

# On the Convergence of Encoder-only Shallow Transformers

**Fanghui Liu**

Department of Computer Science, University of Warwick, UK  
Centre for Discrete Mathematics and its Applications (DIMAP), Warwick

Based on joint work with

[Yongtao Wu (EPFL), Fanghui Liu, Grigorios Chrysos (UW-Madison), Volkan Cevher (EPFL)]

at MILD Seminar, University of British Columbia

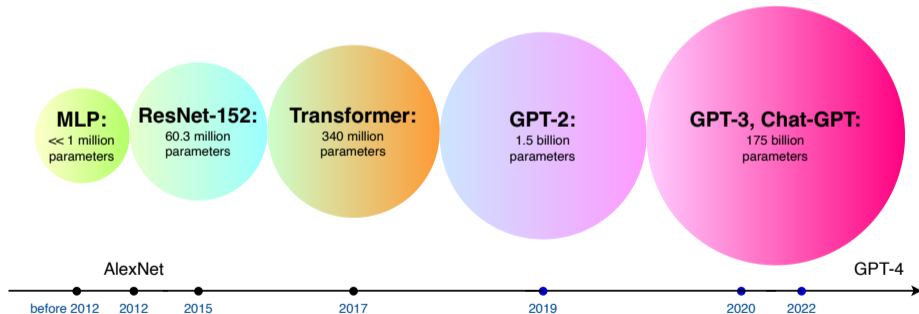


# Over-parameterization: more parameters than training data



```
... This code is not working like I expect - how do I fix it?  
  
def resultNumber(x: AdditionError)  
  def showResult(number: Int)  
    def print(x: Int)  
      def resultNumber() = 1 + resultNumber(x)  
    }  
  }  
  def x = 1 + number(x)  
  cancel()  
  if err == null {  
    return 1 + resultNumber(x)  
  }  
  return null.getError.getMessage() + resultNumber(x)  
}
```

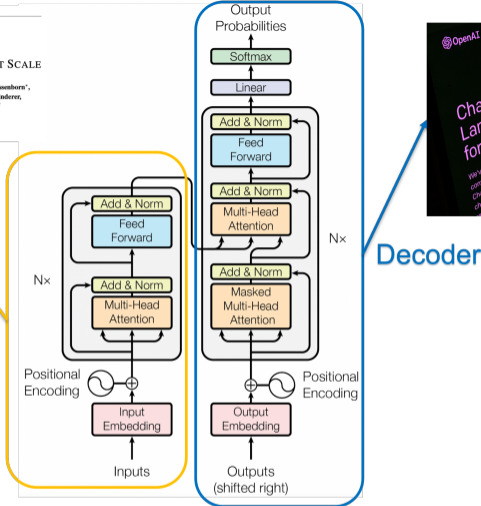
QUESTION: It's difficult to say what's wrong with the code without more context. Can you provide more information about what the code is supposed to do and what isn't working as expected? Also, is this the entire code or just a part of it?



# AN IMAGE IS WORTH 16X16 WORDS: TRANSFORMERS FOR IMAGE RECOGNITION AT SCALE

Alexey Dosovitskiy<sup>1\*</sup>, Lucas Beyer<sup>1</sup>, Alexander Kolesnikov<sup>1</sup>, Dirk Weissenborn<sup>1</sup>,  
Xiaohua Zhai<sup>1</sup>, Thomas Unterthiner<sup>1</sup>, Mostafa Dehghani<sup>1</sup>, Matthias Minderer<sup>1</sup>,  
Georg Heigold<sup>1</sup>, Sylvain Gugli<sup>1</sup>, Jakob Uszkoreit<sup>1</sup>, Neil Houlsby<sup>1\*</sup>  
<sup>\*</sup>equal technical contribution, <sup>1</sup>equal advising  
Google Research, Brain Team  
{dosovitskiy, neilhoulshy}@google.com

Encoder



Decoder



**Figure:** **Left:** Vision Transformer(ViT) [1], based on the encoder of Transformer. **Middle:** Original Transformer [2], with encoder and decoder. **Right:** ChatGPT, based on the decoder of Transformer.

## Self-attention

- ▶ input  $X \in \mathbb{R}^{d_s \times d}$
- ▶  $\sigma_s$ : soft-max (row-wise)
- ▶  $W_Q, W_K, W_V \in \mathbb{R}^{d_m \times d}$
- ▶  $d_s$ : number of tokens
- ▶  $d$ : the feature dimension of each token
- ▶  $d_m$ : width

## Self-attention

- ▶ input  $\mathbf{X} \in \mathbb{R}^{d_s \times d}$
- ▶  $\sigma_s$ : soft-max (row-wise)
- ▶  $\mathbf{W}_Q, \mathbf{W}_K, \mathbf{W}_V \in \mathbb{R}^{d_m \times d}$
- ▶  $d_s$ : number of tokens
- ▶  $d$ : the feature dimension of each token
- ▶  $d_m$ : width

$$\text{Self-attention}(\mathbf{X}) \triangleq \text{Softmax} \left( \tau_0(\mathbf{X}\mathbf{W}_Q^\top) (\mathbf{X}\mathbf{W}_K^\top)^\top \right) (\mathbf{X}\mathbf{W}_V^\top) = \sigma_s \left( \tau_0 \mathbf{X}\mathbf{W}_Q^\top \mathbf{W}_K \mathbf{X}^\top \right) (\mathbf{X}\mathbf{W}_V^\top) .$$

