorname{Softmax}(\mathbf{X}\mathbf{X}^{\mathrm{T}})\mathbf{X} \quad \leftarrow \quad extstyle{\mathsf{Self-Attention}}$$

Goal

Define a transformation from tokens  $\mathbf{x}_1, \dots, \mathbf{x}_N$  to output embeddings  $\mathbf{y}_1, \dots, \mathbf{y}_N$ , where each  $\mathbf{y}_n$  is a weighted combination of ALL  $\mathbf{x}_m$ , with weights reflecting the relevance of  $\mathbf{x}_m$  to  $\mathbf{y}_n$ .

Proposal

$$\mathbf{y}_n = \sum_{m=1}^N \alpha_{nm} \mathbf{x}_m,$$

where  $\alpha_{nm} \geq 0$  are attention weights such that  $\sum_{m=1}^{N} \alpha_{nm} = 1$ .

FIRST APPROACH

$$a_{nm} = \frac{\exp(\mathbf{x}_n^{\mathrm{T}} \mathbf{x}_m)}{\sum_{m'=1}^{N} \exp(\mathbf{x}_n^{\mathrm{T}} \mathbf{x}_{m'})} \leftarrow \text{Softmax}.$$

$$\mathbf{Y} = \operatorname{Softmax}(\mathbf{X}\mathbf{X}^T)\mathbf{X} \quad \leftarrow \quad \text{Self-Attention}.$$

#### LIMITATION

The mapping  $\{\mathbf{x}_n\} \mapsto \{\mathbf{y}_n\}$  lacks learnable parameters and thus cannot adapt to data. Additionally, all features in each  $\mathbf{x}_n$  contribute equally to attention weights, limiting the model's ability to focus selectively on informative features.

PARTIAL SOLUTION

Define

$$\tilde{\mathbf{X}} = \mathbf{X}\mathbf{U}, \quad \mathbf{U} \in \mathbb{R}^{D \times D},$$

where  ${f U}$  is learnable, analogous to a 'layer' in standard neural networks. Then

$$\mathbf{Y} = \operatorname{Softmax}(\mathbf{X}\mathbf{U}\mathbf{U}^{\mathrm{T}}\mathbf{X}^{\mathrm{T}})\mathbf{X}\mathbf{U}.$$

STILL A PROBLEM

Although this has much more flexibility, the matrix

$$\mathbf{X}\mathbf{U}\mathbf{U}^{\mathrm{T}}\mathbf{X}^{\mathrm{T}}$$

#### LIMITATION

The mapping  $\{\mathbf{x}_n\} \mapsto \{\mathbf{y}_n\}$  lacks learnable parameters and thus cannot adapt to data. Additionally, all features in each  $\mathbf{x}_n$  contribute equally to attention weights, limiting the model's ability to focus selectively on informative features.

PARTIAL SOLUTION

Define

$$\tilde{\mathbf{X}} = \mathbf{X}\mathbf{U}, \quad \mathbf{U} \in \mathbb{R}^{D \times D}$$

where  $\mathbf{U}$  is learnable, analogous to a 'layer' in standard neural networks. Then

$$\mathbf{Y} = \operatorname{Softmax}(\mathbf{X}\mathbf{U}\mathbf{U}^{\mathrm{T}}\mathbf{X}^{\mathrm{T}})\mathbf{X}\mathbf{U}.$$

Still a problem

Although this has much more flexibility, the matrix

$$XUU^{T}X^{T}$$

#### LIMITATION

The mapping  $\{\mathbf{x}_n\} \mapsto \{\mathbf{y}_n\}$  lacks learnable parameters and thus cannot adapt to data. Additionally, all features in each  $\mathbf{x}_n$  contribute equally to attention weights, limiting the model's ability to focus selectively on informative features.

#### PARTIAL SOLUTION

Define

$$\tilde{\mathbf{X}} = \mathbf{X}\mathbf{U}, \quad \mathbf{U} \in \mathbb{R}^{D \times D},$$

where  $\mathbf{U}$  is learnable, analogous to a 'layer' in standard neural networks. Then

$$\mathbf{Y} = \operatorname{Softmax}(\mathbf{X}\mathbf{U}\mathbf{U}^{\mathrm{T}}\mathbf{X}^{\mathrm{T}})\mathbf{X}\mathbf{U}.$$

Still a problem

Although this has much more flexibility, the matrix

$$\mathbf{X}\mathbf{U}\mathbf{U}^{\mathrm{T}}\mathbf{X}^{\mathrm{T}}$$

#### LIMITATION

The mapping  $\{\mathbf{x}_n\} \mapsto \{\mathbf{y}_n\}$  lacks learnable parameters and thus cannot adapt to data. Additionally, all features in each  $\mathbf{x}_n$  contribute equally to attention weights, limiting the model's ability to focus selectively on informative features.

#### Partial Solution

Define

$$\tilde{\mathbf{X}} = \mathbf{X}\mathbf{U}, \quad \mathbf{U} \in \mathbb{R}^{D \times D},$$

where  $\mathbf{U}$  is learnable, analogous to a 'layer' in standard neural networks. Then

$$\mathbf{Y} = \operatorname{Softmax}(\mathbf{X}\mathbf{U}\mathbf{U}^{\mathrm{T}}\mathbf{X}^{\mathrm{T}})\mathbf{X}\mathbf{U}.$$

#### STILL A PROBLEM

Although this has much more flexibility, the matrix

$$\mathbf{X}\mathbf{U}\mathbf{U}^{\mathrm{T}}\mathbf{X}^{\mathrm{T}}$$

#### SOLUTION

Vaswani et al. (2017) proposed defining separate query, key, and value matrices each having their own independent linear transformations:

$$\mathbf{Q} = \mathbf{X}\mathbf{W}^{(q)}, \quad \mathbf{K} = \mathbf{X}\mathbf{W}^{(k)}, \quad \mathbf{V} = \mathbf{X}\mathbf{W}^{(v)},$$

where  $\mathbf{W}^{(\mathrm{q})}$ ,  $\mathbf{W}^{(\mathrm{k})} \in \mathbb{R}^{D \times D_{\mathrm{k}}}$  and  $\mathbf{W}^{(\mathrm{v})} \in \mathbb{R}^{D \times D_{\mathrm{v}}}$ , where  $D_{\mathrm{v}}$  governs the dimensionality of the output. Therefore

 $\mathbf{Y} = \operatorname{Softmax}(\mathbf{Q}\mathbf{K}^{T})\mathbf{V} \leftarrow \mathbf{Q}\mathbf{K}^{T}$  are the attention weights.

$$\begin{array}{|c|c|}
\hline
\mathbf{Y} & = & \text{Softmax} & \left\{ \begin{array}{|c|c|}
\hline
\mathbf{Q}\mathbf{K}^{\mathrm{T}} \\
\end{array} \right\} \times \begin{bmatrix} \mathbf{V} \\
\end{array}$$

$$N \times D_{\mathbf{v}} \qquad \qquad N \times D_{\mathbf{v}} \\$$

Figure: Illustration of the evaluation of the output from an attention layer given the query, key, and value matrices Q, K, and V, respectively.

#### SOLUTION

Vaswani et al. (2017) proposed defining separate query, key, and value matrices each having their own independent linear transformations:

$$\mathbf{Q} = \mathbf{X}\mathbf{W}^{(q)}, \quad \mathbf{K} = \mathbf{X}\mathbf{W}^{(k)}, \quad \mathbf{V} = \mathbf{X}\mathbf{W}^{(v)},$$

where  $\mathbf{W}^{(q)}$ ,  $\mathbf{W}^{(k)} \in \mathbb{R}^{D \times D_k}$  and  $\mathbf{W}^{(v)} \in \mathbb{R}^{D \times D_v}$ , where  $D_v$  governs the dimensionality of the output. Therefore

$$\mathbf{Y} = \operatorname{Softmax}(\mathbf{Q}\mathbf{K}^{T})\mathbf{V} \leftarrow \mathbf{Q}\mathbf{K}^{T}$$
 are the attention weights.

$$\begin{array}{|c|c|} \hline \mathbf{Y} & = & \mathrm{Softmax} & \left\{ \begin{array}{|c|c|} \hline \mathbf{Q}\mathbf{K}^{\mathrm{T}} \end{array} \right\} \times \begin{bmatrix} \mathbf{V} \\ \mathbf{V} \end{array} \\
N \times D_{\mathbf{V}} & N \times N & N \times D_{\mathbf{V}} \end{array}$$

Figure: Illustration of the evaluation of the output from an attention layer given the query, key, and value matrices Q, K, and V, respectively.

