

Position: Do pretrained transformers learn in-context by Gradient Descent?

Lingfeng Shen*, **Aayush Mishra***, Daniel Khashabi

I have a dream that one day ...



LLM



this nation will rise up ...

In-Context Learning (ICL)

Input: JHU Output: Baltimore
Input: IITD Output: Delhi
Input: NYU Output:



LLM



New York

Good evening → Guten Abend
Vienna is great → Wien ist großartig
Where is the next ICML? →



LLM



Wo ist das nächste ICML?

Language

Transformers Learn Nonlinear Features In Context: Nonconvex Mean-field Dynamics on the Attention Landscape

Damai Dai^{†,*}

Juno Kim^{1,2} Taiji Suzuki^{1,2}

WHY CAN GPT LEARN IN-CONTEXT? INVESTIGATIONS WITH LINEAR MODELS

Ekin Akyürek^{1,2,a} Dale Schuurmans¹ Jacob Andreas^{*2} Tengyu Ma^{*1,3,b} Denny Zhou^{*1}

Transformers Learn In-Context by Gradient Descent

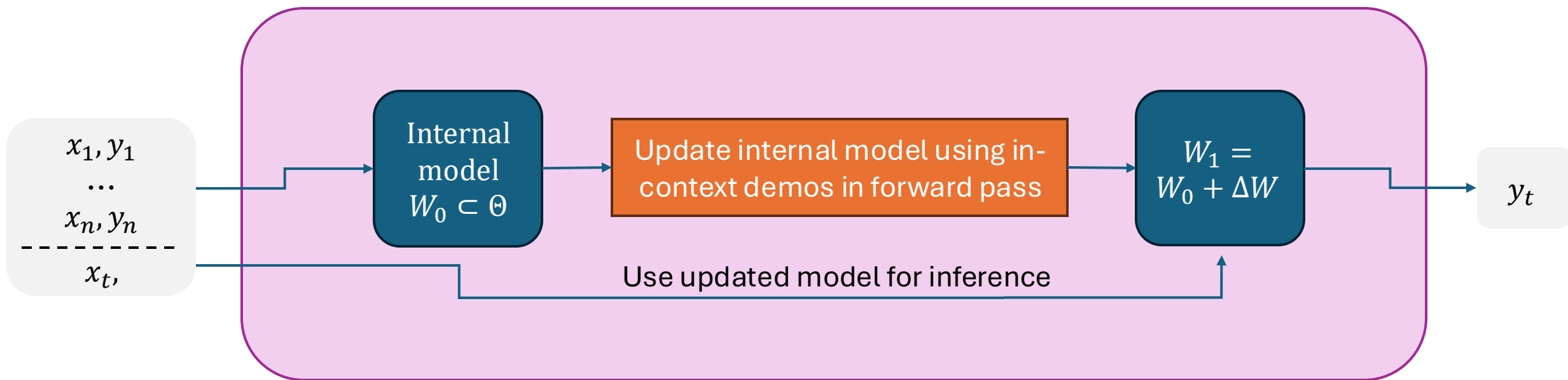
Johannes von Oswald^{1,2} Eyvind Niklasson² Ettore Randazzo² João Sacramento¹
Alexander Mordvintsev² Andrey Zhmoginov² Max Vladymyrov²

1. Dai, Damai, et al. "Why can gpt learn in-context? language models implicitly perform gradient descent as meta-optimizers." arXiv preprint arXiv:2212.10559 (2022).
2. Akyürek, Ekin, et al. "What learning algorithm is in-context learning? investigations with linear models." arXiv preprint arXiv:2211.15661 (2022).
3. Von Oswald, Johannes, et al. "Transformers learn in-context by gradient descent." International Conference on Machine Learning. PMLR, 2023.
4. Kim, Juno, and Taiji Suzuki. "Transformers learn nonlinear features in context: Nonconvex mean-field dynamics on the attention landscape." arXiv preprint arXiv:2402.01258 (2024).

The argument

LLM
(weights Θ)

The argument



But here's the thing ...

These theories make
oversimplified and sometimes
unrealistic assumptions.

