

# In-Context Convergence of Transformers

Yu Huang<sup>1</sup> Yuan Cheng<sup>2</sup> Yingbin Liang<sup>3</sup>

<sup>1</sup>*University of Pennsylvania*

<sup>2</sup>*National University of Singapore*

<sup>3</sup>*The Ohio State University*



# Transformers have revolutionized many domains

Attention is all you need. Vaswani et al, (2017)

PU what are transformers

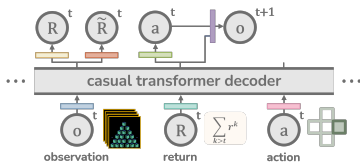
Transformers are a type of deep learning architecture that have revolutionized the field of natural language processing (NLP) and have also found applications in computer vision and other domains. Here's a concise breakdown of transformers:

## Large Language Models (LLMs)



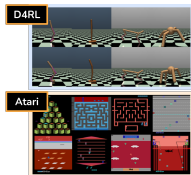
Computer Vision

DALL-E3



## Reinforcement Learning

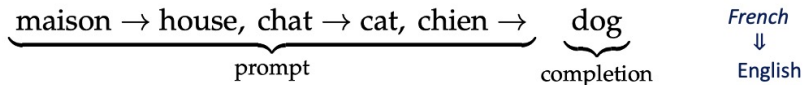
Wu et al. (2023)



# Remarkable emergent ability for LLMs: In-Context learning

---

Given a prompt containing in-context examples, pre-trained LLM responds to new query token appropriately *without further fine-tuning*.

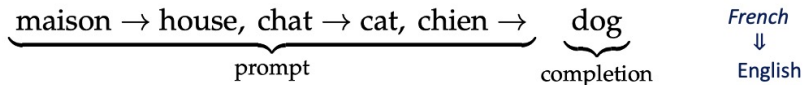


This figure is from Garg et al., (2022)

# Remarkable emergent ability for LLMs: In-Context learning

---

Given a prompt containing in-context examples, pre-trained LLM responds to new query token appropriately *without further fine-tuning*.



This figure is from Garg et al., (2022)

Learning **without** weight updates

# Towards understanding ICL with transformers

---

## Empirical evidence

- Transformers can **in-context** learn **linear** function (Garg et al, 2022)
  - ▶ Sample many  $f \in \mathcal{F}$ , and construct corresponding prompts:

$$(x_1, f(x_1), \dots, x_N, f(x_N), x_{\text{query}})$$

- ▶ Train transformer to predict  $f(x_{\text{query}})$
- ▶ For a new  $f'$  and its prompt: the trained model (**without finetuning**) can predict  $f'(x_{\text{query}})$

# Towards understanding ICL with transformers

---

## Empirical evidence

- Transformers can **in-context** learn **linear** function (Garg et al, 2022)
  - ▶ Sample many  $f \in \mathcal{F}$ , and construct corresponding prompts:

$$(x_1, f(x_1), \dots, x_N, f(x_N), x_{\text{query}})$$

- ▶ Train transformer to predict  $f(x_{\text{query}})$
- ▶ For a new  $f'$  and its prompt: the trained model (**without finetuning**)  
**can** predict  $f'(x_{\text{query}})$

## Going forward **theoretically**

- Most focus on Expressive power or Generalization:  
Oswald et al., (2023); Bai et al., (2023); Li et al., (2023)

# Towards understanding ICL with transformers

---

## Empirical evidence

- Transformers can **in-context** learn **linear** function (Garg et al, 2022)
  - ▶ Sample many  $f \in \mathcal{F}$ , and construct corresponding prompts:

$$(x_1, f(x_1), \dots, x_N, f(x_N), x_{\text{query}})$$

- ▶ Train transformer to predict  $f(x_{\text{query}})$
- ▶ For a new  $f'$  and its prompt: the trained model (**without finetuning**) can predict  $f'(x_{\text{query}})$