#### Final adjustment

To avoid vanishing gradients in softmax, scale the dot product of query and key vectors by  $\sqrt{D_{\mathbf{k}}}$ , yielding the output of the attention layer as:

$$\mathbf{Y} = \operatorname{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) \equiv \operatorname{Softmax}\left(rac{\mathbf{Q}\mathbf{K}^{\mathrm{T}}}{\sqrt{D_{\mathrm{k}}}}
ight)\mathbf{V}$$

Scaled dot-product self-attention

This structure constitutes a single attention head.

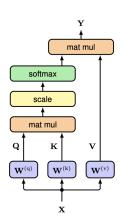


Figure: Information flow. Here 'mat mul' denotes matrix multiplication, and 'scale' refers to the normalization of the argument to the softmax using  $\sqrt{D_{\mathbf{k}}}$ .

#### Modeling multiple dependency patterns

To capture different types of relationships in the data simultaneously we replicate the attention mechanism into multiple heads.

Suppose we have H attention heads indexed by h = 1, ..., H of the form

$$\mathbf{H}_h = \operatorname{Attention}(\mathbf{Q}_h, \mathbf{K}_h, \mathbf{V}_h), \quad \mathbf{Q}_h = \mathbf{X} \mathbf{W}_h^{(q)}, \quad \mathbf{K}_h = \mathbf{X} \mathbf{W}_h^{(k)}, \quad \mathbf{V}_h = \mathbf{X} \mathbf{W}_h^{(v)}.$$

The heads are first concatenated into a single matrix, and the result is then linearly transformed using a matrix  $\mathbf{W}^{(\circ)}$  to give a combined output in the form

$$\mathbf{Y}(\mathbf{X}) = \mathsf{Concat}(\mathbf{H}_1, \dots, \mathbf{H}_H) \mathbf{W}^{(\circ)}, \in \mathbb{R}^{N imes D_{\mathrm{v}}}.$$

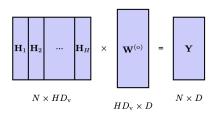


Figure: Network architecture for multi-head attention

#### Modeling multiple dependency patterns

To capture different types of relationships in the data simultaneously we replicate the attention mechanism into multiple heads.

Suppose we have H attention heads indexed by  $h = 1, \ldots, H$  of the form

$$\mathbf{H}_h = \operatorname{Attention}(\mathbf{Q}_h, \mathbf{K}_h, \mathbf{V}_h), \quad \mathbf{Q}_h = \mathbf{X} \mathbf{W}_h^{(q)}, \quad \mathbf{K}_h = \mathbf{X} \mathbf{W}_h^{(k)}, \quad \mathbf{V}_h = \mathbf{X} \mathbf{W}_h^{(v)}.$$

The heads are first concatenated into a single matrix, and the result is then linearly transformed using a matrix  $\mathbf{W}^{(\circ)}$  to give a combined output in the form

$$\mathbf{Y}(\mathbf{X}) = \mathsf{Concat}(\mathbf{H}_1, \dots, \mathbf{H}_H) \mathbf{W}^{(\circ)}, \in \mathbb{R}^{N \times D_{\mathrm{v}}}.$$

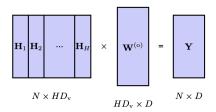


Figure: Network architecture for multi-head attention

#### Modeling multiple dependency patterns

To capture different types of relationships in the data simultaneously we replicate the attention mechanism into multiple heads.

Suppose we have H attention heads indexed by  $h = 1, \dots, H$  of the form

$$\mathbf{H}_h = \operatorname{Attention}(\mathbf{Q}_h, \mathbf{K}_h, \mathbf{V}_h), \quad \mathbf{Q}_h = \mathbf{X}\mathbf{W}_h^{(q)}, \quad \mathbf{K}_h = \mathbf{X}\mathbf{W}_h^{(k)}, \quad \mathbf{V}_h = \mathbf{X}\mathbf{W}_h^{(v)}.$$

The heads are first concatenated into a single matrix, and the result is then linearly transformed using a matrix  $\mathbf{W}^{(\circ)}$  to give a combined output in the form

$$\mathbf{Y}(\mathbf{X}) = \mathsf{Concat}(\mathbf{H}_1, \dots, \mathbf{H}_H) \mathbf{W}^{(\circ)}, \in \mathbb{R}^{N \times D_{\mathbf{v}}}.$$

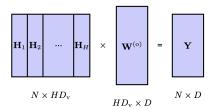


Figure: Network architecture for multi-head attention

#### Modeling multiple dependency patterns

To capture different types of relationships in the data simultaneously we replicate the attention mechanism into multiple heads.

Suppose we have H attention heads indexed by  $h = 1, \dots, H$  of the form

$$\mathbf{H}_h = \operatorname{Attention}(\mathbf{Q}_h, \mathbf{K}_h, \mathbf{V}_h), \quad \mathbf{Q}_h = \mathbf{X}\mathbf{W}_h^{(q)}, \quad \mathbf{K}_h = \mathbf{X}\mathbf{W}_h^{(k)}, \quad \mathbf{V}_h = \mathbf{X}\mathbf{W}_h^{(v)}.$$

The heads are first concatenated into a single matrix, and the result is then linearly transformed using a matrix  $\mathbf{W}^{(\circ)}$  to give a combined output in the form

$$\mathbf{Y}(\mathbf{X}) = \mathsf{Concat}(\mathbf{H}_1, \dots, \mathbf{H}_H) \mathbf{W}^{(\circ)}, \in \mathbb{R}^{N \times D_{\mathbf{V}}}.$$

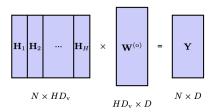


Figure: Network architecture for multi-head attention.

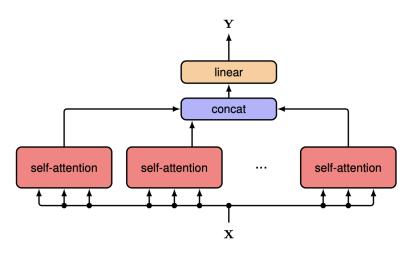


Figure: Information flow in a multi-head attention layer.

#### TRAINING IMPROVEMENTS

To improve training efficiency, the output is followed by a *residual connection* and a *layer normalization* (Ba, Kiros, and Hinton, 2016)

$$\mathbf{Z} = \mathsf{LayerNorm}(\mathbf{Y}(\mathbf{X}) + \mathbf{X}).$$

#### Still work to do

An MLP is added after the attention layer to break the linearity of the output and increase the expressive capabilities of the attention layer

$$\widetilde{\mathbf{X}} = \mathsf{LayerNorm}(\mathrm{MLP}(\mathbf{Z}) + \mathbf{Z}).$$

$$\widetilde{\mathbf{X}} = \mathrm{TF}_{\boldsymbol{\theta}}(\mathbf{X})$$

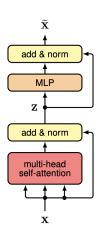


Figure: One layer of the transformer architecture. Here, 'add and norm' denotes a residual connection followed by layer normalization.

#### Training improvements

To improve training efficiency, the output is followed by a *residual connection* and a *layer normalization* (Ba, Kiros, and Hinton, 2016)

$$\mathbf{Z} = \mathsf{LayerNorm}(\mathbf{Y}(\mathbf{X}) + \mathbf{X}).$$

#### STILL WORK TO DO

An MLP is added after the attention layer to break the linearity of the output and increase the expressive capabilities of the attention layer

$$\widetilde{\mathbf{X}} = \mathsf{LayerNorm}(\mathrm{MLP}(\mathbf{Z}) + \mathbf{Z}).$$

$$\widetilde{\mathbf{X}} = \mathrm{TF}_{\boldsymbol{\theta}}(\mathbf{X})$$

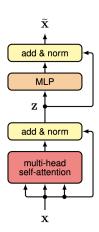


Figure: One layer of the transformer architecture. Here, 'add and norm' denotes a residual connection followed by layer normalization.

