

TRANSFORMERS LEARN HIGHER-ORDER OPTIMIZATION METHODS FOR IN-CONTEXT LEARNING: A STUDY WITH LINEAR MODELS

Anonymous authors

Paper under double-blind review

ABSTRACT

Transformers are remarkably good at *in-context learning* (ICL)—learning from demonstrations without parameter updates—but how they perform ICL remains a mystery. Recent work suggests that Transformers may learn in-context by internally running Gradient Descent (GD), a first-order optimization method. In this paper, we instead demonstrate that Transformers learn to implement higher-order optimization methods to perform ICL. Focusing on in-context linear regression, we show that Transformers learn to implement an algorithm very similar to *Iterative Newton’s Method*, a higher-order optimization method, rather than Gradient Descent. Empirically, we show that predictions from successive Transformer layers closely match different iterations of Newton’s Method *linearly*, with each middle layer roughly computing 3 iterations. In contrast, *exponentially* more GD steps are needed to match an additional Transformers layer; this suggests that Transformers have an comparable rate of convergence with high-order methods, which are exponentially faster than GD. We also show that Transformers can learn in-context on ill-conditioned data, a setting where Gradient Descent struggles but Iterative Newton succeeds. Finally, we show theoretical results which support our empirical findings and have a close correspondence with them: we prove that Transformers can implement k iterations of Newton’s method with $\mathcal{O}(k)$ layers.

1 INTRODUCTION

Transformer neural networks (Vaswani et al., 2017) have become the default architecture for natural language processing (Devlin et al., 2019; Brown et al., 2020; OpenAI, 2023), and have even been adopted by other areas like computer vision (Dosovitskiy et al., 2021). As first demonstrated by GPT-3 (Brown et al., 2020), Transformers excel at *in-context learning* (ICL)—learning from input-output pairs provided as inputs to the model, without updating their model parameters. Through in-context learning, Transformer-based Large Language Models (LLMs) can achieve state-of-the-art few-shot performance across a wide variety of downstream tasks (Rae et al., 2022; Smith et al., 2022; Thoppilan et al., 2022; Chowdhery et al., 2022).

Given the importance of Transformers and ICL, many prior efforts have attempted to understand how Transformers perform in-context learning. Prior work suggests Transformers can approximate linear functions well in-context (Garg et al., 2022). Specifically to linear regression tasks, prior work has tried to understand the ICL mechanism and the dominant hypothesis is that Transformers learn in-context by running optimizations internally through gradient-based algorithms (von Oswald et al., 2022; 2023; Ahn et al., 2023; Dai et al., 2023).

This paper presents strong evidence for a competing hypothesis: Transformers trained to perform in-context linear regression learn to implement a higher-order optimization method rather than a first-order method like Gradient Descent. In particular, Transformers implement a method very similar to Newton-Schulz’s Method, also known as the *Iterative Newton’s Method*, which iteratively improves an estimate of the inverse of the design matrix to compute the optimal weight vector. Across many layers of the Transformer, subsequent layers approximately compute more and more iterations of Newton’s Method, with increasingly better predictions; both eventually converge to the optimal minimum-norm solution found by ordinary least squares (OLS). Interestingly, this mechanism is

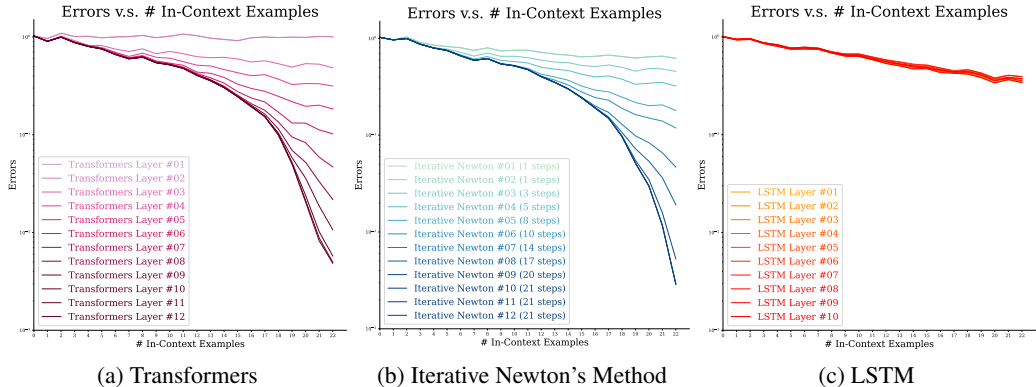


Figure 1: **Progression of Algorithms.** (a) Transformer’s performance improves over the layer index ℓ . (b) Iterative Newton’s performance improves over the number of iterations k , in a way that closely resembles the Transformer. We plot the best-matching k to Transformer’s ℓ following Definition 4. (c) In contrast, LSTM’s performance does not improve from layer to layer.

specific to Transformers: LSTMs do not learn these same higher-order methods, as their predictions do not even improve across layers.