## Self-attention

- ▶ input  $\mathbf{X} \in \mathbb{R}^{d_s \times d}$
- ▶  $\sigma_s$ : soft-max (row-wise)
- ▶  $\mathbf{W}_Q, \mathbf{W}_K, \mathbf{W}_V \in \mathbb{R}^{d_m \times d}$
- ▶  $d_s$ : number of tokens
- ▶  $d$ : the feature dimension of each token
- ▶  $d_m$ : width

$$\text{Self-attention}(\mathbf{X}) \triangleq \text{Softmax} \left( \tau_0 (\mathbf{X} \mathbf{W}_Q^\top) (\mathbf{X} \mathbf{W}_K^\top)^\top \right) (\mathbf{X} \mathbf{W}_V^\top) = \sigma_s \left( \tau_0 \mathbf{X} \mathbf{W}_Q^\top \mathbf{W}_K \mathbf{X}^\top \right) (\mathbf{X} \mathbf{W}_V^\top) .$$

$$\text{input of softmax: } [\tau_0 \mathbf{W}_Q^\top \mathbf{W}_K]_{ij} = \tau_0 \sum_{k=1}^{d_m} [\mathbf{W}_Q^\top]_{ik} [\mathbf{W}_K]_{kj}$$

## Self-attention

- ▶ input  $\mathbf{X} \in \mathbb{R}^{d_s \times d}$
- ▶  $\sigma_s$ : soft-max (row-wise)
- ▶  $\mathbf{W}_Q, \mathbf{W}_K, \mathbf{W}_V \in \mathbb{R}^{d_m \times d}$
- ▶  $d_s$ : number of tokens
- ▶  $d$ : the feature dimension of each token
- ▶  $d_m$ : width

$$\text{Self-attention}(\mathbf{X}) \triangleq \text{Softmax} \left( \tau_0 (\mathbf{X} \mathbf{W}_Q^\top) (\mathbf{X} \mathbf{W}_K^\top)^\top \right) (\mathbf{X} \mathbf{W}_V^\top) = \sigma_s \left( \tau_0 \mathbf{X} \mathbf{W}_Q^\top \mathbf{W}_K \mathbf{X}^\top \right) (\mathbf{X} \mathbf{W}_V^\top) .$$

$$\text{input of softmax: } [\tau_0 \mathbf{W}_Q^\top \mathbf{W}_K]_{ij} = \tau_0 \sum_{k=1}^{d_m} [\mathbf{W}_Q^\top]_{ik} [\mathbf{W}_K]_{kj}$$

- scaling schemes given by  $\mathbf{W}_Q, \mathbf{W}_K$  initialized by standard Gaussian

## Self-attention

- ▶ input  $\mathbf{X} \in \mathbb{R}^{d_s \times d}$
- ▶  $\sigma_s$ : soft-max (row-wise)
- ▶  $\mathbf{W}_Q, \mathbf{W}_K, \mathbf{W}_V \in \mathbb{R}^{d_m \times d}$
- ▶  $d_s$ : number of tokens
- ▶  $d$ : the feature dimension of each token
- ▶  $d_m$ : width

$$\text{Self-attention}(\mathbf{X}) \triangleq \text{Softmax} \left( \tau_0 (\mathbf{X} \mathbf{W}_Q^\top) (\mathbf{X} \mathbf{W}_K^\top)^\top \right) (\mathbf{X} \mathbf{W}_V^\top) = \sigma_s \left( \tau_0 \mathbf{X} \mathbf{W}_Q^\top \mathbf{W}_K \mathbf{X}^\top \right) (\mathbf{X} \mathbf{W}_V^\top) .$$

$$\text{input of softmax: } [\tau_0 \mathbf{W}_Q^\top \mathbf{W}_K]_{ij} = \tau_0 \sum_{k=1}^{d_m} [\mathbf{W}_Q^\top]_{ik} [\mathbf{W}_K]_{kj}$$

◦ scaling schemes given by  $\mathbf{W}_Q, \mathbf{W}_K$  initialized by standard Gaussian

- ▶  $\tau_0 = d_m^{-1/2}$  in the original Transformer [2]:

$$[\tau_0 \mathbf{W}_Q^\top \mathbf{W}_K]_{ij} \text{ has zero-mean and unit variance } \quad \forall i, j \in [d]$$



## Self-attention

- ▶ input  $\mathbf{X} \in \mathbb{R}^{d_s \times d}$
- ▶  $\sigma_s$ : soft-max (row-wise)
- ▶  $\mathbf{W}_Q, \mathbf{W}_K, \mathbf{W}_V \in \mathbb{R}^{d_m \times d}$
- ▶  $d_s$ : number of tokens
- ▶  $d$ : the feature dimension of each token
- ▶  $d_m$ : width

$$\text{Self-attention}(\mathbf{X}) \triangleq \text{Softmax} \left( \tau_0 (\mathbf{X} \mathbf{W}_Q^\top) (\mathbf{X} \mathbf{W}_K^\top)^\top \right) (\mathbf{X} \mathbf{W}_V^\top) = \sigma_s \left( \tau_0 \mathbf{X} \mathbf{W}_Q^\top \mathbf{W}_K \mathbf{X}^\top \right) (\mathbf{X} \mathbf{W}_V^\top).$$

$$\text{input of softmax: } [\tau_0 \mathbf{W}_Q^\top \mathbf{W}_K]_{ij} = \tau_0 \sum_{k=1}^{d_m} [\mathbf{W}_Q^\top]_{ik} [\mathbf{W}_K]_{kj}$$

◦ scaling schemes given by  $\mathbf{W}_Q, \mathbf{W}_K$  initialized by standard Gaussian

- ▶  $\tau_0 = d_m^{-1/2}$  in the original Transformer [2]:

$$[\tau_0 \mathbf{W}_Q^\top \mathbf{W}_K]_{ij} \text{ has zero-mean and unit variance } \quad \forall i, j \in [d]$$

- ▶  $\tau_0 = d_m^{-1}$ : from the neural tangent kernel (NTK) analysis [3] for  $d_m \rightarrow \infty$ .

$$\lim_{d_m \rightarrow \infty} \tau_0 [\mathbf{W}_Q^\top \mathbf{W}_K]^{(ij)} = \lim_{d_m \rightarrow \infty} \frac{1}{d_m} \sum_{k=1}^{d_m} [\mathbf{W}_Q^\top]_{ik} [\mathbf{W}_K]_{kj} = 0.$$