#### TRAINING IMPROVEMENTS

To improve training efficiency, the output is followed by a *residual connection* and a *layer normalization* (Ba, Kiros, and Hinton, 2016)

$$\mathbf{Z} = \mathsf{LayerNorm}(\mathbf{Y}(\mathbf{X}) + \mathbf{X}).$$

#### Still work to do

An MLP is added after the attention layer to break the linearity of the output and increase the expressive capabilities of the attention layer:

$$\widetilde{\mathbf{X}} = \mathsf{LayerNorm}(\mathsf{MLP}(\mathbf{Z}) + \mathbf{Z}).$$

$$\widetilde{\mathbf{X}} = \mathrm{TF}_{\boldsymbol{\theta}}(\mathbf{X})$$

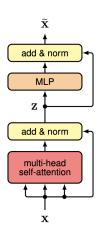


Figure: One layer of the transformer architecture. Here, 'add and norm' denotes a residual connection followed by layer normalization.

#### Training improvements

To improve training efficiency, the output is followed by a *residual connection* and a *layer normalization* (Ba, Kiros, and Hinton, 2016)

$$\mathbf{Z} = \mathsf{LayerNorm}(\mathbf{Y}(\mathbf{X}) + \mathbf{X}).$$

#### Still work to do

An MLP is added after the attention layer to break the linearity of the output and increase the expressive capabilities of the attention layer:

$$\widetilde{\mathbf{X}} = \mathsf{LayerNorm}(\mathrm{MLP}(\mathbf{Z}) + \mathbf{Z}).$$

$$\widetilde{\mathbf{X}} = \mathrm{TF}_{\boldsymbol{\theta}}(\mathbf{X})$$

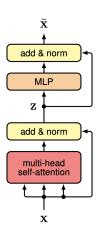


Figure: One layer of the transformer architecture. Here, 'add and norm' denotes a residual connection followed by layer normalization.

#### Training improvements

To improve training efficiency, the output is followed by a *residual connection* and a *layer normalization* (Ba, Kiros, and Hinton, 2016)

$$\mathbf{Z} = \mathsf{LayerNorm}(\mathbf{Y}(\mathbf{X}) + \mathbf{X}).$$

#### Still work to do

An MLP is added after the attention layer to break the linearity of the output and increase the expressive capabilities of the attention layer:

$$\widetilde{\mathbf{X}} = \mathsf{LayerNorm}(\mathsf{MLP}(\mathbf{Z}) + \mathbf{Z}).$$

$$\widetilde{\mathbf{X}} = \mathrm{TF}_{\boldsymbol{\theta}}(\mathbf{X})$$

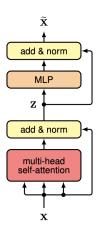


Figure: One layer of the transformer architecture. Here, 'add and norm' denotes a residual connection followed by layer normalization.

Garg et al. (2022) show empirically that standard Transformers can be trained from scratch to perform in-context learning of linear functions.

#### NOTATION

Let  $\mathcal X$  be the input feature space, and  $\mathcal Y$  be the output/label space, such that  $\mathcal F=\{f:\mathcal X\to\mathcal Y\}$  is class of functions between this spaces.

#### IN-CONTEXT LEARNING SETTING

- 1) Sample  $f \sim \mathsf{P}_{\mathcal{F}}$  and N i.i.d. inputs tokens  $\mathbf{x}_1, \dots, \mathbf{x}_N \sim \mathsf{P}_{\mathcal{X}}$ .
- **2)** Generate labels  $\mathbf{y}_i = f(\mathbf{x}_i) \in \mathcal{Y}$  for  $i=1,\ldots,N$ , obtaining  $\{(\mathbf{x}_i,\mathbf{y}_i)\}_{i=1}^N$
- 3) Define a length-i prompt as

$$\mathbf{X}_{\mathsf{prompt}}^{(i)} = (\mathbf{x}_1, \mathbf{y}_1, \dots, \mathbf{x}_{i-1}, \mathbf{y}_{i-1}, \mathbf{x}_i)^{\mathrm{T}},$$

such that the model prediction is  $\hat{\mathbf{y}}_i = \mathbf{W} \operatorname{TF}_{\boldsymbol{\theta}}(\mathbf{X}_{\mathsf{prompt}}^{(i)}) \in \mathcal{Y}, \ i = 1, \dots, N.$ 

4) Training by minimizing the expected loss

$$\min_{\boldsymbol{\theta}} \underset{\mathbf{x}_1, \dots, \mathbf{x}_n \sim \mathsf{P}_{\mathcal{X}}}{\mathbb{E}} \left[ \frac{1}{N} \sum_{i=1}^{N} \ell(\hat{\mathbf{y}}_i, \mathbf{y}_i) \right],$$

Garg et al. (2022) show empirically that standard Transformers can be trained from scratch to perform in-context learning of linear functions.

#### NOTATION

Let  $\mathcal X$  be the input feature space, and  $\mathcal Y$  be the output/label space, such that  $\mathcal F=\{f:\mathcal X\to\mathcal Y\}$  is class of functions between this spaces.

#### IN-CONTEXT LEARNING SETTING

- 1) Sample  $f \sim \mathsf{P}_{\mathcal{F}}$  and N i.i.d. inputs tokens  $\mathbf{x}_1, \dots, \mathbf{x}_N \sim \mathsf{P}_{\mathcal{X}}$ .
- **2)** Generate labels  $\mathbf{y}_i = f(\mathbf{x}_i) \in \mathcal{Y}$  for  $i=1,\ldots,N$ , obtaining  $\{(\mathbf{x}_i,\mathbf{y}_i)\}_{i=1}^N$
- 3) Define a length-i prompt as

$$\mathbf{X}_{\mathsf{prompt}}^{(i)} = (\mathbf{x}_1, \mathbf{y}_1, \dots, \mathbf{x}_{i-1}, \mathbf{y}_{i-1}, \mathbf{x}_i)^{\mathrm{T}},$$

such that the model prediction is  $\hat{\mathbf{y}}_i = \mathbf{W} \operatorname{TF}_{\boldsymbol{\theta}}(\mathbf{X}_{\mathsf{prompt}}^{(i)}) \in \mathcal{Y}, \ i = 1, \dots, N.$ 

4) Training by minimizing the expected loss

$$\min_{\boldsymbol{\theta}} \underset{\mathbf{x}_1, \dots, \mathbf{x}_n \sim \mathsf{P}_{\mathcal{X}}}{\mathbb{E}} \left[ \frac{1}{N} \sum_{i=1}^{N} \ell(\hat{\mathbf{y}}_i, \mathbf{y}_i) \right],$$

Garg et al. (2022) show empirically that standard Transformers can be trained from scratch to perform in-context learning of linear functions.

#### NOTATION

Let  $\mathcal X$  be the input feature space, and  $\mathcal Y$  be the output/label space, such that  $\mathcal F=\{f:\mathcal X\to\mathcal Y\}$  is class of functions between this spaces.

#### IN-CONTEXT LEARNING SETTING

- 1) Sample  $f \sim \mathsf{P}_{\mathcal{F}}$  and N i.i.d. inputs tokens  $\mathbf{x}_1, \dots, \mathbf{x}_N \sim \mathsf{P}_{\mathcal{X}}$ .
- 2) Generate labels  $\mathbf{y}_i = f(\mathbf{x}_i) \in \mathcal{Y}$  for i = 1, ..., N, obtaining  $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^N$
- 3) Define a length-i prompt as

$$\mathbf{X}_{\mathsf{prompt}}^{(i)} = (\mathbf{x}_1, \mathbf{y}_1, \dots, \mathbf{x}_{i-1}, \mathbf{y}_{i-1}, \mathbf{x}_i)^{\mathrm{T}},$$

such that the model prediction is  $\hat{\mathbf{y}}_i = \mathbf{W} \operatorname{TF}_{\theta}(\mathbf{X}_{\mathsf{prompt}}^{(i)}) \in \mathcal{Y}, i = 1, \dots, N.$ 

4) Training by minimizing the expected loss

$$\min_{\boldsymbol{\theta}} \underset{\mathbf{x}_1, \dots, \mathbf{x}_n \sim \mathsf{P}_{\mathcal{X}}}{\mathbb{E}} \left[ \frac{1}{N} \sum_{i=1}^{N} \ell(\hat{\mathbf{y}}_i, \mathbf{y}_i) \right]$$

Garg et al. (2022) show empirically that standard Transformers can be trained from scratch to perform in-context learning of linear functions.