We present both empirical and theoretical evidence for our claims. Empirically, Transformer induced weights and residuals are similar to Iterative Newton and deeper layers match Newton with more iterations (see Figures 1 and 9). Transformers can also handle ill-conditioned problems without requiring significantly more layers, where GD would suffer from slow convergence but Iterative Newton would not. Crucially, Transformers share the same rate of convergence as Iterative Newton and are exponentially faster than GD. Theoretically, we show that Transformer circuits can efficiently implement Iterative Newton, with the number of layers depending linearly on the number of iterations and the dimensionality of the hidden states depending linearly on the dimensionality of the data. Overall, our work provides a mechanistic account of how Transformers perform in-context learning that not only explains model behavior better than previous hypotheses, but also hints at what makes Transformers so well-suited for ICL compared with other neural architectures.

2 RELATED WORK

In-context learning by large language models. GPT-3 (Brown et al., 2020) first showed that Transformer-based large language models can “learn” to perform new tasks from in-context demonstrations (i.e., input-output pairs). Since then, a large body of work in NLP has studied in-context learning, for instance by understanding how the choice and order of demonstrations affects results (Lu et al., 2022; Liu et al., 2022; Rubin et al., 2022; Su et al., 2023; Chang & Jia, 2023; Nguyen & Wong, 2023), studying the effect of label noise (Min et al., 2022c; Yoo et al., 2022; Wei et al., 2023), and proposing methods to improve ICL accuracy (Zhao et al., 2021; Min et al., 2022a;b).

In-context learning beyond natural language. Inspired by the phenomenon of ICL by large language models, subsequent work has studied how Transformers learn in-context beyond NLP tasks. Garg et al. (2022) first investigated Transformers’ ICL abilities for various classical machine learning problems, including linear regression. We largely adopt their linear regression setup in this work. Li et al. (2023) formalize in-context learning as an algorithm learning problem where a Transformer model implicitly constructs a hypothesis function at inference-time and obtain generalization bounds for ICL. Han et al. (2023) suggests that Transformers learn in-context by performing Bayesian inference on prompts, which can be asymptotically interpreted as kernel regression. Tarzanagh et al. (2023a) and Tarzanagh et al. (2023b) show that Transformers can find max-margin solutions for classification tasks and act as support vector machines. Zhang et al. (2023) prove that a linear attention Transformer trained by gradient flow can indeed in-context learn class of linear models. Raventós et al. (2023) explore how diverse pretraining data can enable models to perform ICL on new tasks.

Do Transformers implement Gradient Descent? A growing body of work has suggested that Transformers learn in-context by implementing gradient descent within their internal representations. Akyürek et al. (2022) summarize operations that Transformers can implement, such as multi-

plication and affine transformations, and show that Transformers can implement gradient descent for linear regression using these operations. Concurrently, von Oswald et al. (2022) argue that Transformers learn in-context via gradient descent, where one layer performs one gradient update. In subsequent work, von Oswald et al. (2023) further argue that Transformers are strongly biased towards learning to implement gradient-based optimization routines. Ahn et al. (2023) extend the work of von Oswald et al. (2022) by showing Transformers can learn to implement preconditioned Gradient Descent, where the pre-conditioner can adapt to the data. Bai et al. (2023) provides detailed constructions for how Transformers can implement a range of learning algorithms via gradient descent. Finally, Dai et al. (2023) conduct experiments on NLP tasks and conclude that Transformer-based language models performing ICL behave similarly to models fine-tuned via gradient descent. In this paper, we argue that Transformers actually learn to perform in-context learning by implementing a higher-order optimization method, not gradient descent. Predictions made by different Transformer layers match iterations of higher-order optimization methods better than they match iterations of gradient descent; moreover, Transformers can handle ill-conditioned data, unlike Gradient Descent.

Mechanistic interpretability for Transformers. Our work attempts to understand the mechanism through which Transformers perform in-context learning. Prior work has studied other aspects of Transformers’ internal mechanisms, including reverse-engineering language models (Wang et al., 2022), the grokking phenomenon (Power et al., 2022; Nanda et al., 2023), manipulating attention maps (Hassid et al., 2022), and automated circuit finding (Conmy et al., 2023).

3 PROBLEM SETUP

In this paper, we focus on the following linear regression task. The task involves n examples $\{\mathbf{x}_i, y_i\}_{i=1}^n$ where $\mathbf{x}_i \in \mathbb{R}^d$ and $y_i \in \mathbb{R}$. The examples are generated from the following data generating distribution $P_{\mathcal{D}}$, parameterized by a distribution \mathcal{D} over $(d \times d)$ positive semi-definite matrices. For each sequence of n in-context examples, we first sample a ground-truth weight vector $\mathbf{w}^* \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I}) \in \mathbb{R}^d$ and a matrix $\Sigma \stackrel{\text{i.i.d.}}{\sim} \mathcal{D}$. For $i \in [n]$, we sample each $\mathbf{x}_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, \Sigma)$. The label y_i for each \mathbf{x}_i is given by $y_i = \mathbf{w}^{*\top} \mathbf{x}_i$. Note that for much of our experiments \mathcal{D} is only supported on the identity matrix $\mathbf{I} \in \mathbb{R}^{d \times d}$ and hence $\Sigma = \mathbf{I}$, but we also consider some distributions over ill-conditioned matrices which will give rise to ill-conditioned regression problems.