## Self-attention

- ▶ input  $\mathbf{X} \in \mathbb{R}^{d_s \times d}$
- ▶  $\sigma_s$ : soft-max (row-wise)
- ▶  $\mathbf{W}_Q, \mathbf{W}_K, \mathbf{W}_V \in \mathbb{R}^{d_m \times d}$
- ▶  $d_s$ : number of tokens
- ▶  $d$ : the feature dimension of each token
- ▶  $d_m$ : width

$$\text{Self-attention}(\mathbf{X}) \triangleq \text{Softmax} \left( \tau_0 (\mathbf{X} \mathbf{W}_Q^\top) (\mathbf{X} \mathbf{W}_K^\top)^\top \right) (\mathbf{X} \mathbf{W}_V^\top) = \sigma_s \left( \tau_0 \mathbf{X} \mathbf{W}_Q^\top \mathbf{W}_K \mathbf{X}^\top \right) (\mathbf{X} \mathbf{W}_V^\top).$$

$$\text{input of softmax: } [\tau_0 \mathbf{W}_Q^\top \mathbf{W}_K]_{ij} = \tau_0 \sum_{k=1}^{d_m} [\mathbf{W}_Q^\top]_{ik} [\mathbf{W}_K]_{kj}$$

◦ scaling schemes given by  $\mathbf{W}_Q, \mathbf{W}_K$  initialized by standard Gaussian

- ▶  $\tau_0 = d_m^{-1/2}$  in the original Transformer [2]:

$$[\tau_0 \mathbf{W}_Q^\top \mathbf{W}_K]_{ij} \text{ has zero-mean and unit variance } \quad \forall i, j \in [d]$$

- ▶  $\tau_0 = d_m^{-1}$ : from the neural tangent kernel (NTK) analysis [3] for  $d_m \rightarrow \infty$ .

$$\lim_{d_m \rightarrow \infty} \tau_0 [\mathbf{W}_Q^\top \mathbf{W}_K]^{(ij)} = \lim_{d_m \rightarrow \infty} \frac{1}{d_m} \sum_{k=1}^{d_m} [\mathbf{W}_Q^\top]_{ik} [\mathbf{W}_K]_{kj} = 0.$$

Softmax becomes a pooling layer!

## Previous attempts on scaling in theory

$$[\tau_0 \mathbf{W}_Q^\top \mathbf{W}_K]_{ij}$$

◦ scaling schemes given by  $\mathbf{W}_Q, \mathbf{W}_K$  initialized by standard Gaussian

▶  $\tau_0 = d_m^{-1/2}$  in the original Transformer [2]:

$$[\tau_0 \mathbf{W}_Q^\top \mathbf{W}_K]_{ij} \text{ has zero-mean and unit variance } \forall i, j \in [d]$$

## Previous attempts on scaling in theory

$$[\tau_0 \mathbf{W}_Q^\top \mathbf{W}_K]_{ij}$$

◦ scaling schemes given by  $\mathbf{W}_Q, \mathbf{W}_K$  initialized by standard Gaussian

▶  $\tau_0 = d_m^{-1/2}$  in the original Transformer [2]:

$[\tau_0 \mathbf{W}_Q^\top \mathbf{W}_K]_{ij}$  has zero-mean and unit variance  $\forall i, j \in [d]$

$$[4]: \underbrace{\text{Softmax}}_{\text{ReLU}} \left( \tau_0 (\mathbf{XW}_Q^\top) (\mathbf{XW}_K^\top)^\top \right) (\mathbf{XW}_V^\top)$$

## Previous attempts on scaling in theory

$$[\tau_0 \mathbf{W}_Q^\top \mathbf{W}_K]_{ij}$$

◦ scaling schemes given by  $\mathbf{W}_Q, \mathbf{W}_K$  initialized by standard Gaussian

- ▶  $\tau_0 = d_m^{-1/2}$  in the original Transformer [2]:

$$[\tau_0 \mathbf{W}_Q^\top \mathbf{W}_K]_{ij} \text{ has zero-mean and unit variance } \forall i, j \in [d]$$

$$[4]: \underbrace{\text{Softmax}}_{\text{ReLU}} \left( \tau_0 (\mathbf{X} \mathbf{W}_Q^\top) (\mathbf{X} \mathbf{W}_K^\top)^\top \right) (\mathbf{X} \mathbf{W}_V^\top)$$

- ▶  $\tau_0 = d_m^{-1}$ : from the neural tangent kernel (NTK) analysis [3] for  $d_m \rightarrow \infty$ .

$$\lim_{d_m \rightarrow \infty} \tau_0 [\mathbf{W}_Q^\top \mathbf{W}_K]^{(ij)} = \lim_{d_m \rightarrow \infty} \frac{1}{d_m} \sum_{k=1}^{d_m} [\mathbf{W}_Q^\top]_{ik} [\mathbf{W}_K]_{kj} = 0.$$

$$[5]: \text{ setting } \mathbf{W}_Q = \mathbf{W}_K$$

## Question

*How can we do analysis of Transformers under a realistic setting?*

## Question

*How can we do analysis of Transformers under a realistic setting?*

even though

- ▶ a shallow Transformer
- ▶ an encoder-only shallow Transformer
- ▶ global convergence
- ▶ under the lazy training regime

## Question

*How can we do analysis of Transformers under a realistic setting?*

even though

- ▶ a shallow Transformer
- ▶ an encoder-only shallow Transformer
- ▶ global convergence
- ▶ under the lazy training regime

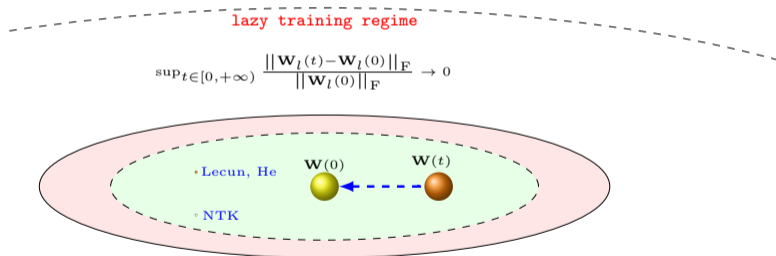


Figure: Training dynamics of two-layer ReLU NNs with infinite width under different initializations [3, 6, 7].