The functional nature of ICL
(and its equivalence to GD) **remains unclear.**

Evolution of this theory

1. *In-context Learning can be **interpreted** as implicit finetuning.* [Dai+, 2022]

Show that transformer attention has a dual form of gradient descent:

$$\mathcal{F}(\mathbf{x}) = (\mathbf{W}_0 + \Delta\mathbf{W})\mathbf{x} = \mathbf{W}_0\mathbf{x} + \text{LinearAttn}(\mathbf{E}, \mathbf{X}', \mathbf{x})$$

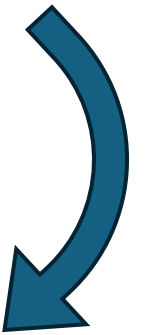
2. *Hand crafted transformer weights that simulate gradient descent.* [Akyurek+, 2022]

These weights can **imitate** GD in the forward pass of the transformer.

3. *Weights found by optimization match the construction.* [Oswald+, 2023]

Compare actual trained weights with their new construction.

Claim: Gradient-based optimization and attention-based in-context learning are **equivalent**.



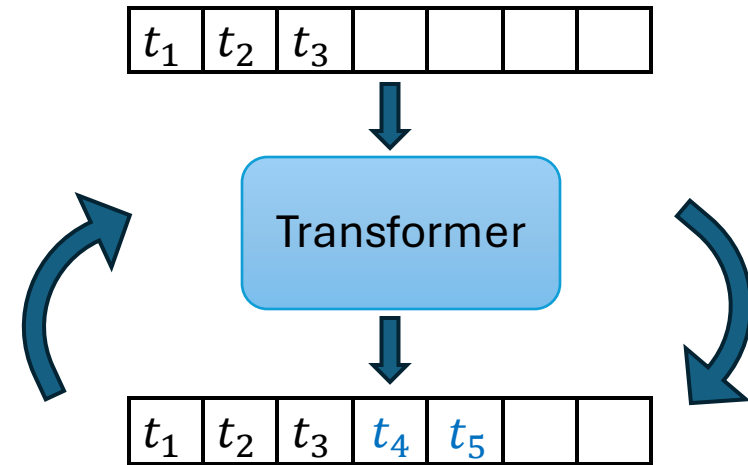
The “ICL Objective” problem

ICL

Given an unsupervised corpus of tokens $\{t_1, t_2, \dots, t_n\}$, causal language modeling (CLM) objective is used to train the model (with a context window of size k):

$$\arg \max_{\Theta} \sum_i \log P(t_i | t_{i-k}, \dots, t_{i-1}; \Theta)$$

- Models trained on **unstructured** sequences
- **Emergent** phenomenon



Transformers, pretrained with **CLM objective**, yield **emergent** ICL.

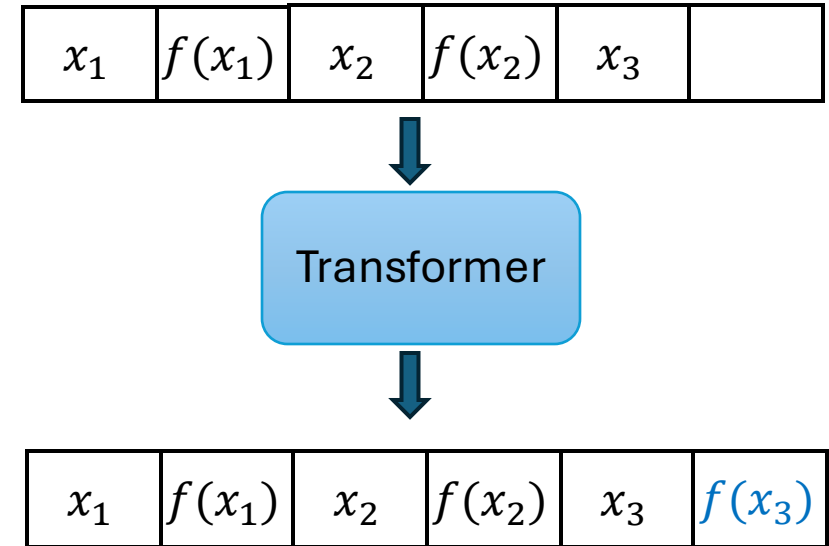
Our new terminology

$\widehat{\text{ICL}}$

Given an input domain $x \sim X$, and a function class $f \sim F$, **ICL objective** is used to train the model by giving it structured paired inputs:

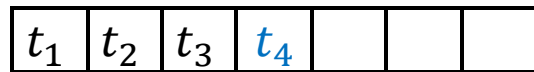
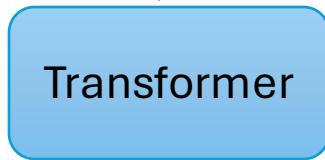
$$\arg \max_{\theta} \sum \log P(f(x_{n+1}) | x_1, f(x_1), \dots, x_n \circ f(x_n) \circ x_{n+1}; \theta)$$

- Models trained on **structured** sequences
- **Non-Emergent** phenomenon



Transformers trained with **ICL objective**, yield **non-emergent** meta learning.

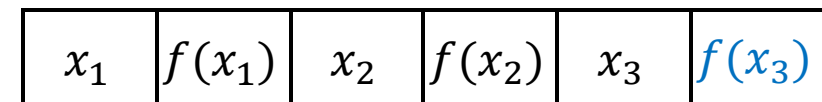
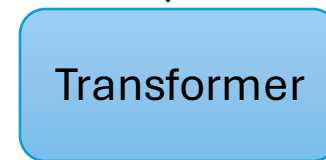
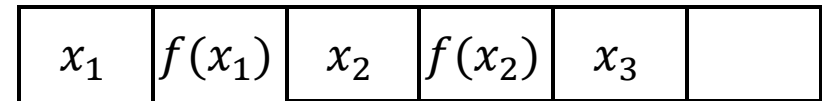
ICL



Emergent

\neq

$\widehat{\text{ICL}}$



Non-emergent

What they
imply

$\text{ICL} \approx \text{GD equivalence}$

For any Transformer weights resulting from self-supervised pretraining and **for any** well-defined task, ICL is algorithmically equivalent to GD.

What they
show

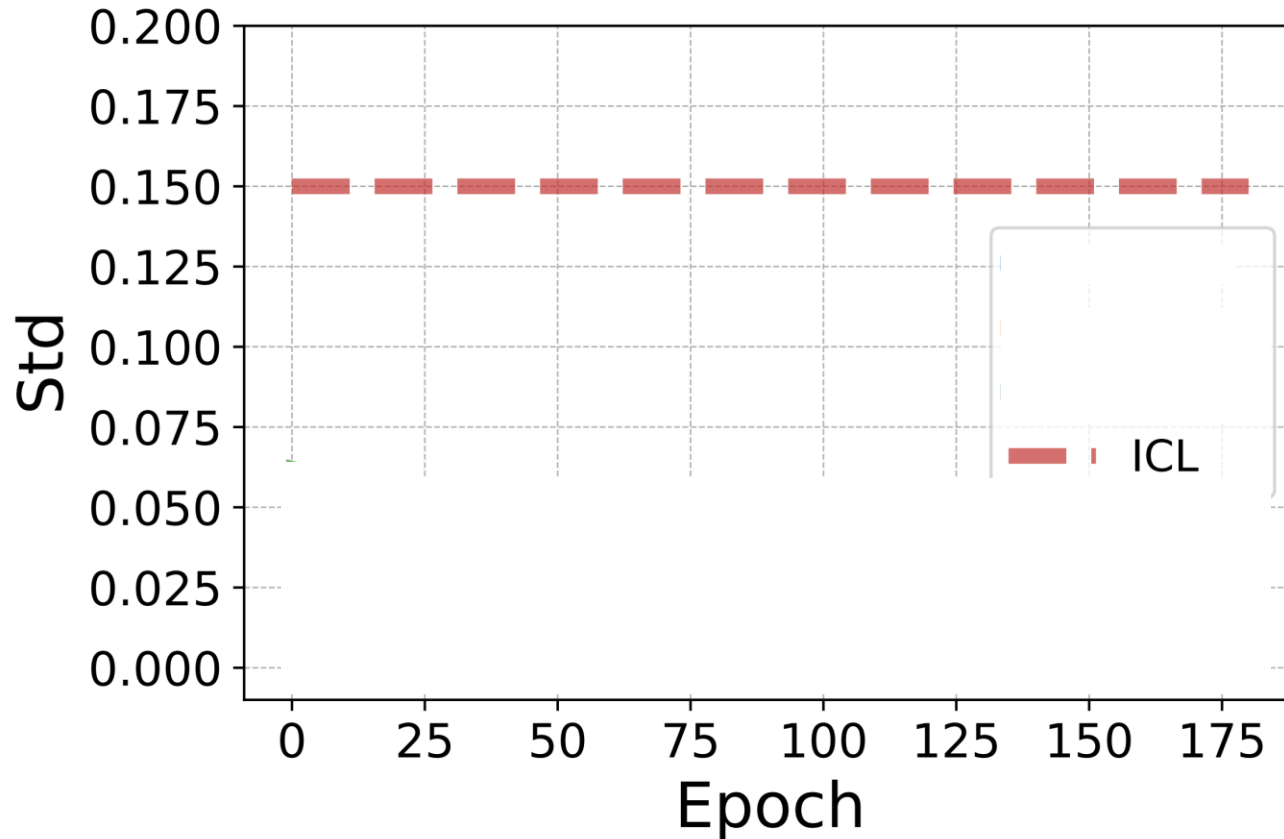
$\widehat{\text{ICL}} \approx \text{GD equivalence}$