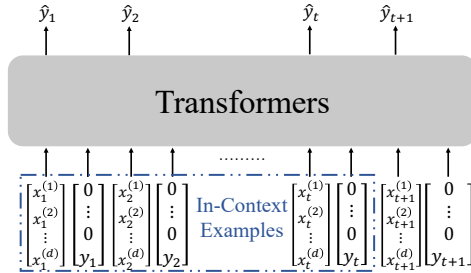


Figure 2: Illustration of how Transformers are trained to do in-context linear regression.

3.1 SOLVING LINEAR REGRESSION WITH TRANSFORMERS

We will use neural network models such as Transformers to solve this linear regression task. As shown in Figure 2, at time step $t + 1$, the model sees the first t in-context examples $\{\mathbf{x}_i, y_i\}_{i=1}^t$, and then makes predictions for \mathbf{x}_{t+1} , whose label y_{t+1} is not observed by the Transformers model.

We randomly initialize our models and then train them on the linear regression task to make predictions for every number of in-context examples t , where $t \in [n]$. Training and test data are both drawn from $P_{\mathcal{D}}$. To make the input prompts contain both \mathbf{x}_i and y_i , we follow same the setup as Garg et al. (2022)’s to zero-pad y_i ’s, and use the same decoder-only GPT-2 model (Radford et al., 2019) with softmax activation and causal attention mask (discussed later at definition 1).

We now present the key mathematical details for the Transformer architecture, and how they can be used for in-context learning. First, the causal attention mask enforces that attention heads can only attend to hidden states of previous time steps, and is defined as follows.

Definition 1 (Causal Attention Layer). An *causal attention layer* with M heads and activation function σ is denoted as Attn on any input sequence $\mathbf{H} = [\mathbf{h}_1, \dots, \mathbf{h}_N] \in \mathbb{R}^{D \times N}$, where D is the dimension of hidden states and N is the sequence length. In the vector form,

$$\tilde{\mathbf{h}}_t = [\text{Attn}(\mathbf{H})]_t = \mathbf{h}_t + \sum_{m=1}^M \sum_{j=1}^t \sigma(\langle \mathbf{Q}_m \mathbf{h}_t, \mathbf{K}_m \mathbf{h}_j \rangle) \cdot \mathbf{V}_m \mathbf{h}_j. \quad (1)$$

Vaswani et al. (2017) originally proposed the Transformer architecture with the Softmax activation function for the attention layers. Later works have found that replacing $\text{Softmax}(\cdot)$ with $\frac{1}{t}\text{ReLU}(\cdot)$ does not hurt model performance (Cai et al., 2022; Shen et al., 2023; Wortsman et al., 2023). The Transformers architecture is defined by putting together attention layers with feed forward layers:

Definition 2 (Transformers). *An L -layer decoder-based transformer with Causal Attention Layers is denoted as TF_θ and is a composition of a MLP Layer (with a skip connection) and a Causal Attention Layers. For input sequence $\mathbf{H}^{(0)}$, the transformers ℓ -th hidden layer is given by*

$$\text{TF}_\theta^\ell(\mathbf{H}^{(0)}) := \mathbf{H}^{(\ell)} = \text{MLP}_{\theta_{\text{mlp}}^{(\ell)}} \left(\text{Attn}_{\theta_{\text{attn}}^{(\ell)}}(\mathbf{H}^{(\ell-1)}) \right) \quad (2)$$

where $\theta = \{\theta_{\text{mlp}}^{(\ell)}, \theta_{\text{attn}}^{(\ell)}\}_{\ell=1}^L$ and $\theta_{\text{attn}}^{(\ell)} = \{\mathbf{Q}_m^{(\ell)}, \mathbf{K}_m^{(\ell)}, \mathbf{V}_m^{(\ell)}\}_{m=1}^M$ consists of M heads at layer ℓ .

In particular for the linear regression task, Transformers perform in-context learning as follows

Definition 3 (Transformers for Linear Regression). *Given in-context examples $\{\mathbf{x}_1, y_1, \dots, \mathbf{x}_t, y_t\}$, Transformers make predictions on a query example \mathbf{x}_{t+1} through a readout layer parameterized as $\theta_{\text{readout}} = \{\mathbf{u}, v\}$, and the prediction $\hat{y}_{t+1}^{\text{TF}}$ is given by*

$$\hat{y}_{t+1}^{\text{TF}} := \text{ReadOut} \left[\underbrace{\text{TF}_\theta^L(\{\mathbf{x}_1, y_1, \dots, \mathbf{x}_t, y_t, \mathbf{x}_{t+1}\})}_{\mathbf{H}^{(L)}} \right] = \mathbf{u}^\top \mathbf{H}_{:,2t+1}^{(L)} + v. \quad (3)$$

3.2 STANDARD METHODS FOR SOLVING LINEAR REGRESSION

Our central research question is:

Does the algorithm Transformers learn for linear regression resemble any known algorithm?

Here we discuss various known algorithms we compare Transformers with.