#### NOTATION

Let  $\mathcal X$  be the input feature space, and  $\mathcal Y$  be the output/label space, such that  $\mathcal F=\{f:\mathcal X\to\mathcal Y\}$  is class of functions between this spaces.

#### IN-CONTEXT LEARNING SETTING

- 1) Sample  $f \sim P_{\mathcal{F}}$  and N i.i.d. inputs tokens  $\mathbf{x}_1, \dots, \mathbf{x}_N \sim P_{\mathcal{X}}$ .
- 2) Generate labels  $\mathbf{y}_i = f(\mathbf{x}_i) \in \mathcal{Y}$  for i = 1, ..., N, obtaining  $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^N$
- 3) Define a length-i prompt as

$$\mathbf{X}_{\mathsf{prompt}}^{(i)} = (\mathbf{x}_1, \mathbf{y}_1, \dots, \mathbf{x}_{i-1}, \mathbf{y}_{i-1}, \mathbf{x}_i)^{\mathrm{T}},$$

such that the model prediction is  $\hat{\mathbf{y}}_i = \mathbf{W} \operatorname{TF}_{\boldsymbol{\theta}}(\mathbf{X}_{\mathsf{prompt}}^{(i)}) \in \mathcal{Y}, \ i = 1, \dots, N$ 

4) Training by minimizing the expected loss

$$\min_{\boldsymbol{\theta}} \underset{\mathbf{x}_1, \dots, \mathbf{x}_n \sim \mathsf{P}_{\mathcal{X}}}{\mathbb{E}} \left[ \frac{1}{N} \sum_{i=1}^{N} \ell(\hat{\mathbf{y}}_i, \mathbf{y}_i) \right]$$

Garg et al. (2022) show empirically that standard Transformers can be trained from scratch to perform in-context learning of linear functions.

#### NOTATION

Let  $\mathcal X$  be the input feature space, and  $\mathcal Y$  be the output/label space, such that  $\mathcal F=\{f:\mathcal X\to\mathcal Y\}$  is class of functions between this spaces.

#### IN-CONTEXT LEARNING SETTING

- 1) Sample  $f \sim P_{\mathcal{F}}$  and N i.i.d. inputs tokens  $\mathbf{x}_1, \dots, \mathbf{x}_N \sim P_{\mathcal{X}}$ .
- 2) Generate labels  $\mathbf{y}_i = f(\mathbf{x}_i) \in \mathcal{Y}$  for i = 1, ..., N, obtaining  $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^N$ .
- 3) Define a length-i prompt as

$$\mathbf{X}_{\mathsf{prompt}}^{(i)} = (\mathbf{x}_1, \mathbf{y}_1, \dots, \mathbf{x}_{i-1}, \mathbf{y}_{i-1}, \mathbf{x}_i)^{\mathrm{T}},$$

such that the model prediction is  $\hat{\mathbf{y}}_i = \mathbf{W} \operatorname{TF}_{\boldsymbol{\theta}}(\mathbf{X}_{\mathsf{prompt}}^{(i)}) \in \mathcal{Y}, \ i = 1, \dots, N$ 

4) Training by minimizing the expected loss

$$\min_{\boldsymbol{\theta}} \underset{\mathbf{x}_{1},...,\mathbf{x}_{n} \sim \mathsf{P}_{\mathcal{X}}}{\mathbb{E}} \left[ \frac{1}{N} \sum_{i=1}^{N} \ell(\hat{\mathbf{y}}_{i}, \mathbf{y}_{i}) \right]$$

where  $\ell$  is an appropriate loss function.

Garg et al. (2022) show empirically that standard Transformers can be trained from scratch to perform in-context learning of linear functions.

#### NOTATION

Let  $\mathcal X$  be the input feature space, and  $\mathcal Y$  be the output/label space, such that  $\mathcal F=\{f:\mathcal X\to\mathcal Y\}$  is class of functions between this spaces.

#### IN-CONTEXT LEARNING SETTING

- 1) Sample  $f \sim P_{\mathcal{F}}$  and N i.i.d. inputs tokens  $\mathbf{x}_1, \dots, \mathbf{x}_N \sim P_{\mathcal{X}}$ .
- 2) Generate labels  $\mathbf{y}_i = f(\mathbf{x}_i) \in \mathcal{Y}$  for i = 1, ..., N, obtaining  $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^N$ .
- 3) Define a length-i prompt as

$$\mathbf{X}_{\mathsf{prompt}}^{(i)} = (\mathbf{x}_1, \mathbf{y}_1, \dots, \mathbf{x}_{i-1}, \mathbf{y}_{i-1}, \mathbf{x}_i)^{\mathrm{T}},$$

such that the model prediction is  $\hat{\mathbf{y}}_i = \mathbf{W} \operatorname{TF}_{\boldsymbol{\theta}}(\mathbf{X}_{\mathsf{prompt}}^{(i)}) \in \mathcal{Y}, \ i = 1, \dots, N.$ 

4) Training by minimizing the expected loss

$$\min_{\boldsymbol{\theta}} \underset{\substack{\mathbf{x}_1, \dots, \mathbf{x}_n \sim \mathsf{P}_{\mathcal{X}} \\ f \sim \mathsf{P}_{\mathcal{T}}}}{\mathbb{E}} \left[ \frac{1}{N} \sum_{i=1}^{N} \ell(\hat{\mathbf{y}}_i, \mathbf{y}_i) \right]$$

where  $\ell$  is an appropriate loss function.

Garg et al. (2022) show empirically that standard Transformers can be trained from scratch to perform in-context learning of linear functions.

#### NOTATION

Let  $\mathcal X$  be the input feature space, and  $\mathcal Y$  be the output/label space, such that  $\mathcal F=\{f:\mathcal X\to\mathcal Y\}$  is class of functions between this spaces.

#### IN-CONTEXT LEARNING SETTING

- 1) Sample  $f \sim P_{\mathcal{F}}$  and N i.i.d. inputs tokens  $\mathbf{x}_1, \dots, \mathbf{x}_N \sim P_{\mathcal{X}}$ .
- 2) Generate labels  $\mathbf{y}_i = f(\mathbf{x}_i) \in \mathcal{Y}$  for  $i = 1, \dots, N$ , obtaining  $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^N$ .
- 3) Define a length-i prompt as

$$\mathbf{X}_{\mathsf{prompt}}^{(i)} = (\mathbf{x}_1, \mathbf{y}_1, \dots, \mathbf{x}_{i-1}, \mathbf{y}_{i-1}, \mathbf{x}_i)^T,$$

such that the model prediction is  $\hat{\mathbf{y}}_i = \mathbf{W} \operatorname{TF}_{\boldsymbol{\theta}}(\mathbf{X}_{\mathsf{prompt}}^{(i)}) \in \mathcal{Y}, \ i = 1, \dots, N.$ 

4) Training by minimizing the expected loss

$$\min_{\boldsymbol{\theta}} \ \underset{\boldsymbol{\mathbf{x}}_{1}, \dots, \boldsymbol{\mathbf{x}}_{n} \sim \mathsf{P}_{\mathcal{X}}}{\mathbb{E}} \left[ \frac{1}{N} \sum_{i=1}^{N} \ell(\hat{\mathbf{y}}_{i}, \mathbf{y}_{i}) \right],$$

where  $\ell$  is an appropriate loss function.