## Problem setting: encoder-only shallow Transformer

$$\mathbf{A}_1 = \text{Self-attention}(\mathbf{X}) \triangleq \sigma_s \left( \tau_0(\mathbf{X}\mathbf{W}_Q^\top) (\mathbf{X}\mathbf{W}_K^\top)^\top \right) (\mathbf{X}\mathbf{W}_V^\top),$$

$$\mathbf{A}_2 = \tau_1 \sigma_r(\mathbf{A}_1 \mathbf{W}_H), \quad \mathbf{a}_3 = \varphi(\mathbf{A}_2), \quad f(\mathbf{X}; \boldsymbol{\theta}) = \mathbf{a}_3^\top \mathbf{w}_O.$$

- ▶ Input:  $\mathbf{X} \in \mathbb{R}^{d_s \times d}$  ( $d_s$  is the number of tokens and  $d$  is the feature dimension of each token)

## Problem setting: encoder-only shallow Transformer

$$\mathbf{A}_1 = \text{Self-attention}(\mathbf{X}) \triangleq \sigma_s \left( \tau_0(\mathbf{X}\mathbf{W}_Q^\top) (\mathbf{X}\mathbf{W}_K^\top)^\top \right) (\mathbf{X}\mathbf{W}_V^\top),$$

$$\mathbf{A}_2 = \tau_1 \sigma_r(\mathbf{A}_1 \mathbf{W}_H), \quad \mathbf{a}_3 = \varphi(\mathbf{A}_2), \quad f(\mathbf{X}; \boldsymbol{\theta}) = \mathbf{a}_3^\top \mathbf{w}_O.$$

- ▶ Input:  $\mathbf{X} \in \mathbb{R}^{d_s \times d}$  ( $d_s$  is the number of tokens and  $d$  is the feature dimension of each token)
- ▶ A *self-attention layer*:  $\sigma_s$  is the row-wise softmax function and the learnable parameters are  $\mathbf{W}_Q, \mathbf{W}_K, \mathbf{W}_V \in \mathbb{R}^{d_m \times d}$ .

## Problem setting: encoder-only shallow Transformer

$$\mathbf{A}_1 = \text{Self-attention}(\mathbf{X}) \triangleq \sigma_s \left( \tau_0(\mathbf{X}\mathbf{W}_Q^\top) (\mathbf{X}\mathbf{W}_K^\top)^\top \right) (\mathbf{X}\mathbf{W}_V^\top),$$

$$\mathbf{A}_2 = \tau_1 \sigma_r(\mathbf{A}_1 \mathbf{W}_H), \quad \mathbf{a}_3 = \varphi(\mathbf{A}_2), \quad f(\mathbf{X}; \boldsymbol{\theta}) = \mathbf{a}_3^\top \mathbf{w}_O.$$

- ▶ Input:  $\mathbf{X} \in \mathbb{R}^{d_s \times d}$  ( $d_s$  is the number of tokens and  $d$  is the feature dimension of each token)
- ▶ A *self-attention layer*:  $\sigma_s$  is the row-wise softmax function and the learnable parameters are  $\mathbf{W}_Q, \mathbf{W}_K, \mathbf{W}_V \in \mathbb{R}^{d_m \times d}$ .
- ▶ A *feed-forward ReLU layer*:  $\sigma_r$  is the ReLU activation function; the learnable parameter is  $\mathbf{W}_H \in \mathbb{R}^{d_m \times d_m}$ . We assume  $\mathbf{W}_H = \mathbf{I}$ .
- ▶ An *average pooling layer*:  $\varphi$  indicates the column-wise average pooling.
- ▶ An *output layer* with learnable parameter  $\mathbf{w}_O \in \mathbb{R}^{d_m}$ .

## Problem setting: encoder-only shallow Transformer

$$\begin{aligned} \mathbf{A}_1 &= \text{Self-attention}(\mathbf{X}) \triangleq \sigma_s \left( \tau_0(\mathbf{X}\mathbf{W}_Q^\top) (\mathbf{X}\mathbf{W}_K^\top)^\top \right) (\mathbf{X}\mathbf{W}_V^\top), \\ \mathbf{A}_2 &= \tau_1 \sigma_r(\mathbf{A}_1 \mathbf{W}_H), \quad \mathbf{a}_3 = \varphi(\mathbf{A}_2), \quad f(\mathbf{X}; \boldsymbol{\theta}) = \mathbf{a}_3^\top \mathbf{w}_O. \end{aligned}$$

- ▶ Input:  $\mathbf{X} \in \mathbb{R}^{d_s \times d}$  ( $d_s$  is the number of tokens and  $d$  is the feature dimension of each token)
- ▶ A *self-attention layer*:  $\sigma_s$  is the row-wise softmax function and the learnable parameters are  $\mathbf{W}_Q, \mathbf{W}_K, \mathbf{W}_V \in \mathbb{R}^{d_m \times d}$ .
- ▶ A *feed-forward ReLU layer*:  $\sigma_r$  is the ReLU activation function; the learnable parameter is  $\mathbf{W}_H \in \mathbb{R}^{d_m \times d_m}$ . We assume  $\mathbf{W}_H = \mathbf{I}$ .
- ▶ An *average pooling layer*:  $\varphi$  indicates the column-wise average pooling.
- ▶ An *output layer* with learnable parameter  $\mathbf{w}_O \in \mathbb{R}^{d_m}$ .