For a **given** well-defined task, **there exist** Transformer weights such that $\widehat{\text{ICL}}$ is algorithmically equivalent to GD.

Transformers have the **expressive capacity** to simulate GD.

This **does not imply** that LLMs **actually do** simulate it.

Evidence against $ICL \approx GD$: Order sensitivity



- ICL is known to be highly order-sensitive [Lu+].
- GD is order-insensitive.
contradicts [Oswald+]
- Variants of GD are still not as sensitive as ICL.
undermines [Akyurek+]

ICL is likely **not** equivalent to GD based on order sensitivity.

Evidence against $ICL \approx GD$: Weight Sparsity

- Hand constructed weights and inputs in [2] are **highly sparse**.
- Constructions of [3] are similarly sparse.

$$H^{(0)} = \begin{bmatrix} \cdots & \boxed{0} & y_i & \boxed{0} & \cdots \end{bmatrix} \quad W_e = \begin{pmatrix} \boxed{I^{(d+1) \times (d+1)}} & \boxed{0} \\ \boxed{0} & \boxed{0} \end{pmatrix}$$

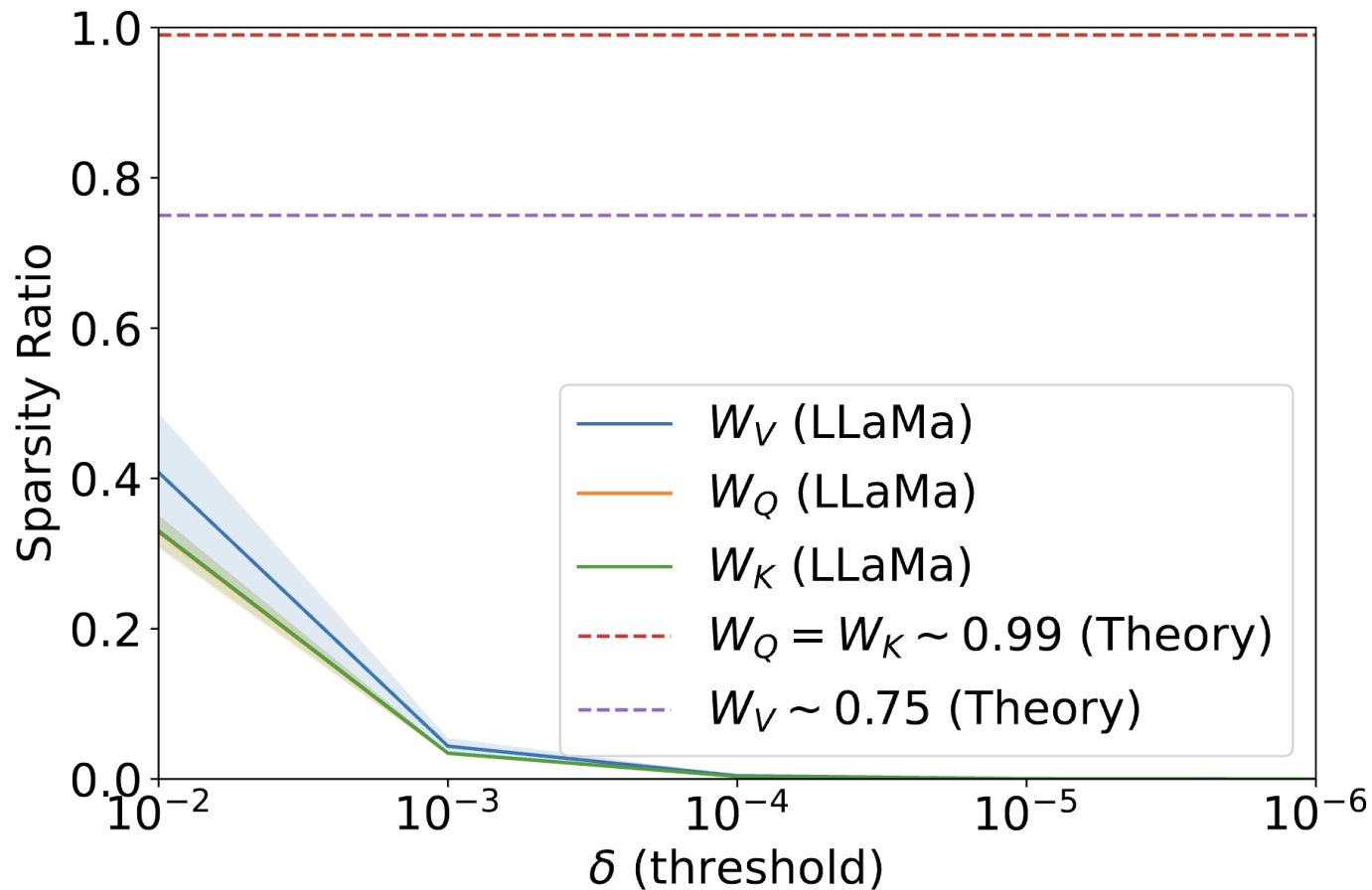
$$W_V = \begin{pmatrix} \boxed{0} & \boxed{0} \\ W_0 & \boxed{-I_y} \end{pmatrix}$$

$$W_K^l = \begin{pmatrix} \boxed{0} & \cdots \\ \vdots & \vdots \\ \boxed{I^{p \times p}} & \boxed{0^{p \times p}} \\ \vdots & \vdots \end{pmatrix} \quad W_Q^l = \begin{pmatrix} \boxed{0} & \cdots \\ \vdots & \vdots \\ \boxed{0^{p \times p}} & \boxed{I^{p \times p}} \\ \vdots & \vdots \end{pmatrix}$$

$$W_K = W_Q = \begin{pmatrix} \boxed{I_x} & \boxed{0} \\ \boxed{0} & \boxed{0} \end{pmatrix}$$