## Example (Linear functions)

Consider the class of linear functions

$$\mathcal{F} = \{ f : \mathbb{R}^d \to \mathbb{R} \mid f(\mathbf{x}) = \mathbf{w}^T \mathbf{x}, \, \mathbf{w} \in \mathbb{R}^d \}$$

Sample

$$\mathbf{w} \sim \mathsf{P}_{\mathcal{F}} := \mathsf{N}_d(\mathbf{0}, \mathbf{I}_d), \quad \mathbf{x}_i \overset{\mathsf{i.i.d.}}{\sim} \mathsf{P}_{\mathcal{X}} := \mathsf{N}_d(\mathbf{0}, \mathbf{I}_d), \quad i = 1, \dots, N,$$

independently. For each f determined by  $\mathbf{w}$ , generate training pairs

$$\mathbf{v}_i = f(\mathbf{x}_i) = \mathbf{w}^{\mathrm{T}} \mathbf{x}_i.$$

Then, train the Transformer to predict  $\hat{\mathbf{y}}_i = \mathbf{W} \, \mathrm{TF}_{m{ heta}}(\mathbf{X}_{\mathsf{prompt}}^{(i)})$  where

$$\mathbf{X}_{\mathsf{prompt}}^{(i)} = (\mathbf{x}_1, \mathbf{y}_1, \dots, \mathbf{x}_{i-1}, \mathbf{y}_{i-1}, \mathbf{x}_i),$$

$$\min_{\boldsymbol{\theta}} \mathbb{E}_{\substack{\mathbf{x}_1, \dots, \mathbf{x}_N \sim \mathbf{N}_d(\mathbf{0}, \mathbf{I}_d) \\ f \sim \mathbf{N}_d(\mathbf{0}, \mathbf{I}_d)}} \left[ \frac{1}{N} \sum_{i=1}^{N} (\hat{\mathbf{y}}_i - \mathbf{y}_i)^2 \right]$$

## Example (Linear functions)

Consider the class of linear functions

$$\mathcal{F} = \{ f : \mathbb{R}^d \to \mathbb{R} \mid f(\mathbf{x}) = \mathbf{w}^{\mathrm{T}} \mathbf{x}, \, \mathbf{w} \in \mathbb{R}^d \}.$$

Sample

$$\mathbf{w} \sim \mathsf{P}_{\mathcal{F}} := \mathsf{N}_d(\mathbf{0}, \mathbf{I}_d), \quad \mathbf{x}_i \overset{\mathsf{i.i.d.}}{\sim} \mathsf{P}_{\mathcal{X}} := \mathsf{N}_d(\mathbf{0}, \mathbf{I}_d), \quad i = 1, \dots, N,$$

independently. For each f determined by  $\mathbf{w}$ , generate training pairs

$$\mathbf{y}_i = f(\mathbf{x}_i) = \mathbf{w}^{\mathrm{T}} \mathbf{x}_i$$

Then, train the Transformer to predict  $\hat{\mathbf{y}}_i = \mathbf{W} \operatorname{TF}_{\boldsymbol{\theta}}(\mathbf{X}_{\mathsf{prompt}}^{(i)})$  where

$$\mathbf{X}_{\mathsf{prompt}}^{(i)} = (\mathbf{x}_1, \mathbf{y}_1, \dots, \mathbf{x}_{i-1}, \mathbf{y}_{i-1}, \mathbf{x}_i)$$

$$\min_{\boldsymbol{\theta}} \underset{\mathbf{x}_1, \dots, \mathbf{x}_N \sim \mathbf{N}_d(\mathbf{0}, \mathbf{I}_d)}{\mathbb{E}} \left[ \frac{1}{N} \sum_{i=1}^{N} (\hat{\mathbf{y}}_i - \mathbf{y}_i)^2 \right]$$

Example (Linear functions)

Consider the class of linear functions

$$\mathcal{F} = \{ f : \mathbb{R}^d \to \mathbb{R} \mid f(\mathbf{x}) = \mathbf{w}^{\mathrm{T}} \mathbf{x}, \, \mathbf{w} \in \mathbb{R}^d \}.$$

Sample

$$\mathbf{w} \sim \mathsf{P}_{\mathcal{F}} := \mathsf{N}_d(\mathbf{0}, \mathbf{I}_d), \quad \mathbf{x}_i \overset{\mathsf{i.i.d.}}{\sim} \mathsf{P}_{\mathcal{X}} := \mathsf{N}_d(\mathbf{0}, \mathbf{I}_d), \quad i = 1, \dots, N,$$

independently. For each f determined by w, generate training pairs

$$\mathbf{y}_i = f(\mathbf{x}_i) = \mathbf{w}^{\mathrm{T}} \mathbf{x}_i$$

Then, train the Transformer to predict  $\hat{\mathbf{y}}_i = \mathbf{W} \, ext{TF}_{m{ heta}}(\mathbf{X}_{\mathsf{prompt}}^{(i)})$  where

$$\mathbf{X}_{\mathsf{prompt}}^{(i)} = (\mathbf{x}_1, \mathbf{y}_1, \dots, \mathbf{x}_{i-1}, \mathbf{y}_{i-1}, \mathbf{x}_i)$$

$$\min_{\boldsymbol{\theta}} \mathbb{E}_{\substack{\mathbf{x}_1, \dots, \mathbf{x}_N \sim \mathbf{N}_d(\mathbf{0}, \mathbf{I}_d) \\ f \sim \mathbf{N}_d(\mathbf{0}, \mathbf{I}_d)}} \left[ \frac{1}{N} \sum_{i=1}^{N} (\hat{\mathbf{y}}_i - \mathbf{y}_i)^2 \right]$$

Example (Linear functions)

Consider the class of linear functions

$$\mathcal{F} = \{ f : \mathbb{R}^d \to \mathbb{R} \mid f(\mathbf{x}) = \mathbf{w}^{\mathrm{T}} \mathbf{x}, \, \mathbf{w} \in \mathbb{R}^d \}.$$

Sample

$$\mathbf{w} \sim \mathsf{P}_{\mathcal{F}} := \mathsf{N}_d(\mathbf{0}, \mathbf{I}_d), \quad \mathbf{x}_i \overset{\mathsf{i.i.d.}}{\sim} \mathsf{P}_{\mathcal{X}} := \mathsf{N}_d(\mathbf{0}, \mathbf{I}_d), \quad i = 1, \dots, N,$$

independently. For each f determined by w, generate training pairs

$$\mathbf{y}_i = f(\mathbf{x}_i) = \mathbf{w}^{\mathrm{T}} \mathbf{x}_i.$$

Then, train the Transformer to predict  $\hat{\mathbf{y}}_i = \mathbf{W} \, \mathrm{TF}_{m{ heta}}(\mathbf{X}_{\mathsf{prompt}}^{(i)})$  where

$$\mathbf{X}_{\mathsf{prompt}}^{(i)} = (\mathbf{x}_1, \mathbf{y}_1, \dots, \mathbf{x}_{i-1}, \mathbf{y}_{i-1}, \mathbf{x}_i)$$

$$\min_{\theta} \underset{\substack{\mathbf{x}_1, \dots, \mathbf{x}_N \sim \mathbf{N}_d(\mathbf{0}, \mathbf{I}_d) \\ f \sim \mathbf{N}_d(\mathbf{0}, \mathbf{I}_d)}}{\mathbb{E}} \left[ \frac{1}{N} \sum_{i=1}^N (\hat{\mathbf{y}}_i - \mathbf{y}_i)^2 \right]$$

Example (Linear functions)

Consider the class of linear functions

$$\mathcal{F} = \{ f : \mathbb{R}^d \to \mathbb{R} \mid f(\mathbf{x}) = \mathbf{w}^{\mathrm{T}} \mathbf{x}, \, \mathbf{w} \in \mathbb{R}^d \}.$$

Sample

$$\mathbf{w} \sim \mathsf{P}_{\mathcal{F}} := \mathsf{N}_d(\mathbf{0}, \mathbf{I}_d), \quad \mathbf{x}_i \overset{\mathsf{i.i.d.}}{\sim} \mathsf{P}_{\mathcal{X}} := \mathsf{N}_d(\mathbf{0}, \mathbf{I}_d), \quad i = 1, \dots, N,$$

independently. For each f determined by w, generate training pairs

$$\mathbf{y}_i = f(\mathbf{x}_i) = \mathbf{w}^{\mathrm{T}} \mathbf{x}_i.$$

Then, train the Transformer to predict  $\hat{\mathbf{y}}_i = \mathbf{W} \, \mathrm{TF}_{m{ heta}}(\mathbf{X}_{\mathsf{prompt}}^{(i)})$  where

$$\mathbf{X}_{\mathsf{prompt}}^{(i)} = (\mathbf{x}_1, \mathbf{y}_1, \dots, \mathbf{x}_{i-1}, \mathbf{y}_{i-1}, \mathbf{x}_i),$$

$$\min_{\boldsymbol{\theta}} \mathbb{E}_{\mathbf{x}_1, \dots, \mathbf{x}_N \sim N_d(\mathbf{0}, \mathbf{I}_d)} \left[ \frac{1}{N} \sum_{i=1}^N (\hat{\mathbf{y}}_i - \mathbf{y}_i)^2 \right].$$

In this setting, the Transformer learns to approximate the least squares estimator.