Initialization	$\eta_O$	$\eta_V$	$\eta_Q$	$\eta_K$	$\tau_1$
LeCun	$d_m^{-1}$	$d^{-1}$	$d^{-1}$	$d^{-1}$	1
He	$2d_m^{-1}$	$2d^{-1}$	$2d^{-1}$	$2d^{-1}$	1
NTK	1	1	1	1	$d_m^{-1/2}$

## Training by gradient descent

- ▶ data  $\{(X_i, y_i)\}_{i=1}^n$  with  $\mathbf{y} = [y_1, y_2, \dots, y_n]^\top$
- ▶ model output  $\mathbf{f}(\boldsymbol{\theta}) := [f(X_1; \boldsymbol{\theta}), f(X_2; \boldsymbol{\theta}), \dots, f(X_n; \boldsymbol{\theta})]^\top$

$$\ell(\boldsymbol{\theta}) = \frac{1}{2} \|\mathbf{f}(\boldsymbol{\theta}) - \mathbf{y}\|_2^2$$

## Training by gradient descent

- ▶ data  $\{(X_i, y_i)\}_{i=1}^n$  with  $\mathbf{y} = [y_1, y_2, \dots, y_n]^\top$
- ▶ model output  $\mathbf{f}(\boldsymbol{\theta}) := [f(X_1; \boldsymbol{\theta}), f(X_2; \boldsymbol{\theta}), \dots, f(X_n; \boldsymbol{\theta})]^\top$

$$\ell(\boldsymbol{\theta}) = \frac{1}{2} \|\mathbf{f}(\boldsymbol{\theta}) - \mathbf{y}\|_2^2$$

---

### Algorithm 2: Gradient descent training

---

**Input:** data  $(X_i, y_i)_{i=1}^n$ , step size  $\gamma$ .

Initialize weights as follows:

$$\boldsymbol{\theta}^0 := \{\mathbf{W}_Q^0, \mathbf{W}_K^0, \mathbf{W}_V^0, \mathbf{W}_O^0\}.$$

**for**  $t = 0$  **to**  $t' - 1$  **do**

$$\mathbf{W}_Q^{t+1} = \mathbf{W}_Q^t - \gamma \cdot \nabla_{\mathbf{W}_Q} \ell(\boldsymbol{\theta}^t), \quad \mathbf{W}_K^{t+1} = \mathbf{W}_K^t - \gamma \cdot \nabla_{\mathbf{W}_K} \ell(\boldsymbol{\theta}^t),$$

$$\mathbf{W}_V^{t+1} = \mathbf{W}_V^t - \gamma \cdot \nabla_{\mathbf{W}_V} \ell(\boldsymbol{\theta}^t), \quad \mathbf{W}_O^{t+1} = \mathbf{W}_O^t - \gamma \cdot \nabla_{\mathbf{W}_O} \ell(\boldsymbol{\theta}^t).$$

**end for**

**Output:** the model based on  $\boldsymbol{\theta}^{t'}$ .

---

## Assumptions on data

### Assumption (Bounded data)

*The input data is bounded  $\|\mathbf{X}\|_F \leq \sqrt{d_s} C_x$  with some positive constant  $C_x$ .*

- The embedding of token can have a unit norm [8] independent of  $d$ .

## Assumptions on data

### Assumption (Bounded data)

*The input data is bounded  $\|X\|_F \leq \sqrt{d_s}C_x$  with some positive constant  $C_x$ .*

- The embedding of token can have a unit norm [8] independent of  $d$ .

### Assumption

*The input data  $X$  has full row rank.*

- language tasks: added with positional embedding
- ViT: image grouped by patch and can be uncorrelated



## Assumptions on data

### Assumption (Bounded data)

The input data is bounded  $\|X\|_F \leq \sqrt{d_s} C_x$  with some positive constant  $C_x$ .

- The embedding of token can have a unit norm [8] independent of  $d$ .

### Assumption

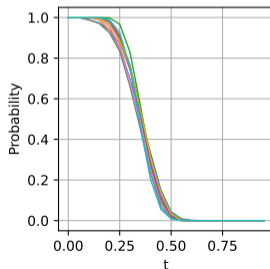
The input data  $X$  has full row rank.

- language tasks: added with positional embedding
- ViT: image grouped by patch and can be uncorrelated

### Assumption (different data have smaller similarity)

For any data pair  $(X_i, X_j)$ , with  $i \neq j$  and  $i, j \in [n]$ , then we assume that  $\mathbb{P} \left( \left| \langle X_i^\top X_i, X_j^\top X_j \rangle \right| \geq t \right) \leq \exp(-t^{\hat{c}})$  with some constant  $\hat{c} > 0$ .

- larger  $\hat{c} \Rightarrow$  more diverse data  $\Rightarrow$  more separable
- validated on a language IMDB dataset (sampling with 100 normalized sentences with embedding)



## Main results: Global convergence

Theorem (Under  $\tau_0 = d_m^{-1/2}$ )

Assume  $d_m \geq d$ , under LeCun/He (NTK) initialization and  $d_m = \tilde{\Omega}(n^3)$  ( $d_m = \tilde{\Omega}(n^2)$ ), with probability at least  $1 - 8e^{-\frac{d_m}{2}} - \delta - \exp(-\Omega(n-1)^{-\hat{c}}d_s^{-1})$  for proper  $\delta$ , choosing the step-size  $\gamma \leq \mathcal{O}(n^{-\frac{1}{2}})$ , then the GD training of the Transformer converges to a global minimum as follows:

$$\ell(\boldsymbol{\theta}^t) \leq \left(1 - \gamma \frac{\alpha^2}{2}\right)^t \ell(\boldsymbol{\theta}^0), \quad \text{for } t \geq 0. \quad (2)$$