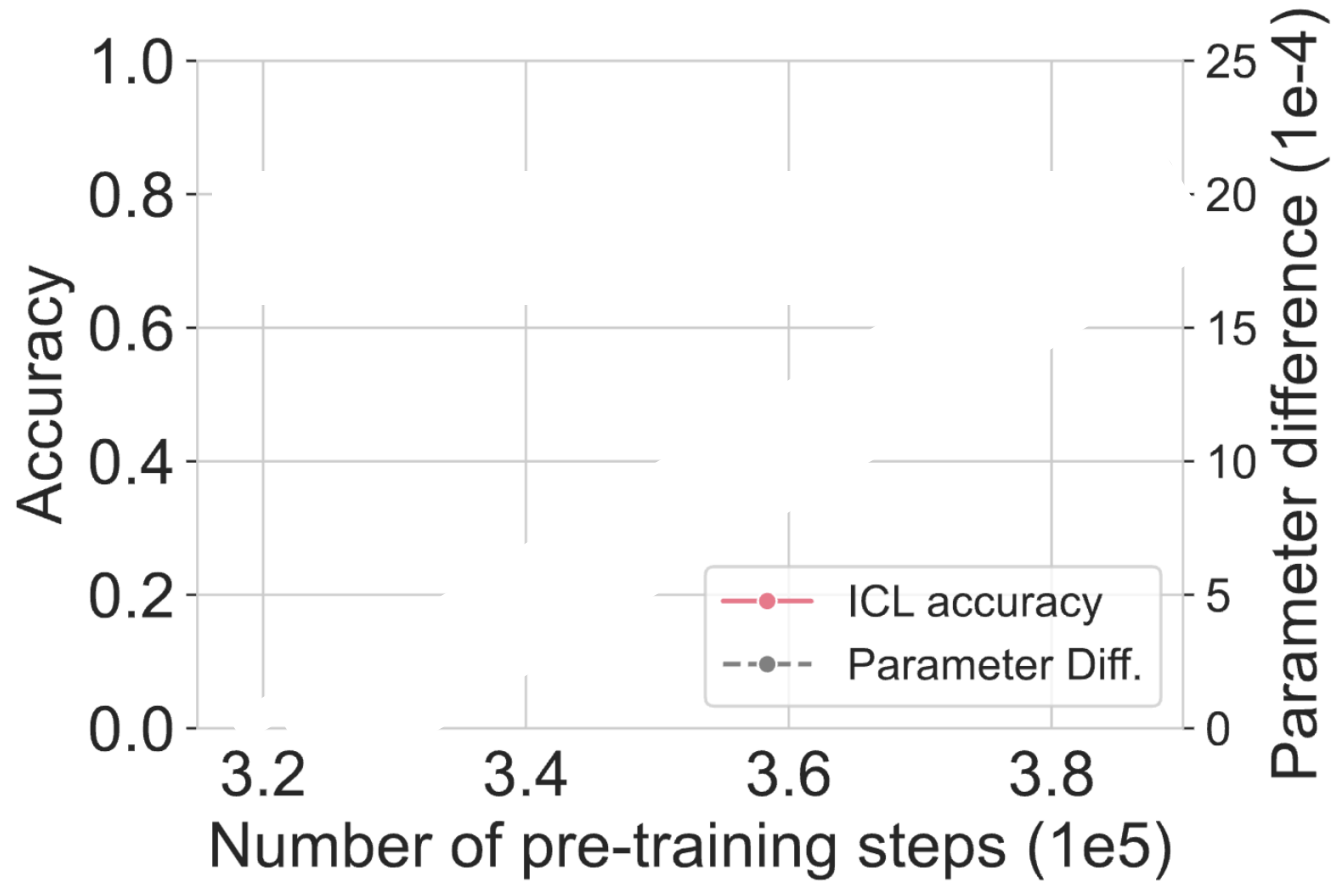
$$P = \frac{\eta}{N} \boxed{I}$$

Evidence against $ICL \approx GD$: Weight Sparsity



Real LLMs are rather dense.