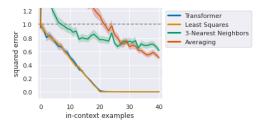


Figure: Normalized squared error of the Transformer as a function of the number of in-context examples (Garg et al., 2022).

Bai et al. (2023) show<sup>1</sup> that Transformers can implement a broad class of standard ML algorithms in context, such as least squares, Ridge regression, Lasso, convex risk minimization for GLMs.

<sup>&</sup>lt;sup>1</sup>Possibly the most mathematically formal paper on this topic

In this setting, the Transformer learns to approximate the least squares estimator.

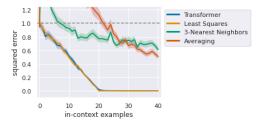


Figure: Normalized squared error of the Transformer as a function of the number of in-context examples (Garg et al., 2022).

Bai et al. (2023) show<sup>1</sup> that Transformers can implement a broad class of standard ML algorithms in context, such as least squares, Ridge regression, Lasso, convex risk minimization for GLMs.

<sup>&</sup>lt;sup>1</sup>Possibly the most mathematically formal paper on this topic.

# Large Language Models

#### Decoder transformers

These can be used as generative models that create output sequences of tokens. we will focus on a class of models called GPT which stands for generative pre-trained transformer.

#### Goai

Construct an autoregressive model of the form defined by

$$p(\mathbf{x}_1, \dots, \mathbf{x}_N) = \prod_{n=1}^N p(\mathbf{x}_n \mid \mathbf{x}_{1:n-1}), \quad \mathbf{x}_{1:n-1} = (\mathbf{x}_1, \dots, \mathbf{x}_{n-1})$$

where conditional distributions  $p(\mathbf{x}_n \mid \mathbf{x}_{1:n-1})$  are expressed using a Transformer neural network that is learned from data.

TDEA

$$\mathbf{x}_1, \dots, \mathbf{x}_N \overset{\mathsf{Transformer layers}}{\Rightarrow} \widetilde{\mathbf{x}}_1, \dots, \widetilde{\mathbf{x}}_N$$

Each output needs to represent a probability distribution over the class of tokens

$$\mathbf{Y} = \operatorname{Softmax}(\widetilde{\mathbf{X}}\mathbf{W}^{(p)}), \quad \mathbf{W}^{(p)} \in \mathbb{R}^{D \times K}$$

# Large Language Models

#### Decoder transformers

These can be used as generative models that create output sequences of tokens. we will focus on a class of models called GPT which stands for generative pre-trained transformer.

GOAL

Construct an autoregressive model of the form defined by

$$p(\mathbf{x}_1, \dots, \mathbf{x}_N) = \prod_{n=1}^{N} p(\mathbf{x}_n \mid \mathbf{x}_{1:n-1}), \quad \mathbf{x}_{1:n-1} = (\mathbf{x}_1, \dots, \mathbf{x}_{n-1}),$$

where conditional distributions  $p(\mathbf{x}_n \mid \mathbf{x}_{1:n-1})$  are expressed using a Transformer neural network that is learned from data.

IDEA

$$\mathbf{x}_1, \dots, \mathbf{x}_N \overset{\mathsf{Transformer layers}}{\Rightarrow} \widetilde{\mathbf{x}}_1, \dots, \widetilde{\mathbf{x}}_N$$

Each output needs to represent a probability distribution over the class of tokens

$$\mathbf{Y} = \text{Softmax}(\widetilde{\mathbf{X}}\mathbf{W}^{(\text{p})}), \quad \mathbf{W}^{(\text{p})} \in \mathbb{R}^{D \times K}$$

# Large Language Models

#### Decoder transformers

These can be used as generative models that create output sequences of tokens. we will focus on a class of models called GPT which stands for generative pre-trained transformer.

Goal

Construct an autoregressive model of the form defined by

$$p(\mathbf{x}_1,\ldots,\mathbf{x}_N) = \prod_{n=1}^N p(\mathbf{x}_n \mid \mathbf{x}_{1:n-1}), \quad \mathbf{x}_{1:n-1} = (\mathbf{x}_1,\ldots,\mathbf{x}_{n-1}),$$

where conditional distributions  $p(\mathbf{x}_n \mid \mathbf{x}_{1:n-1})$  are expressed using a Transformer neural network that is learned from data.

IDEA

$$\mathbf{x}_1, \dots, \mathbf{x}_N \quad \stackrel{\mathsf{Transformer layers}}{\Rightarrow} \quad \widetilde{\mathbf{x}}_1, \dots, \widetilde{\mathbf{x}}_N,$$

Each output needs to represent a probability distribution over the class of tokens:

$$\mathbf{Y} = \operatorname{Softmax}(\widetilde{\mathbf{X}}\mathbf{W}^{(p)}), \quad \mathbf{W}^{(p)} \in \mathbb{R}^{D \times K}.$$

#### EMBEDDING AND POSITIONAL ENCODING

We begin by defining a dictionary

$$\mathcal{V} = \{w_1, w_2, \dots, w_K\},\$$

which is a fixed, finite set of words indexed from 1 to K. A word  $v_i$  is mapped to an index in the dictionary using a tokenization function  $\tau: \mathcal{V} \to \{1, \dots, K\}$ , such that

$$\mathbf{v}(\mathbf{v}_i) = k_i$$
 with  $\mathbf{v}_i = w_{k_i}$ .

This index is then used to construct a **one-hot encoded vector**  $\mathbf{o}_i \in \{0,1\}^K$ , where the  $k_i$ -th component is set to 1 and all others are 0.

$$\mathbf{x}_i = \mathbf{E}\mathbf{o}_i + \mathbf{r}_i$$

#### EMBEDDING AND POSITIONAL ENCODING

We begin by defining a dictionary

$$\mathcal{V} = \{w_1, w_2, \dots, w_K\},\$$

which is a fixed, finite set of words indexed from 1 to K. A word  $\mathbf{v}_i$  is mapped to an index in the dictionary using a tokenization function  $\tau: \mathcal{V} \to \{1, \dots, K\}$ , such that

$$\tau(\mathbf{v}_i) = k_i$$
 with  $\mathbf{v}_i = w_{k_i}$ .

This index is then used to construct a one-hot encoded vector  $\mathbf{o}_i \in \{0,1\}^K$ , where the  $k_i$ -th component is set to 1 and all others are 0.

$$\mathbf{x}_i = \mathbf{Eo}_i + \mathbf{r}_i$$

#### Embedding and positional encoding

We begin by defining a dictionary

$$\mathcal{V} = \{w_1, w_2, \dots, w_K\},\$$

which is a fixed, finite set of words indexed from 1 to K. A word  $\mathbf{v}_i$  is mapped to an index in the dictionary using a tokenization function  $\tau: \mathcal{V} \to \{1, \dots, K\}$ , such that

$$\tau(\mathbf{v}_i) = k_i$$
 with  $\mathbf{v}_i = w_{k_i}$ .

This index is then used to construct a **one-hot encoded vector**  $\mathbf{o}_i \in \{0,1\}^K$ , where the  $k_i$ -th component is set to 1 and all others are 0.

$$\mathbf{x}_i = \mathbf{Eo}_i + \mathbf{r}_i$$

#### EMBEDDING AND POSITIONAL ENCODING

We begin by defining a dictionary

$$\mathcal{V} = \{w_1, w_2, \dots, w_K\},\$$

which is a fixed, finite set of words indexed from 1 to K. A word  $\mathbf{v}_i$  is mapped to an index in the dictionary using a tokenization function  $\tau: \mathcal{V} \to \{1, \dots, K\}$ , such that

$$\tau(\mathbf{v}_i) = k_i$$
 with  $\mathbf{v}_i = w_{k_i}$ .