## Going forward **theoretically**

- Most focus on Expressive power or Generalization:  
Oswald et al., (2023); Bai et al., (2023); Li et al., (2023)
- Training dynamics of **linear** attention:  
Zhang et al., (2023); Mahankali et al., (2023); Ahn et al., (2023)

# How do *softmax*-based transformers trained via gradient descent learn in-context?

- Transformers are based on **softmax** attention mechanism.

$$\text{softmax}\left(\frac{\begin{matrix} \mathbf{Q} \\ \text{2x2 grid} \end{matrix} \times \begin{matrix} \mathbf{K}^T \\ \text{2x2 grid} \end{matrix}}{\sqrt{d_k}}\right) \begin{matrix} \mathbf{V} \\ \text{2x2 grid} \end{matrix} = \begin{matrix} \mathbf{Z} \\ \text{2x2 grid} \end{matrix}$$





**Our contribution:** first step towards in-context learning dynamics of the 1-layer softmax transformer

- Transformers are based on **softmax** attention mechanism.

$$\text{softmax}\left(\frac{\begin{matrix} \mathbf{Q} \\ \text{3x3 grid} \end{matrix} \times \begin{matrix} \mathbf{K}^T \\ \text{3x3 grid} \end{matrix}}{\sqrt{d_k}}\right) \begin{matrix} \mathbf{V} \\ \text{3x3 grid} \end{matrix} = \begin{matrix} \mathbf{Z} \\ \text{3x3 grid} \end{matrix}$$



# ICL framework

---

- **Prompt:**  $P = (x_1, f(x_1), \dots, x_N, f(x_N), x_{\text{query}})$ :
  - ▶ **Linear task:**  $f(x) = \langle w, x \rangle, w \sim \mathcal{D}_\Omega$
  - ▶ **IID data:**  $\{x_i\} \cup \{x_{\text{query}}\} \stackrel{i.i.d.}{\sim} D_{\mathcal{X}}$
- **Goal:** Predict  $\hat{y}_{\text{query}} \approx f(x_{\text{query}})$ .

# ICL framework

---

- **Prompt:**  $P = (x_1, f(x_1), \dots, x_N, f(x_N), x_{\text{query}})$ :
  - ▶ **Linear task:**  $f(x) = \langle w, x \rangle$ ,  $w \sim \mathcal{D}_\Omega$
  - ▶ **IID data:**  $\{x_i\} \cup \{x_{\text{query}}\} \stackrel{i.i.d.}{\sim} \mathcal{D}_\mathcal{X}$
- **Goal:** Predict  $\hat{y}_{\text{query}} \approx f(x_{\text{query}})$ .
- **Task Distribution  $\mathcal{D}_\Omega$ :**
  - ▶  $w \stackrel{i.i.d.}{\sim} \mathcal{D}_\Omega$  with zero mean and covariance  $\mathbf{I}_{d \times d}$

# ICL framework

- **Prompt:**  $P = (x_1, f(x_1), \dots, x_N, f(x_N), x_{\text{query}})$ :
  - ▶ **Linear task:**  $f(x) = \langle w, x \rangle, w \sim \mathcal{D}_\Omega$
  - ▶ **IID data:**  $\{x_i\} \cup \{x_{\text{query}}\} \stackrel{i.i.d.}{\sim} \mathcal{D}_\mathcal{X}$
- **Goal:** Predict  $\hat{y}_{\text{query}} \approx f(x_{\text{query}})$ .

- **Task Distribution  $\mathcal{D}_\Omega$ :**

- ▶  $w \stackrel{i.i.d.}{\sim} \mathcal{D}_\Omega$  with zero mean and covariance  $\mathbf{I}_{d \times d}$

- **Data Distribution  $\mathcal{D}_\mathcal{X}$ :**

- ▶  $K$  distinct features:

$$v_k \in \mathbb{R}^d, \|v_k\| = 1 \text{ for } k \in [K], v_i \perp v_j \text{ for } i \neq j$$

- ▶  $x = v_k$  with prob  $p_k$ , where  $p_k \in (0, 1)$  and  $\sum_{k \in [K]} p_k = 1$ .