For any time step t , let $\mathbf{X}^{(t)} = [\mathbf{x}_1 \ \dots \ \mathbf{x}_t]^\top$ be the data matrix and $\mathbf{y}^{(t)} = [y_1 \ \dots \ y_t]^\top$ be the labels for all the datapoints seen so far. Note that since t can be smaller than the data dimension d , $\mathbf{X}^{(t)}$ can be singular. We now consider various algorithms for making predictions for \mathbf{x}_{t+1} based on $\mathbf{X}^{(t)}$ and $\mathbf{y}^{(t)}$. When it is clear from context, we will drop the superscript and refer to $\mathbf{X}^{(t)}$ and $\mathbf{y}^{(t)}$ as \mathbf{X} and \mathbf{y} , where \mathbf{X} and \mathbf{y} correspond to all the datapoints seen so far.

Ordinary Least Squares. This method finds the minimum-norm solution to the objective:

$$\mathcal{L}(\mathbf{w} \mid \mathbf{X}, \mathbf{y}) = \frac{1}{2n} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2. \quad (4)$$

The Ordinary Least Squares (OLS) solution has a closed form given by the Normal Equations:

$$\hat{\mathbf{w}}^{\text{OLS}} = (\mathbf{X}^\top \mathbf{X})^\dagger \mathbf{X}^\top \mathbf{y} \quad (5)$$

where $\mathbf{S} := \mathbf{X}^\top \mathbf{X}$ and \mathbf{S}^\dagger is the pseudo-inverse (Moore, 1920) of \mathbf{S} .

Gradient Descent. Gradient descent (GD) finds the weight vector $\hat{\mathbf{w}}^{\text{GD}}$ by initializing $\hat{\mathbf{w}}_0^{\text{GD}} = \mathbf{0}$ and using the iterative update rule:

$$\hat{\mathbf{w}}_{k+1}^{\text{GD}} = \hat{\mathbf{w}}_k^{\text{GD}} - \eta \nabla_{\mathbf{w}} \mathcal{L}(\hat{\mathbf{w}}_k^{\text{GD}} \mid \mathbf{X}, \mathbf{y}). \quad (6)$$

It is known that Gradient Descent requires $\mathcal{O}(\kappa(\mathbf{S}) \log(1/\epsilon))$ steps to converge to an ϵ error where $\kappa(\mathbf{S}) = \frac{\lambda_{\max}(\mathbf{S})}{\lambda_{\min}(\mathbf{S})}$ is the *condition number*. Thus, when $\kappa(\mathbf{S})$ is large, Gradient Descent converges slowly (Boyd & Vandenberghe, 2004).

Online Gradient Descent. While GD computes the gradient with respect to the full data matrix \mathbf{X} at each iteration, Online Gradient Descent (OGD) is an online algorithm that only computes gradients on the newly received data point $\{\mathbf{x}_k, y_k\}$ at step k :

$$\hat{\mathbf{w}}_{k+1}^{\text{OGD}} = \hat{\mathbf{w}}_k^{\text{OGD}} - \eta_k \nabla_{\mathbf{w}} \mathcal{L}(\hat{\mathbf{w}}_k^{\text{OGD}} \mid \mathbf{x}_k, y_k). \quad (7)$$

Picking $\eta_k = \frac{1}{\|\mathbf{x}_k\|_2^2}$ ensures that the new weight vector $\hat{\mathbf{w}}_{k+1}^{\text{OGD}}$ makes zero error on $\{\mathbf{x}_k, y_k\}$.

Iterative Newton’s Method. This method finds the weight vector $\hat{\mathbf{w}}^{\text{Newton}}$ by iteratively apply Newton’s method to finding the pseudo inverse of $\mathbf{S} = \mathbf{X}^\top \mathbf{X}$ (Schulz, 1933; Ben-Israel, 1965).

$$\begin{aligned} \mathbf{M}_0 &= \alpha \mathbf{S}, \text{ where } \alpha = \frac{2}{\|\mathbf{S} \mathbf{S}^\top\|_2}, & \hat{\mathbf{w}}_0^{\text{Newton}} &= \mathbf{M}_0 \mathbf{X}^\top \mathbf{y}, \\ \mathbf{M}_{k+1} &= 2\mathbf{M}_k - \mathbf{M}_k \mathbf{S} \mathbf{M}_k, & \hat{\mathbf{w}}_{k+1}^{\text{Newton}} &= \mathbf{M}_{k+1} \mathbf{X}^\top \mathbf{y}. \end{aligned} \quad (8)$$

This computes an approximation of the psuedo inverse using the moments of \mathbf{S} . In contrast to GD, the Iterative Newton’s method only requires $\mathcal{O}(\log \kappa(\mathbf{S}) + \log \log(1/\epsilon))$ steps to converge to an ϵ error (Soderstrom & Stewart, 1974; Pan & Schreiber, 1991). Note that this is exponentially faster than the convergence rate of Gradient Descent.

3.3 LSTM

While our primary goal is to analyze Transformers, we also consider the LSTM architecture (Hochreiter & Schmidhuber, 1997) to understand whether Transformers learn different algorithms than other neural sequence models trained to do linear regression. In particular, we train a unidirectional L -layer LSTM, which generates a sequence of hidden states $\mathbf{H}^{(\ell)}$ for each layer ℓ , similarly to an L -layer Transformer. As with Transformers, we add a readout layer that predicts the $\hat{y}_{t+1}^{\text{LSTM}}$ from the final hidden state at the final layer, $\mathbf{H}_{:,2t+1}^{(L)}$.