This index is then used to construct a **one-hot encoded vector**  $\mathbf{o}_i \in \{0,1\}^K$ , where the  $k_i$ -th component is set to 1 and all others are 0.

$$\mathbf{x}_i = \mathbf{E}\mathbf{o}_i + \mathbf{r}_i.$$

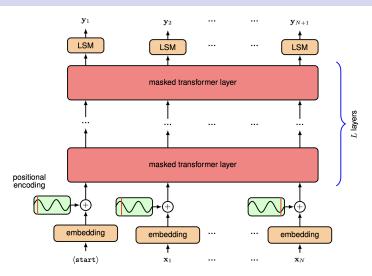


Figure: Architecture of a GPT decoder transformer network. Here 'LSM' stands for linear-softmax and denotes a linear transformation whose learnable parameters are shared across the token positions, followed by a softmax activation function.

## How do we apply GPT to generate sequences of extreme events?<sup>2</sup>

- Let  $(X_1, \ldots, X_d) \in \mathcal{X}^d$  be a chain of extreme events, e.g.  $\mathcal{X} = \mathbb{R}_+$ .
- ▶ Define a finite partition  $\{\mathcal{C}_k\}_{k=1}^K$  of  $\mathcal{X}$ , and a one-hot tokenization function  $\tau:\mathcal{X}\to\{0,1\}^K$  defined as

$$\tau(X_i) = (I(X_i \in \mathcal{C}_1), \dots, I(X_i \in \mathcal{C}_K)), \quad i = 1, \dots, d$$

▶ Each token is then projected into an embedding space using a learnable matrix  $\mathbf{E} \in \mathbb{R}^{D \times K}$ , positional encoding  $\mathbf{r}_i$  and extremal encoding<sup>3</sup> are added

$$\mathbf{s}_i = \mathbf{E}\tau(X_i) + \mathbf{r}_i + \mathbf{e}_i.$$

$$P(X_1 \in C_{k_1}, \dots, X_d \in C_{k_d}) = \prod_{i=1}^d P(X_i \in C_{k_i} \mid X_{1:i-1})$$

<sup>&</sup>lt;sup>2</sup>Ongoing work with Miguel de Carvalho and Ronny Vallejos.

e.g. based on extreme quantile regression models

How do we apply GPT to generate sequences of extreme events?<sup>2</sup>

- Let  $(X_1, \ldots, X_d) \in \mathcal{X}^d$  be a chain of extreme events, e.g.  $\mathcal{X} = \mathbb{R}_+$ .
- ▶ Define a finite partition  $\{\mathcal{C}_k\}_{k=1}^K$  of  $\mathcal{X}$ , and a one-hot tokenization function  $\tau:\mathcal{X}\to\{0,1\}^K$  defined as

$$\tau(X_i) = (I(X_i \in \mathcal{C}_1), \dots, I(X_i \in \mathcal{C}_K)), \quad i = 1, \dots, d.$$

▶ Each token is then projected into an embedding space using a learnable matrix  $\mathbf{E} \in \mathbb{R}^{D \times K}$ , positional encoding  $\mathbf{r}_i$  and extremal encoding<sup>3</sup> are added

$$\mathbf{s}_i = \mathbf{E}\tau(X_i) + \mathbf{r}_i + \mathbf{e}_i.$$

$$P(X_1 \in C_{k_1}, \dots, X_d \in C_{k_d}) = \prod_{i=1}^d P(X_i \in C_{k_i} \mid X_{1:i-1}).$$

<sup>&</sup>lt;sup>2</sup>Ongoing work with Miguel de Carvalho and Ronny Vallejos.

e.g. based on extreme quantile regression models

How do we apply GPT to generate sequences of extreme events?<sup>2</sup>

- ▶ Let  $(X_1, ..., X_d) \in \mathcal{X}^d$  be a chain of extreme events, e.g.  $\mathcal{X} = \mathbb{R}_+$ .
- ▶ Define a finite partition  $\{\mathcal{C}_k\}_{k=1}^K$  of  $\mathcal{X}$ , and a one-hot tokenization function  $\tau: \mathcal{X} \to \{0,1\}^K$  defined as

$$\tau(X_i) = (I(X_i \in \mathcal{C}_1), \dots, I(X_i \in \mathcal{C}_K)), \quad i = 1, \dots, d.$$

▶ Each token is then projected into an embedding space using a learnable matrix  $\mathbf{E} \in \mathbb{R}^{D \times K}$ , positional encoding  $\mathbf{r}_i$  and extremal encoding<sup>3</sup> are added

$$\mathbf{s}_i = \mathbf{E}\tau(X_i) + \mathbf{r}_i + \mathbf{e}_i.$$

$$P(X_1 \in C_{k_1}, \dots, X_d \in C_{k_d}) = \prod_{i=1}^d P(X_i \in C_{k_i} \mid X_{1:i-1})$$

<sup>&</sup>lt;sup>2</sup>Ongoing work with Miguel de Carvalho and Ronny Vallejos.

e.g. based on extreme quantile regression models

How do we apply GPT to generate sequences of extreme events?<sup>2</sup>

- ▶ Let  $(X_1, ..., X_d) \in \mathcal{X}^d$  be a chain of extreme events, e.g.  $\mathcal{X} = \mathbb{R}_+$ .
- ▶ Define a finite partition  $\{\mathcal{C}_k\}_{k=1}^K$  of  $\mathcal{X}$ , and a one-hot tokenization function  $\tau: \mathcal{X} \to \{0,1\}^K$  defined as

$$\tau(X_i) = (I(X_i \in \mathcal{C}_1), \dots, I(X_i \in \mathcal{C}_K)), \quad i = 1, \dots, d.$$

Each token is then projected into an embedding space using a learnable matrix  $\mathbf{E} \in \mathbb{R}^{D \times K}$ , positional encoding  $\mathbf{r}_i$  and extremal encoding<sup>3</sup> are added

$$\mathbf{s}_i = \mathbf{E}\tau(X_i) + \mathbf{r}_i + \mathbf{e}_i.$$

$$P(X_1 \in C_{k_1}, \dots, X_d \in C_{k_d}) = \prod_{i=1}^d P(X_i \in C_{k_i} \mid X_{1:i-1})$$

<sup>&</sup>lt;sup>2</sup>Ongoing work with Miguel de Carvalho and Ronny Vallejos.

<sup>&</sup>lt;sup>3</sup>e.g. based on extreme quantile regression models.

How do we apply GPT to generate sequences of extreme events?<sup>2</sup>

- ▶ Let  $(X_1, ..., X_d) \in \mathcal{X}^d$  be a chain of extreme events, e.g.  $\mathcal{X} = \mathbb{R}_+$ .
- ▶ Define a finite partition  $\{\mathcal{C}_k\}_{k=1}^K$  of  $\mathcal{X}$ , and a one-hot tokenization function  $\tau: \mathcal{X} \to \{0,1\}^K$  defined as

$$\tau(X_i) = (I(X_i \in \mathcal{C}_1), \dots, I(X_i \in \mathcal{C}_K)), \quad i = 1, \dots, d.$$

▶ Each token is then projected into an embedding space using a learnable matrix  $\mathbf{E} \in \mathbb{R}^{D \times K}$ , positional encoding  $\mathbf{r}_i$  and extremal encoding<sup>3</sup> are added

$$\mathbf{s}_i = \mathbf{E}\tau(X_i) + \mathbf{r}_i + \mathbf{e}_i.$$

$$P(X_1 \in C_{k_1}, \dots, X_d \in C_{k_d}) = \prod_{i=1}^d P(X_i \in C_{k_i} \mid X_{1:i-1}).$$

<sup>&</sup>lt;sup>2</sup>Ongoing work with Miguel de Carvalho and Ronny Vallejos.

<sup>&</sup>lt;sup>3</sup>e.g. based on extreme quantile regression models.

How do we apply GPT to generate sequences of extreme events?