Evidence against $ICL \approx GD$: Weight Evolution



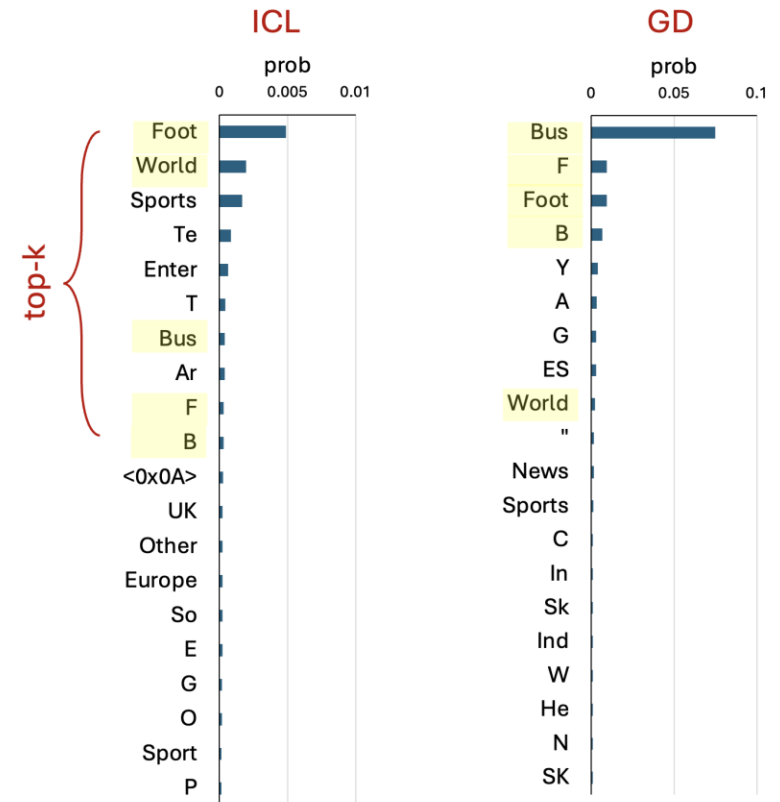
ICL ability **remains stable** even when weights keep evolving.

To claim $ICL \approx GD$, showing it for **a single sparse choice of parameters** is **not** enough.