# Transformer Architecture

---

- Embeddings

$$E = E(P) = \begin{pmatrix} x_1 & x_2 & \cdots & x_N & x_{\text{query}} \\ y_1 & y_2 & \cdots & y_N & 0 \end{pmatrix} \in \mathbb{R}^{(d+1) \times (N+1)}.$$

# Transformer Architecture

- Embeddings

$$E = E(P) = \begin{pmatrix} x_1 & x_2 & \cdots & x_N & x_{\text{query}} \\ y_1 & y_2 & \cdots & y_N & 0 \end{pmatrix} \in \mathbb{R}^{(d+1) \times (N+1)}.$$

- One-layer transformer:

- ▶ Self-attention mechanism:  $W^V E \cdot \text{softmax} \left( (W^K E)^\top W^Q E \right)$ ,
- ▶ Mask:  $W^V M(E)$ ,  $W^K M(E)$
- ▶ Reparameterization:  $\theta = (\nu, Q)$  (Anh et al., 2023, Zhang et al., 2023)

$$W^V = \begin{pmatrix} 0_{d \times d} & 0_d \\ 0_d^\top & \nu \end{pmatrix}, \quad W^{KQ} = \begin{pmatrix} Q & 0_d \\ 0_d^\top & 0 \end{pmatrix}$$

- Consolidate Key and Query
- Fixed  $\nu = 1$

- Nealy no loss of optimality!

# Transformer Architecture

---

- **Output:**

$$\begin{aligned}\widehat{y}_{\text{query}}^{(t)} &= [M(E^y) \cdot \text{softmax} \left( M(E^x)^\top Q^{(t)} E^x \right)]_{N+1} \\ &= \sum_{i \in [N]} \mathbf{attn}_i^{(t)} y_i = \sum_{k \in [K]} \mathbf{Attn}_k^{(t)} \langle w, v_k \rangle.\end{aligned}$$

# Transformer Architecture

---

- **Output:**

$$\begin{aligned}\widehat{y}_{\text{query}}^{(t)} &= [M(E^y) \cdot \text{softmax} \left( M(E^x)^\top Q^{(t)} E^x \right)]_{N+1} \\ &= \sum_{i \in [N]} \mathbf{attn}_i^{(t)} y_i = \sum_{k \in [K]} \mathbf{Attn}_k^{(t)} \langle w, v_k \rangle.\end{aligned}$$

- for the  $i$ -th token:  $\mathbf{attn}_i^{(t)} = \frac{e^{E_i^x \top Q^{(t)} E_{N+1}^x}}{\sum_{j \in [N]} e^{E_j^x \top Q^{(t)} E_{N+1}^x}}$ .
- for the  $k$ -th features:  $\mathbf{Attn}_k^{(t)} = \sum_{i \in [N], x_i = v_k} \mathbf{attn}_i^{(t)}$ .



# Training Settings

---

- **Loss Function:**

$$L(\theta) = \frac{1}{2} \mathbb{E}_{w \sim \mathcal{D}_\Omega, \{x_i\}_{i=1}^N \cup \{x_{\text{query}}\} \sim \mathcal{D}_X^{N+1}} \left[ (\hat{y}_{\text{query}} - \langle w, x_{\text{query}} \rangle)^2 \right]$$

- **Training Algorithm:**  $Q^{(0)}$  initialize as  $\mathbf{0}_{d \times d}$ , with GD update.
- **Prediction Error:**

$$\mathcal{L}_k(\theta) = \frac{1}{2} \mathbb{E} \left[ (\hat{y}_{\text{query}} - \langle w, x_{\text{query}} \rangle)^2 \mid x_{\text{query}} = v_k \right].$$

*performance measure*

# ICL with imbalanced features

**Imbalanced Cases:** One dominant feature  $v_1$ :  $p_1 = \Theta(1)$ ;  
Under-represented  $v_k$ :  $p_k = \Theta\left(\frac{1}{K}\right)$ .