3.4 MEASURING ALGORITHMIC SIMILARITY

We propose two metrics to measure the similarity between linear regression algorithms.

Similarity of Errors. For a linear regression algorithm \mathcal{A} , let $\mathcal{A}(\mathbf{x}_{t+1} \mid \{\mathbf{x}_i, y_i\}_{i=1}^t)$ denote its prediction on the $(t+1)$ -th in-context example \mathbf{x}_{t+1} after observing the first t examples (see Figure 2). We write $\mathcal{A}(\mathbf{x}_{t+1}) := \mathcal{A}(\mathbf{x}_{t+1} \mid \{\mathbf{x}_i, y_i\}_{i=1}^t)$ for brevity. The errors (i.e., residuals) on the data sequence are:¹

$$\mathcal{E}(\mathcal{A} \mid \{\mathbf{x}_i, y_i\}_{i=1}^{n+1}) = \left[\mathcal{A}(\mathbf{x}_2) - y_2, \dots, \mathcal{A}(\mathbf{x}_{n+1}) - y_{n+1} \right]^\top \in \mathbb{R}^n. \quad (9)$$

For any two algorithms \mathcal{A}_a and \mathcal{A}_b , their similarity of errors, corresponding to the metric $\mathcal{C}(\cdot, \cdot)$, is

$$\text{SimE}(\mathcal{A}_a, \mathcal{A}_b) = \mathbb{E}_{\{\mathbf{x}_i, y_i\}_{i=1}^{n+1} \sim P_{\mathcal{D}}} \mathcal{C}\left(\mathcal{E}(\mathcal{A}_a \mid \{\mathbf{x}_i, y_i\}_{i=1}^{n+1}), \mathcal{E}(\mathcal{A}_b \mid \{\mathbf{x}_i, y_i\}_{i=1}^{n+1})\right) \quad (10)$$

where we use the cosine similarity as our correlation metric $\mathcal{C}(\mathbf{u}, \mathbf{v}) = \frac{\langle \mathbf{u}, \mathbf{v} \rangle}{\|\mathbf{u}\|_2 \|\mathbf{v}\|_2}$. Here n is the total number of in-context examples and $P_{\mathcal{D}}$ is the data generation process discussed previously.

Similarity of Induced Weights. All standard algorithms for linear regression estimate a weight vector $\hat{\mathbf{w}}$. While neural ICL models like Transformers do not explicitly learn such a weight vector, similar to Akyürek et al. (2022), we can *induce* an implicit weight vector $\tilde{\mathbf{w}}$ learned by any algorithm \mathcal{A} by fitting a weight vector to its predictions. To do this, for any fixed sequence of t in-context examples $\{\mathbf{x}_i, y_i\}_{i=1}^t$, we sample $T \gg d$ query examples $\tilde{\mathbf{x}}_k \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, \Sigma)$, where $k \in [T]$. For this fixed sequence of in-context examples $\{\mathbf{x}_i, y_i\}_{i=1}^t$, we create T in-context prediction tasks and use the algorithm \mathcal{A} to make predictions $\mathcal{A}(\tilde{\mathbf{x}}_k \mid \{\mathbf{x}_i, y_i\}_{i=1}^t)$. We define the induced data matrix and labels as

$$\tilde{\mathbf{X}} = \begin{bmatrix} \tilde{\mathbf{x}}_1^\top \\ \vdots \\ \tilde{\mathbf{x}}_T^\top \end{bmatrix} \quad \tilde{\mathbf{Y}} = \begin{bmatrix} \mathcal{A}(\tilde{\mathbf{x}}_1 \mid \{\mathbf{x}_i, y_i\}_{i=1}^t) \\ \vdots \\ \mathcal{A}(\tilde{\mathbf{x}}_T \mid \{\mathbf{x}_i, y_i\}_{i=1}^t) \end{bmatrix}. \quad (11)$$

Then, the induced weight vector for \mathcal{A} and these t in-context examples is:

$$\tilde{\mathbf{w}}_t(\mathcal{A}) := \tilde{\mathbf{w}}_t(\mathcal{A} \mid \{\mathbf{x}_i, y_i\}_{i=1}^t) = (\tilde{\mathbf{X}}^\top \tilde{\mathbf{X}})^{-1} \tilde{\mathbf{X}}^\top \tilde{\mathbf{Y}}. \quad (12)$$

¹the indices start from 2 to $n+1$ because we evaluate all cases where t can choose from $1, \dots, n$.