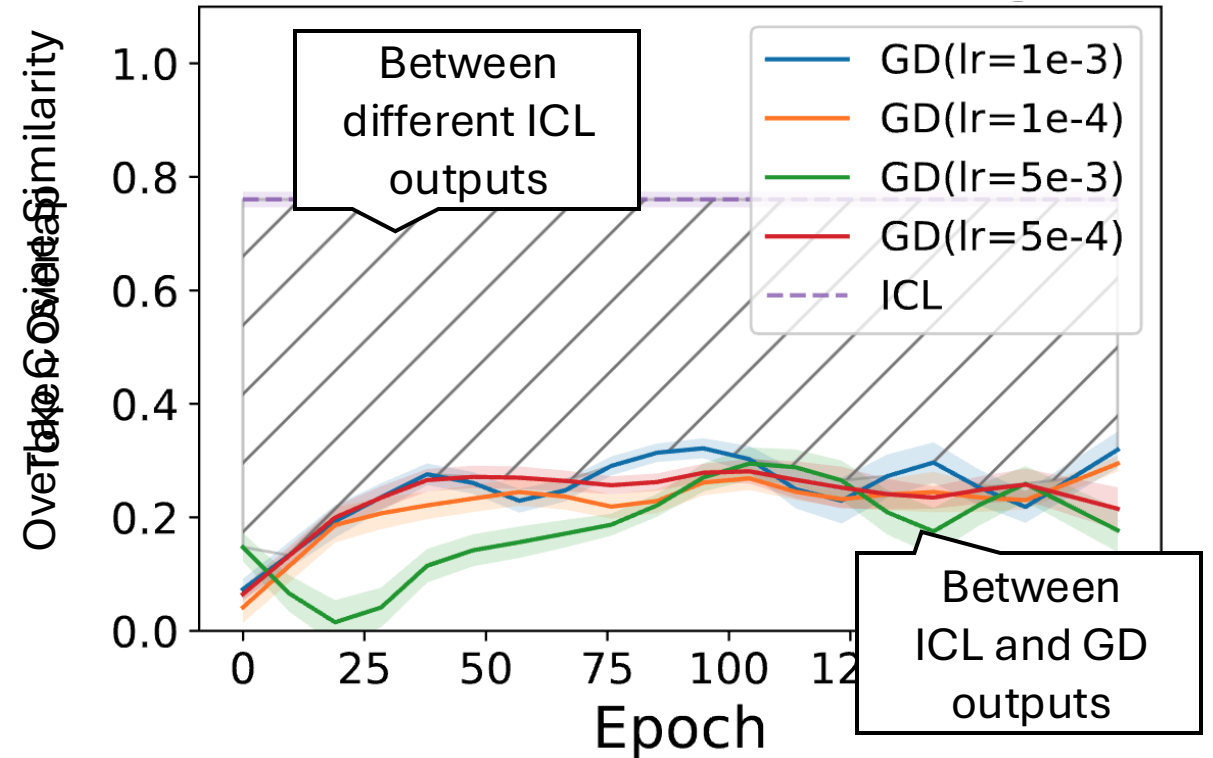
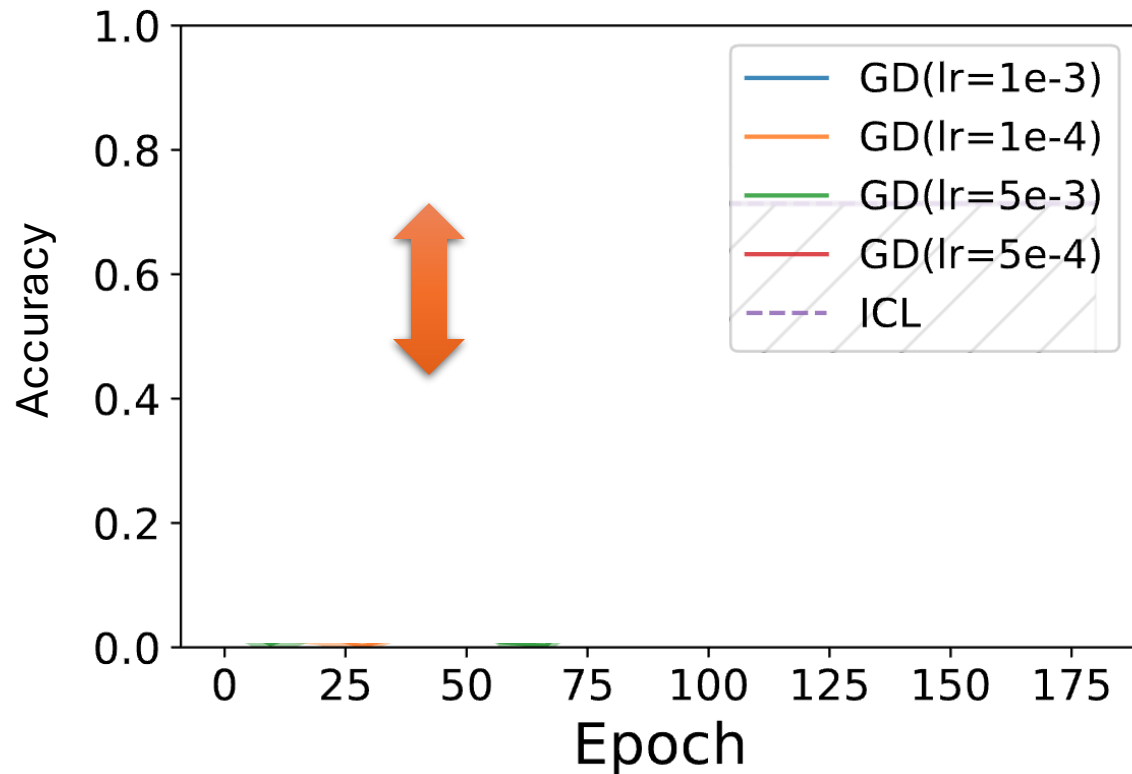
Evidence against $ICL \approx GD$: Performance

We use three coarse-to-fine metrics to compare ICL and GD:

1. **Accuracy:** compare top predicted tokens against the ground truth (instead of top predicted label).
2. **Token Overlap:** compare overlap in top-K tokens of two distributions.
3. **Overlap Cosine Similarity:** compare the individual agreement between top-K tokens.



Evidence against $ICL \approx GD$: Outputs



ICL performs **differently** and does **not align** with GD.

Summary

1. Recent works studying ICL **do not align with emergent ICL in LLMs.**
In-Context Learning → Learning to Learn, Meta Learning, etc.
Transformers learn in-context by gradient descent → Transformers can perform gradient descent in their forward pass when trained appropriately.
2. Expressivity of the Transformer architecture to simulate GD **does not imply** that LLMs actually do it.
We present arguments and evidence against the [current] $ICL \approx GD$ equivalence theory.
3. Maintain **parallels to real world** settings when developing.
It is OK to study ICL in a simpler setting like Linear Regression,
but need to find a corresponding pretraining distribution to elicit ICL.

Thank you!



Paper



Contact