## Theorem (Prediction Error Converges (Informal))

For  $0 < \epsilon < 1$ ,  $N \geq \text{poly}(K)$ ,  $\text{polylog}(K) \gg \log\left(\frac{1}{\epsilon}\right)$ , prediction error:

1. **Dominant feature  $v_1$ :** with at most  $T_1 = O\left(\frac{\log(\epsilon^{-1/2})}{\eta\epsilon}\right)$  GD iterations,  
 $\mathcal{L}_1(\theta^{(T_1)}) \leq \mathcal{L}_1^* + \epsilon$ .

2. **Under-represented features  $v_k$ :** with at most

$T_k = O\left(\frac{\log(K)K^2}{\eta} + \frac{K \log\left(K\epsilon^{-\frac{1}{2}}\right)}{\epsilon\eta}\right)$  GD iterations,  $\mathcal{L}_k(\theta^{(T_k)}) \leq \mathcal{L}_k^* + \epsilon$ .

# ICL with imbalanced features

**Imbalanced Cases:** One dominant feature  $v_1$ :  $p_1 = \Theta(1)$ ;  
Under-represented  $v_k$ :  $p_k = \Theta\left(\frac{1}{K}\right)$ .

## Theorem (Prediction Error Converges (Informal))

For  $0 < \epsilon < 1$ ,  $N \geq \text{poly}(K)$ ,  $\text{polylog}(K) \gg \log\left(\frac{1}{\epsilon}\right)$ , prediction error:

1. **Dominant feature  $v_1$ :** with at most  $T_1 = O\left(\frac{\log(\epsilon^{-1/2})}{\eta\epsilon}\right)$  GD iterations,  
 $\mathcal{L}_1(\theta^{(T_1)}) \leq \mathcal{L}_1^* + \epsilon$ .

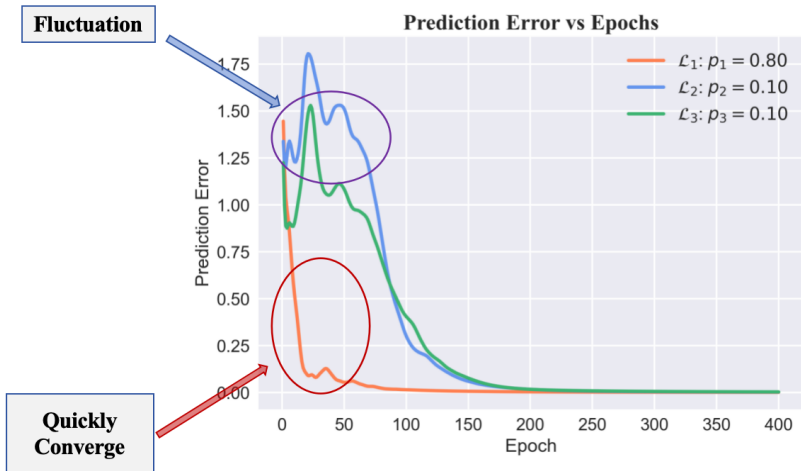
2. **Under-represented features  $v_k$ :** with at most

$T_k = O\left(\frac{\log(K)K^2}{\eta} + \frac{K \log\left(K\epsilon^{-\frac{1}{2}}\right)}{\epsilon\eta}\right)$  GD iterations,  $\mathcal{L}_k(\theta^{(T_k)}) \leq \mathcal{L}_k^* + \epsilon$ .