We measure the similarity between two algorithms \mathcal{A}_a and \mathcal{A}_b by measuring the similarity of induced weight vectors $\tilde{\mathbf{w}}_t(\mathcal{A}_a)$ and $\tilde{\mathbf{w}}_t(\mathcal{A}_b)$. We define the similarity of induced weights between two algorithms as

$$\text{SimW}(\mathcal{A}_a, \mathcal{A}_b) = \mathbb{E}_{\{\mathbf{x}_i, y_i\}_{i=1}^n \sim P_{\mathcal{D}}} \frac{1}{n} \sum_{t=1}^n \mathcal{C}(\tilde{\mathbf{w}}_t(\mathcal{A}_a | \{\mathbf{x}_i, y_i\}_{i=1}^t), \tilde{\mathbf{w}}_t(\mathcal{A}_b | \{\mathbf{x}_i, y_i\}_{i=1}^t)).$$

Each algorithm we consider has its own hyper-parameter(s), for example the number of iterations for Iterative Newton and Gradient Descent, and the number of layers for Transformers (see Section 4.1). When comparing two algorithms, given a choice of hyper-parameters for the first algorithm, we compare with the hyper-parameters for the second algorithm that maximize algorithmic similarity. In other words, we measure whether there exists hyperparameters for the second algorithm that make the two algorithms are similar.

Definition 4 (Hyper-Parameter Matching). *Let \mathcal{M} be the metric for evaluating similarities between two algorithms \mathcal{A}_a and \mathcal{A}_b , which have hyper-parameters $p_a \in \mathcal{P}_a$ and $p_b \in \mathcal{P}_b$, respectively. For a given choice of p_a , We define the best-matching hyper-parameters of algorithm \mathcal{A}_b for \mathcal{A}_a as:*

$$p_b^{\mathcal{M}}(p_a) := \arg \min_{p_b \in \mathcal{P}_b} \mathcal{M}(\mathcal{A}_a(\cdot | p_a), \mathcal{A}_b(\cdot | p_b)). \quad (13)$$

The matching processes can be visualized as heatmaps as shown in Figure 3, where best-matching hyperparameters are highlighted. This enables us to better compare the rate of convergence of algorithms and we will discuss these results further in §4.

4 EXPERIMENTAL EVIDENCE

We mainly work on the Transformers-based GPT-2 model with 12 layers with 8 heads per layer. Alternative configurations with fewer heads per layer also support our findings and we defer them to Appendix A. We initially focus on isotropic cases where $\Sigma = \mathbf{I}$ and later consider ill-conditioned Σ in §4.3. Our training setup is exactly the same as Garg et al. (2022): models are trained with at most $n = 40$ in-context examples for $d = 20$ (with the same learning rate, batch size etc.). We claim that Transformers learn high-order optimization methods in-context.

4.1 TRANSFORMERS IMPROVE PROGRESSIVELY OVER LAYERS

Many known algorithms for linear regression, including GD, OGD, and Iterative Newton, are *iterative*: their performance progressively improves as they perform more iterations, eventually converging to a final solution. How could a Transformer implement such an iterative algorithm? von Oswald et al. (2022) propose that deeper *layers* of the Transformer may correspond to more iterations of an iterative method; in particular, they show that there exist Transformer parameters such that each attention layer performs one step of Gradient Descent.

Following this intuition, we first investigate whether the predictions of a trained Transformer improve as the layer index ℓ increases. For each layer of hidden states $\mathbf{H}^{(\ell)}$ (defined in definition 2), we re-train the ReadOut layer to predict y_t for each t ; the new predictions are given by $\text{ReadOut}^{(\ell)}[\mathbf{H}^{(\ell)}]$. Thus for each input prompt, there are L Transformer predictions parameterized by layer index ℓ . All parameters besides the Readout layer parameters are kept frozen.

As shown in Figure 1a, as we increase the layer index ℓ , the prediction performance improves progressively. Hence, Transformers progressively improve their predictions over layers ℓ , similar to how iterative algorithms improve over steps.

4.2 TRANSFORMERS ARE MORE SIMILAR TO ITERATIVE NEWTON’S METHOD

Next, we test the more specific hypothesis that the iterative updates performed across Transformer layers are similar to the iterative updates for known iterative algorithms. Specifically, we test whether each layer ℓ of the Transformer corresponds to performing k steps of some iterative algorithm, for some k depending on ℓ . We focus here on Gradient Descent and Iterative Newton’s Method; we will discuss online algorithms in Section 4.4.