## Main results: Global convergence

### Theorem (Under $\tau_0 = d_m^{-1/2}$ )

Assume  $d_m \geq d$ , under LeCun/He (NTK) initialization and  $d_m = \tilde{\Omega}(n^3)$  ( $d_m = \tilde{\Omega}(n^2)$ ), with probability at least  $1 - 8e^{-\frac{d_m}{2}} - \delta - \exp(-\Omega(n-1)^{-\hat{c}}d_s^{-1})$  for proper  $\delta$ , choosing the step-size  $\gamma \leq \mathcal{O}(n^{-\frac{1}{2}})$ , then the GD training of the Transformer converges to a global minimum as follows:

$$\ell(\boldsymbol{\theta}^t) \leq \left(1 - \gamma \frac{\alpha^2}{2}\right)^t \ell(\boldsymbol{\theta}^0), \quad \text{for } t \geq 0. \quad (2)$$

### Theorem (Under $\tau_0 = d_m^{-1}$ )

Under the NTK initialization, denote  $\boldsymbol{\Phi}^* := \frac{1}{d_s} [\mathbf{X}_1^\top \mathbf{1}_{d_s}, \dots, \mathbf{X}_n^\top \mathbf{1}_{d_s}]^\top \in \mathbb{R}^{n \times d}$ , the limiting kernel matrix will depend on  $\boldsymbol{\Phi}^*$ , and with  $d_m = \Omega(n)$ , the GD training of Transformer converges as Eq. (2).

## Main results: Global convergence

### Theorem (Under $\tau_0 = d_m^{-1/2}$ )

Assume  $d_m \geq d$ , under LeCun/He (NTK) initialization and  $d_m = \tilde{\Omega}(n^3)$  ( $d_m = \tilde{\Omega}(n^2)$ ), with probability at least  $1 - 8e^{-\frac{d_m}{2}} - \delta - \exp(-\Omega(n-1)^{-\hat{c}}d_s^{-1})$  for proper  $\delta$ , choosing the step-size  $\gamma \leq \mathcal{O}(n^{-\frac{1}{2}})$ , then the GD training of the Transformer converges to a global minimum as follows:

$$\ell(\boldsymbol{\theta}^t) \leq \left(1 - \gamma \frac{\alpha^2}{2}\right)^t \ell(\boldsymbol{\theta}^0), \quad \text{for } t \geq 0. \quad (2)$$

### Theorem (Under $\tau_0 = d_m^{-1}$ )

Under the NTK initialization, denote  $\Phi^* := \frac{1}{d_s} [X_1^\top \mathbf{1}_{d_s}, \dots, X_n^\top \mathbf{1}_{d_s}]^\top \in \mathbb{R}^{n \times d}$ , the limiting kernel matrix will depend on  $\Phi^*$ , and with  $d_m = \Omega(n)$ , the GD training of Transformer converges as Eq. (2).

**Remark:** 1) dimension missing: self-attention layer becomes  $XW_V^\top$

2)  $\tau_0 = d_m^{-1}$  and NTK initialization make Transformer

- ▶ enter into the lazy training regime easier
- ▶ require less over-parameterization requirement

## Proof framework

Polyak-Lojasiewicz (PL) inequality + Lipschitz continuous of gradient, defining  $\mathbf{F}_{\text{pre}} := \frac{\partial f(\mathbf{X})}{\partial \mathbf{w}_o} \in \mathbb{R}^{n \times d_m}$

$$\|\nabla \ell(\boldsymbol{\theta})\|_2^2 \geq 2\lambda_{\min}(\mathbf{F}_{\text{pre}}\mathbf{F}_{\text{pre}}^\top)\ell(\boldsymbol{\theta})$$

## Proof framework

Polyak-Lojasiewicz (PL) inequality + Lipschitz continuous of gradient, defining  $\mathbf{F}_{\text{pre}} := \frac{\partial f(\mathbf{X})}{\partial \mathbf{w}_o} \in \mathbb{R}^{n \times d_m}$

$$\|\nabla \ell(\boldsymbol{\theta})\|_2^2 \geq 2\lambda_{\min}(\mathbf{F}_{\text{pre}}\mathbf{F}_{\text{pre}}^\top)\ell(\boldsymbol{\theta})$$

$$\ell(\boldsymbol{\theta}^{t+1}) \leq \ell(\boldsymbol{\theta}^t) - \frac{\gamma}{2}\lambda_{\min}(\mathbf{F}_{\text{pre}}\mathbf{F}_{\text{pre}}^\top)\|\mathbf{f}^t - \mathbf{y}\|_2^2 \leq (1 - \frac{\gamma\alpha}{2})\ell(\boldsymbol{\theta}^t)$$

## Proof framework

Polyak-Lojasiewicz (PL) inequality + Lipschitz continuous of gradient, defining  $\mathbf{F}_{\text{pre}} := \frac{\partial \mathbf{f}(X)}{\partial \mathbf{w}_o} \in \mathbb{R}^{n \times d_m}$

$$\|\nabla \ell(\boldsymbol{\theta})\|_2^2 \geq 2\lambda_{\min}(\mathbf{F}_{\text{pre}}\mathbf{F}_{\text{pre}}^\top)\ell(\boldsymbol{\theta})$$

$$\ell(\boldsymbol{\theta}^{t+1}) \leq \ell(\boldsymbol{\theta}^t) - \frac{\gamma}{2}\lambda_{\min}(\mathbf{F}_{\text{pre}}\mathbf{F}_{\text{pre}}^\top)\|\mathbf{f}^t - \mathbf{y}\|_2^2 \leq (1 - \frac{\gamma\alpha}{2})\ell(\boldsymbol{\theta}^t)$$