- Global optimal:  $\mathcal{L}_k^* = \Theta(e^{-\text{poly}(K)})$ .
- Nearly optimal prediction error for both **under-represented features** and the dominant feature.
- **Stage-wise Convergence!**

# ICL with imbalanced features

## Stage-wise Convergence



# ICL with imbalanced features

---

## Theorem (Attention score concentrates (Informal))

For any  $0 < \epsilon < 1$ ,  $N \geq \text{poly}(K)$ ,  $\text{polylog}(K) \gg \log(\frac{1}{\epsilon})$ , for attention score,  $x_{\text{query}} = v_k$ , after  $T_k$ , w.h.p

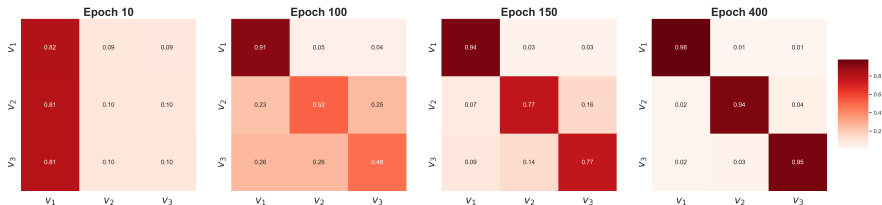
$$(1 - \text{Attn}_k^{(T_k)})^2 \leq O(\epsilon).$$

# ICL with imbalanced features

## Theorem (Attention score concentrates (Informal))

For any  $0 < \epsilon < 1$ ,  $N \geq \text{poly}(K)$ ,  $\text{polylog}(K) \gg \log(\frac{1}{\epsilon})$ , for attention score,  $x_{\text{query}} = v_k$ , after  $T_k$ , w.h.p

$$(1 - \text{Attn}_k^{(T_k)})^2 \leq O(\epsilon).$$



Attention scores Heatmap

# ICL with imbalanced features

## Theorem (Attention score concentrates (Informal))

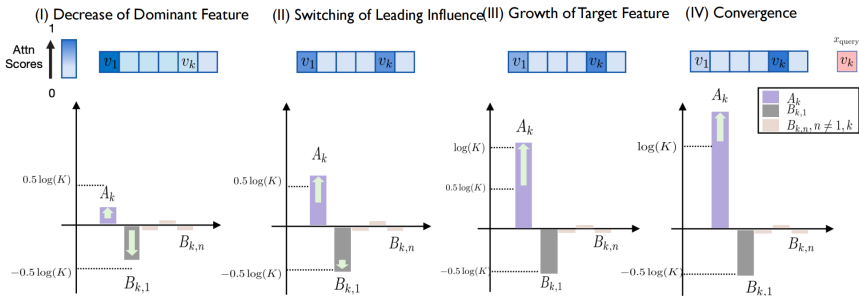
For any  $0 < \epsilon < 1$ ,  $N \geq \text{poly}(K)$ ,  $\text{polylog}(K) \gg \log(\frac{1}{\epsilon})$ , for attention score,  $x_{\text{query}} = v_k$ , after  $T_k$ , w.h.p

$$(1 - \text{Attn}_k^{(T_k)})^2 \leq O(\epsilon).$$

**In-context Ability:** given a test prompt from any new task  $w$  (possibly unseen), model can still well approximate test query.

$$y_{\text{query}}^{(T^*)} = \text{Attn}_k^{(T^*)} \langle w, v_k \rangle + \sum_{m \neq k} \text{Attn}_m^{(T^*)} \langle w, v_m \rangle \approx \langle w, v_k \rangle.$$

# Four-phase Behavior of Under-presented Features



Four-phase learning dynamics of under-represented features

*Bilinear attention weight*

$A_k$  : weight of query token and its target feature  
 $B_{k,n}$  : weight of query token and off-target features



## Concluding remarks

---

- Analyzing the training dynamics of a one-layer transformer with **softmax** attention trained by GD for in-context learning.

# Concluding remarks

---

- Analyzing the training dynamics of a one-layer transformer with **softmax** attention trained by GD for in-context learning.
- **Take away message:**
  - ▶ Stage-wise convergence
  - ▶ Attention concentration → In-context ability..
  - ▶ Novel analysis of **phase decomposition**.

Thanks & Questions