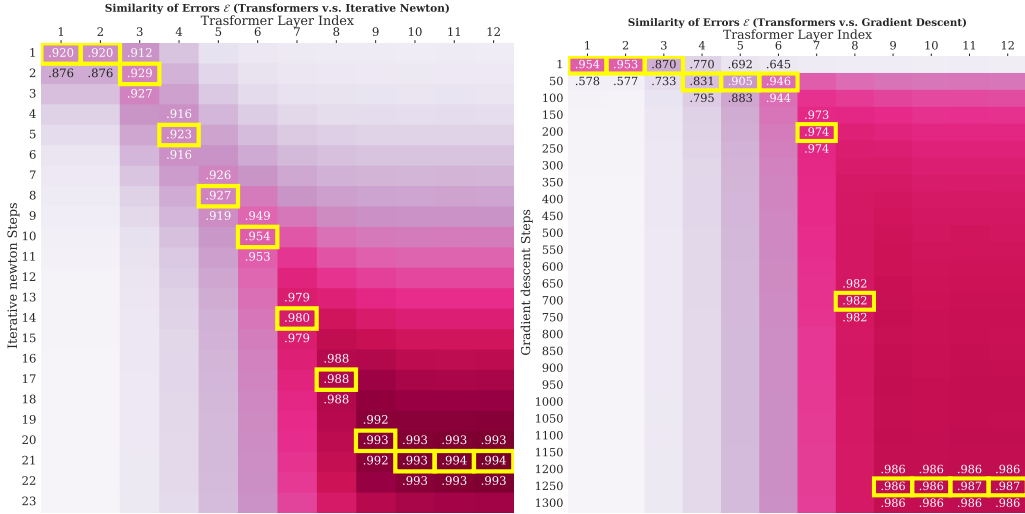


Figure 3: **Heatmaps of Similarity.** The best matching hyper-parameters are highlighted in yellow. See Figure 8 for an additional heatmap where Gradient Descent’s steps are shown in log scale.

For each layer ℓ of the Transformer, we measure the best-matching similarity (see Def. 4) with candidate iterative algorithms with the optimal choice of the number of steps k . As shown in Figure 3, the Transformer has very high error similarity with Iterative Newton’s method at all layers. Moreover, we see a clear *linear* trend between layer 3 and layer 9 of the Transformer, where each layer appears to compute roughly 3 additional iterations of the Iterative Newton’s method. This trend only stops at the last few layers because both algorithms converge to the OLS solution; Newton is known to converge to OLS (see §3.2), and we verify in Appendix A.1 that the last few layers of the Transformer also basically compute OLS (see Figure 10 in the Appendix). We observe the same trends when using similarity of induced weights as our similarity metric (see Figure 7 in the Appendix),

In contrast, even though GD has a comparable similarity with the Transformers at later layers, their best matching follows an *exponential* trend. For well-conditioned problems, to achieve ϵ error, the convergence rate of GD is $\mathcal{O}(\log(1/\epsilon))$ while the convergence rate of Newton is $\mathcal{O}(\log \log(1/\epsilon))$ (see §3.2). Therefore the convergence rate of Newton is exponentially faster than GD. Transformer’s *linear* correspondence with Newton and its *exponential* correspondence with GD provides strong evidence that the rate of convergence of Transformers is similar to Newton, i.e., $\mathcal{O}(\log \log(1/\epsilon))$.

Overall, we conclude that a Transformer trained to perform in-context linear regression learns to implement an algorithm that is very similar to Iterative Newton’s method, not Gradient Descent. Starting at layer 3, subsequent layers of the Transformer compute more and more iterations of Iterative Newton’s method. This algorithm successfully solves the linear regression problem, as it converges to the optimal OLS solution in the final layers.

4.3 TRANSFORMERS PERFORM WELL ON ILL-CONDITIONED DATA

We repeat the same experiments with data $x_i \stackrel{i.i.d.}{\sim} \mathcal{N}(\mathbf{0}, \Sigma)$ sampled from an ill-condition covariance matrix Σ with condition number $\kappa(\Sigma) = 100$, and eigenbasis chosen uniformly at random. The first $d/2$ eigenvalues of Σ are 100, the last $d/2$ are 1.

As shown in Figure 4, the Transformer model’s performance still closely matches Iterative Newton’s Method with 21 iterations, same as when $\Sigma = I$ (see layer 10-12 in Figure 3). The convergence of higher-order methods has a mild logarithmic dependence on the condition number since they correct for the curvature. On the other hand, Gradient Descent’s convergence is affected polynomially by conditioning. As $\kappa(\Sigma)$ increase from 1 to 100, the number steps required for GD’s convergence increases

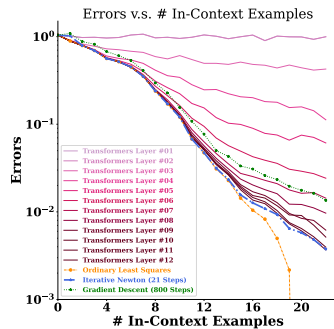


Figure 4: Transformers performance on ill-conditioned data.

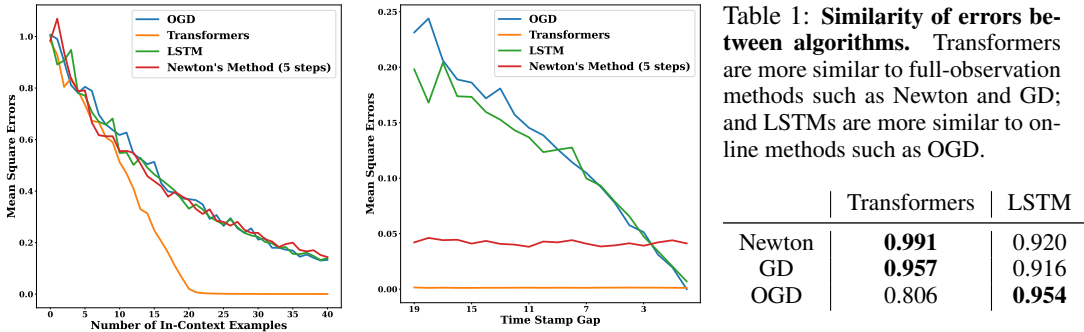


Figure 5: In the left figure, we measure model predictions with normalized MSE. Though LSTM is seemingly most similar to Newton’s Method with only 5 steps, neither algorithm converges yet. OGD also has a similar trend as LSTM. In the center figure, we measure model’s forgetting phenomenon (see Appendix A.4 for explanations), and find both Transformers and not-converged Newton have better memorization than LSTM and OGD. In the right table, we find that Transformers are more similar to Newton and GD than LSTMs do while LSTM is significantly more similar to OGD.