### Lemma (minimum eigenvalue estimation)

Let  $\Phi = [X_1\boldsymbol{\beta}_{1,1}, X_2\boldsymbol{\beta}_{1,2}, \dots, X_n\boldsymbol{\beta}_{1,n}]^\top \in \mathbb{R}^{n \times d}$  where  $\boldsymbol{\beta}_{1,i}$  is the  $i$ -th output of softmax, then under our assumptions, we have

$$\eta_V/d_s \lesssim \lambda_0 := \lambda_{\min}(\mathbb{E}_{\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \eta_V I_d)}[\sigma_r(\Phi\mathbf{w})\sigma_r(\Phi\mathbf{w})^\top]) \lesssim \eta_V d_s \quad \text{w.h.p}$$

## Proof framework

Polyak-Lojasiewicz (PL) inequality + Lipschitz continuous of gradient, defining  $\mathbf{F}_{\text{pre}} := \frac{\partial \mathbf{f}(X)}{\partial \mathbf{w}_o} \in \mathbb{R}^{n \times d_m}$

$$\|\nabla \ell(\boldsymbol{\theta})\|_2^2 \geq 2\lambda_{\min}(\mathbf{F}_{\text{pre}}\mathbf{F}_{\text{pre}}^\top)\ell(\boldsymbol{\theta})$$

$$\ell(\boldsymbol{\theta}^{t+1}) \leq \ell(\boldsymbol{\theta}^t) - \frac{\gamma}{2}\lambda_{\min}(\mathbf{F}_{\text{pre}}\mathbf{F}_{\text{pre}}^\top)\|\mathbf{f}^t - \mathbf{y}\|_2^2 \leq (1 - \frac{\gamma\alpha}{2})\ell(\boldsymbol{\theta}^t)$$

### Lemma (minimum eigenvalue estimation)

Let  $\Phi = [X_1\boldsymbol{\beta}_{1,1}, X_2\boldsymbol{\beta}_{1,2}, \dots, X_n\boldsymbol{\beta}_{1,n}]^\top \in \mathbb{R}^{n \times d}$  where  $\boldsymbol{\beta}_{1,i}$  is the  $i$ -th output of softmax, then under our assumptions, we have

$$\eta_V/d_s \lesssim \lambda_0 := \lambda_{\min}(\mathbb{E}_{\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \eta_V I_d)}[\sigma_r(\Phi\mathbf{w})\sigma_r(\Phi\mathbf{w})^\top]) \lesssim \eta_V d_s \quad w.h.p$$

### Proof.

- ▶ Hermite expansion:  $\lambda_0 > \lambda_{\min}[\Phi\Phi^\top]$
- ▶ Gershgorin circle theorem:  $\lambda_{\min}[\Phi\Phi^\top] \geq \Omega(\|\boldsymbol{\beta}_{1,k}\|_2^2)$

□



## Discussion on $\alpha$

under LeCun initialization, we have  $\alpha^2 \geq d_m \lambda_0 / 4 \geq d_m \eta_V \mu(\sigma_r)^2 \Theta(\|\boldsymbol{\beta}_{1,k}\|_2^2)$

- ▶  $\tau = d_m^{-1/2}$ , we have  $\|\boldsymbol{\beta}_{1,k}\|_2^2 \geq 1/d_s$
- ▶  $\tau = d_m^{-1}$ , we have  $\|\boldsymbol{\beta}_{1,k}\|_2^2 \approx 1/d_s$

## Discussion on $\alpha$

under LeCun initialization, we have  $\alpha^2 \geq d_m \lambda_0 / 4 \geq d_m \eta_V \mu(\sigma_r)^2 \Theta(\|\boldsymbol{\beta}_{1,k}\|_2^2)$

▶  $\tau = d_m^{-1/2}$ , we have  $\|\boldsymbol{\beta}_{1,k}\|_2^2 \geq 1/d_s$

▶  $\tau = d_m^{-1}$ , we have  $\|\boldsymbol{\beta}_{1,k}\|_2^2 \approx 1/d_s$

different initializations:  $\alpha^2 \geq \tau_1^2 \eta_V d_m \Omega(1/d)$

▶ LeCun/He initialization:  $\alpha^2 \geq \Omega(d_m/d)$

▶ NTK initialization:  $\alpha^2 \geq \Omega(1/d)$

## Discussion on $\alpha$

under LeCun initialization, we have  $\alpha^2 \geq d_m \lambda_0 / 4 \geq d_m \eta_V \mu(\sigma_r)^2 \Theta(\|\beta_{1,k}\|_2^2)$

▶  $\tau = d_m^{-1/2}$ , we have  $\|\beta_{1,k}\|_2^2 \geq 1/d_s$

▶  $\tau = d_m^{-1}$ , we have  $\|\beta_{1,k}\|_2^2 \approx 1/d_s$

different initializations:  $\alpha^2 \geq \tau_1^2 \eta_V d_m \Omega(1/d)$

▶ LeCun/He initialization:  $\alpha^2 \geq \Omega(d_m/d)$

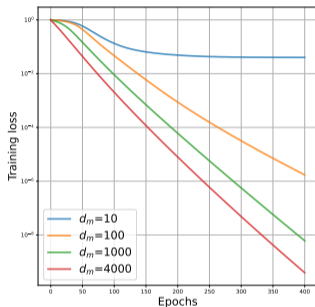
▶ NTK initialization:  $\alpha^2 \geq \Omega(1/d)$

architectures under LeCun initialization:

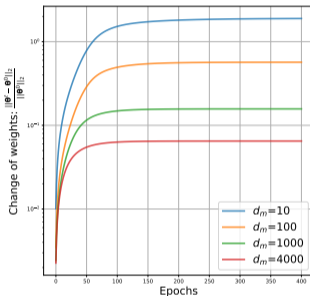
▶ self-attention + two-layer ReLU NN:  $\Omega(n^3)$  over-parameterization

▶ three-layer ReLU NN:  $\Omega(n^3)$  over-parameterization

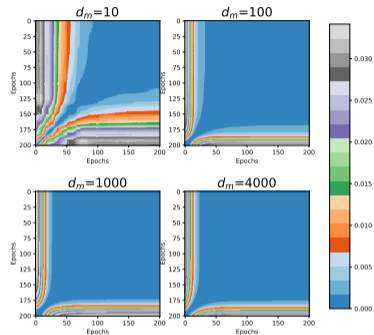
## Experimental validations (width matters)



(a) Convergence curve.