Each conditional  $P(X_i \in \mathcal{C}_k \mid X_{i-L:i-1})$  is approximated by a **decoder Transformer**  $^4$   $\mathrm{TF}_{\boldsymbol{\theta}}$ , which maps the input embeddings  $(\mathbf{s}_{i-L}, \ldots, \mathbf{s}_{i-1})$ :

$$\widetilde{\mathbf{S}}_i = \mathrm{TF}_{\boldsymbol{\theta}}(\mathbf{S}_i), \quad \text{with } \mathbf{S}_i = (\mathbf{s}_{i-L}, \dots, \mathbf{s}_{i-1})^{\mathrm{T}},$$

This representation is then passed through a learnable linear projection  $\mathbf{W} \in \mathbb{R}^{K \times D}$  to obtain the logits  $\mathbf{z}_i = (z_{i1}, \dots, z_{iK}) \in \mathbb{R}^K$ :

$$\mathbf{z}_i = \mathbf{W}\widetilde{\mathbf{S}}_i$$

Finally, the softmax layer converts the logits into probabilities:

$$P_{\theta}(X_i \in \mathcal{C}_k \mid X_{i-L:i-1}) = \frac{\exp(z_{ik})}{\sum_{j=1}^K \exp(z_{ij})}, \quad k = 1, \dots, K$$

▶ The model is trained by minimizing the empirical negative log-likelihood over a large dataset  $\mathcal{D} = \{(X_1^{(t)}, \dots, X_d^{(t)})\}_{t=1}^N$ :

$$\hat{\ell}(\theta) = -\frac{1}{Nd} \sum_{t=1}^{N} \sum_{i=1}^{d} \log P_{\theta}(X_{i}^{(t)} \in \mathcal{C}_{k_{i}^{(t)}} \mid X_{i-L:i-1}^{(t)})$$

<sup>&</sup>lt;sup>4</sup>Just like ChatGPT!

How do we apply GPT to generate sequences of extreme events?

Each conditional  $P(X_i \in \mathcal{C}_k \mid X_{i-L:i-1})$  is approximated by a **decoder** Transformer<sup>4</sup> TF<sub>\theta</sub>, which maps the input embeddings  $(\mathbf{s}_{i-L}, \dots, \mathbf{s}_{i-1})$ :

$$\widetilde{\mathbf{S}}_i = \mathrm{TF}_{m{ heta}}(\mathbf{S}_i), \quad \text{with } \mathbf{S}_i = (\mathbf{s}_{i-L}, \dots, \mathbf{s}_{i-1})^{\mathrm{T}},$$

This representation is then passed through a learnable linear projection  $\mathbf{W} \in \mathbb{R}^{K \times D}$  to obtain the logits  $\mathbf{z}_i = (z_{i1}, \dots, z_{iK}) \in \mathbb{R}^K$ :

$$\mathbf{z}_i = \mathbf{W}\widetilde{\mathbf{S}}_i.$$

Finally, the softmax layer converts the logits into probabilities:

$$P_{\theta}(X_i \in C_k \mid X_{i-L:i-1}) = \frac{\exp(z_{ik})}{\sum_{j=1}^K \exp(z_{ij})}, \quad k = 1, \dots, K$$

▶ The model is trained by minimizing the empirical negative log-likelihood over a large dataset  $\mathcal{D} = \{(X_1^{(t)}, \dots, X_d^{(t)})\}_{t=1}^N$ :

$$\hat{\ell}(\theta) = -\frac{1}{Nd} \sum_{t=1}^{N} \sum_{i=1}^{d} \log P_{\theta}(X_i^{(t)} \in \mathcal{C}_{k_i^{(t)}} \mid X_{i-L:i-1}^{(t)}).$$

<sup>&</sup>lt;sup>4</sup>Just like ChatGPT!

How do we apply GPT to generate sequences of extreme events?

Each conditional  $P(X_i \in \mathcal{C}_k \mid X_{i-L:i-1})$  is approximated by a **decoder Transformer**<sup>4</sup>  $\mathrm{TF}_{\boldsymbol{\theta}}$ , which maps the input embeddings  $(\mathbf{s}_{i-L},\ldots,\mathbf{s}_{i-1})$ :

$$\widetilde{\mathbf{S}}_i = \mathrm{TF}_{m{ heta}}(\mathbf{S}_i), \quad \text{with } \mathbf{S}_i = (\mathbf{s}_{i-L}, \dots, \mathbf{s}_{i-1})^{\mathrm{T}},$$

This representation is then passed through a learnable linear projection  $\mathbf{W} \in \mathbb{R}^{K \times D}$  to obtain the logits  $\mathbf{z}_i = (z_{i1}, \dots, z_{iK}) \in \mathbb{R}^K$ :

$$\mathbf{z}_i = \mathbf{W}\widetilde{\mathbf{S}}_i$$
.

Finally, the softmax layer converts the logits into probabilities:

$$P_{\theta}(X_i \in C_k \mid X_{i-L:i-1}) = \frac{\exp(z_{ik})}{\sum_{j=1}^K \exp(z_{ij})}, \quad k = 1, \dots, K.$$

The model is trained by minimizing the empirical negative log-likelihood over a large dataset  $\mathcal{D} = \{(X_1^{(t)}, \dots, X_d^{(t)})\}_{t=1}^{N}$ :

$$\hat{\ell}(\theta) = -\frac{1}{Nd} \sum_{t=1}^{N} \sum_{i=1}^{d} \log P_{\theta}(X_i^{(t)} \in \mathcal{C}_{k_i^{(t)}} \mid X_{i-L:i-1}^{(t)}).$$

<sup>&</sup>lt;sup>4</sup>Just like ChatGPT!

How do we apply GPT to generate sequences of extreme events?

Each conditional  $P(X_i \in \mathcal{C}_k \mid X_{i-L:i-1})$  is approximated by a **decoder** Transformer<sup>4</sup>  $\mathrm{TF}_{\theta}$ , which maps the input embeddings  $(\mathbf{s}_{i-L}, \ldots, \mathbf{s}_{i-1})$ :

$$\widetilde{\mathbf{S}}_i = \mathrm{TF}_{\boldsymbol{\theta}}(\mathbf{S}_i), \quad \text{with } \mathbf{S}_i = (\mathbf{s}_{i-L}, \dots, \mathbf{s}_{i-1})^{\mathrm{T}},$$

This representation is then passed through a learnable linear projection  $\mathbf{W} \in \mathbb{R}^{K \times D}$  to obtain the logits  $\mathbf{z}_i = (z_{i1}, \dots, z_{iK}) \in \mathbb{R}^K$ :

$$\mathbf{z}_i = \mathbf{W}\widetilde{\mathbf{S}}_i$$
.

Finally, the softmax layer converts the logits into probabilities:

$$P_{\theta}(X_i \in \mathcal{C}_k \mid X_{i-L:i-1}) = \frac{\exp(z_{ik})}{\sum_{j=1}^K \exp(z_{ij})}, \quad k = 1, \dots, K.$$

▶ The model is trained by minimizing the empirical negative log-likelihood over a large dataset  $\mathcal{D} = \{(X_1^{(t)}, \dots, X_d^{(t)})\}_{t=1}^N$ :

$$\hat{\ell}(\boldsymbol{\theta}) = -\frac{1}{Nd} \sum_{t=1}^{N} \sum_{i=1}^{d} \log P_{\boldsymbol{\theta}}(X_i^{(t)} \in \mathcal{C}_{k_i^{(t)}} \mid X_{i-L:i-1}^{(t)}).$$

<sup>&</sup>lt;sup>4</sup>Just like ChatGPT!

## References

#### References

- 1) Vaswani et al. (2017). Attention is All You Need. arXiv:1706.03762.
- 2) Ba et al. (2016). Layer Normalization. arXiv:1607.06450.
- 3) Garg et al. (2022). What Can Transformers Learn In-Context? A Case Study of Simple Function Classes. arXiv:2208.01066v3.
- 4) Akyürek et al. (2023). What Learning Algorithm Is In-Context Learning? Investigations with Linear Models. arXiv:2301.05207.
- Bai et al. (2023). Transformers as Statisticians: Provable In-Context Learning with In-Context Algorithm Selection. arXiv:2306.11644
- 6) Bishop & Bishop (2024). Deep Learning: Foundations and Concepts. Springer.

# Exploring Transformer-Based Models for Cascading Extremes

## Clemente Ferrer

clemente.ferrer@usm.cl

Department of Mathematics Universidad Técnica Federico Santa María

Joint work with Miguel de Carvalho & Ronny Vallejos

Workshop on Generative AI for Extreme Events

Edinburgh, UK June 13, 2025