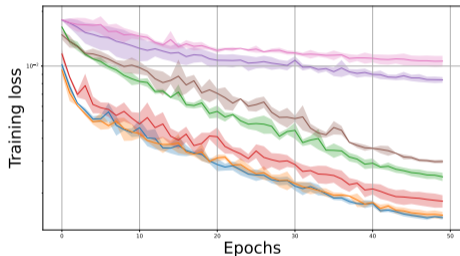
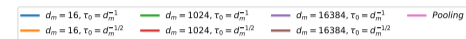
(b) Weight movement.



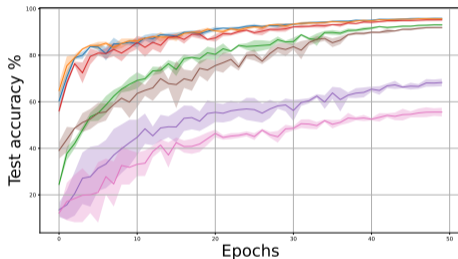
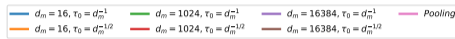
(c) Kernel distance.

**Figure:** Visualization of the training process of Transformers with LeCun initialization and  $\tau_0 = d_m^{-1/2}$  scaling on synthetic data. (a) Linear convergence. (b) Rate of change of the weights during training. Observe that the weights change very slowly after the 50<sup>th</sup> epoch. (c) Evolution of the NTK during the training. The result mirrors the plot (b) and demonstrates how the kernel varies significantly at the beginning of the training and remains approximately constant later. As the width increases, the empirical NTK becomes more stable.

# Separation between $d_m^{-1}$ and $d_m^{-1/2}$



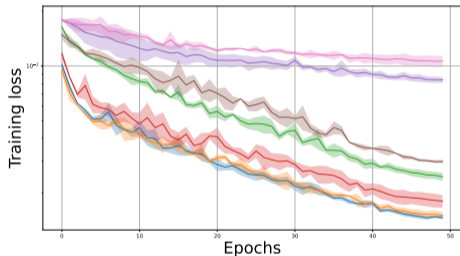
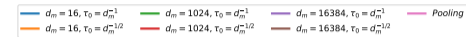
(a) Training loss



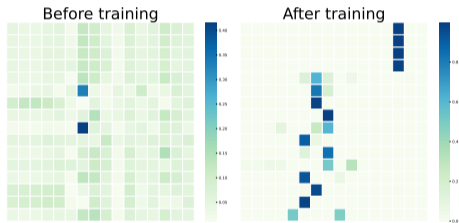
(b) Accuracy

Figure: Results on MNIST dataset trained by ViT with different scaling schemes.

# Separation between $d_m^{-1}$ and $d_m^{-1/2}$



(a) Training loss



(b) Attention map,  $d_m = 16384$ .

Figure: Results on MNIST dataset trained by ViT with different scaling schemes.

## Conclusion

- ▶ scaling factor  $\tau_0$ :  $d_m^{-1/2}$  vs.  $d_m^{-1}$
- ▶ initializations: LeCun/He vs. NTK

# Conclusion

- ▶ scaling factor  $\tau_0$ :  $d_m^{-1/2}$  vs.  $d_m^{-1}$
- ▶ initializations: LeCun/He vs. NTK

## Future direction

- ▶ Architecture: benefits of attention
- ▶ Optimization objective: implicit bias
- ▶ Application: in-context learning, chain-of-thought reasoning



Thanks for your attention!

Q & A

my homepage [www.lfhsgre.org](http://www.lfhsgre.org) for more information!

## References I

- [1] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby.  
An image is worth 16x16 words: Transformers for image recognition at scale.  
2021.  
(Cited on page 3.)
- [2] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin.  
Attention is all you need.  
In *Advances in Neural Information Processing Systems*, pages 5998–6008, 2017.  
(Cited on pages 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, and 13.)
- [3] Arthur Jacot, Franck Gabriel, and Clément Hongler.  
Neural tangent kernel: Convergence and generalization in neural networks.  
In *Advances in Neural Information Processing Systems*, pages 8571–8580, 2018.  
(Cited on pages 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, and 16.)
- [4] Jiri Hron, Yasaman Bahri, Jascha Sohl-Dickstein, and Roman Novak.  
Infinite attention: NNGP and NTK for deep attention networks.  
2020.  
(Cited on pages 11, 12, and 13.)

## References II

- [5] Greg Yang.  
Tensor programs II: Neural tangent kernel for any architecture.  
*arXiv preprint arXiv:2006.14548*, 2020.  
(Cited on pages 11, 12, and 13.)
- [6] Lenaïc Chizat, Edouard Oyallon, and Francis Bach.  
On lazy training in differentiable programming.  
In *Advances in Neural Information Processing Systems*, pages 2933–2943, 2019.  
(Cited on pages 14, 15, and 16.)
- [7] Tao Luo, Zhi-Qin John Xu, Zheng Ma, and Yaoyu Zhang.  
Phase diagram for two-layer relu neural networks at infinite-width limit.  
*Journal of Machine Learning Research*, 22(71):1–47, 2021.  
(Cited on pages 14, 15, and 16.)
- [8] Hongkang Li, Meng Wang, Sijia Liu, and Pin-Yu Chen.  
A theoretical understanding of shallow vision transformers: Learning, generalization, and sample complexity.  
2023.  
(Cited on pages 23, 24, and 25.)