significantly (see Fig. 4 where GD requires 800 steps to converge), making it impossible for a 12-layer Transformers to implement these many gradient updates. We also note that preconditioning the data by $(\mathbf{X}^\top \mathbf{X})^\dagger$ can make the data well-conditioned, but since the eigenbasis is chosen uniformly at random, with high probability there is no sparse pre-conditioner or any fixed pre-conditioner which works across the data distribution. Computing $(\mathbf{X}^\top \mathbf{X})^\dagger$ appears to be as hard as computing the OLS solution (Eq. 4)—in fact Sharan et al. (2019) conjecture that first-order methods such as gradient descent and its variants cannot avoid polynomial dependencies in condition number in the ill-conditioned case.

See Appendix A.2 for detailed experiments on ill-conditioned problems. These experiments on ill-conditioned data further strengthen our hypothesis that Transformers are learning to perform higher-order optimization methods in-context, not Gradient Descent.

4.4 LSTM IS MORE SIMILAR TO OGD THAN TRANSFORMERS

As discussed in Section 3, LSTM is an alternative auto-regressive model widely used before the introduction of Transformers. Thus, a natural research question is: *If Transformers can learn in-context, can LSTMs do so as well? If so, do they learn the same algorithms?* To answer this question, we train a 10-layer LSTM model, with 5.3M parameters, in an identical manner to the Transformers (with 9.5M parameters) studied in the previous sections.²

Figure 5 plots the mean squared error of Transformers, LSTMs, and other standard methods as a function of the number of in-context (i.e., training) examples provided. While LSTMs can also learn linear regression in-context, they have much higher mean-squared error than Transformers. Their error rate is similar to Iterative Newton’s Method after only 5 iterations, a point where it is far from converging to the OLS solution.

LSTMs’ inferior performance to Transformers can be explained by the inability of LSTMs to use deeper layers to improve their predictions. Figure 1 shows that LSTM performance does not improve across layers—a readout head fine-tuned for the first layer makes equally good predictions as the full 10-layer model. Thus, LSTMs seem poorly equipped to fully implement iterative algorithms.

Finally, we show that LSTMs behave more like an online learning algorithm than Transformers. In particular, its predictions are biased towards getting more recent training examples correct, as opposed to earlier examples, as shown in Figure 5. This property makes LSTMs similar to online Gradient Descent. In contrast, five steps of Newton’s method has the same error on average for recent and early examples, showing that the LSTM implements a very different algorithm from a few iterations of Newton. Similarly, Table 1 shows that LSTMs are more similar to OGD than

²While the LSTM has fewer parameters than the Transformer, we found in preliminary experiments that increasing the hidden dimension or number of layers in the LSTM would not substantively change our results.

Transformers are, whereas Transformers are more similar to Newton and GD than LSTMs. We hypothesize that since LSTMs have limited memory, they must learn in a roughly online fashion; in contrast, Transformers’ attention heads can access the entire sequence of past examples, enabling it to learn more complex algorithms.

5 MECHANISTIC EVIDENCE

Our empirical evidence demonstrates that Transformers behave much more similarly to Iterative Newton’s than to Gradient Descent. Iterative Newton is a higher-order optimization method, and is algorithmically more involved than Gradient Descent. We begin by first examining this difference in complexity. As discussed in Section 3, the updates for Iterative Newton are of the form,

$$\hat{\mathbf{w}}_{k+1}^{\text{Newton}} = \mathbf{M}_{k+1} \mathbf{X}^\top \mathbf{y} \quad \text{where} \quad \mathbf{M}_{k+1} = 2\mathbf{M}_k - \mathbf{M}_k \mathbf{S} \mathbf{M}_k \quad (14)$$

and $\mathbf{M}_0 = \alpha \mathbf{S}$ for some $\alpha > 0$. We can express \mathbf{M}_k in terms of powers of \mathbf{S} by expanding iteratively, for example $\mathbf{M}_1 = 2\alpha \mathbf{S} - 4\alpha^2 \mathbf{S}^3$, $\mathbf{M}_2 = 4\alpha \mathbf{S} - 12\alpha^2 \mathbf{S}^3 + 16\alpha^3 \mathbf{S}^5 - 16\alpha^4 \mathbf{S}^7$, and in general $\mathbf{M}_k = \sum_{s=1}^{2^{k+1}-1} \beta_s \mathbf{S}^s$ for some $\beta_s \in \mathbb{R}$ (see Appendix B.3 for detailed calculations). Note that k steps of Iterative Newton’s requires computing $\Omega(2^k)$ moments of \mathbf{S} . Let us contrast this with Gradient Descent. Gradient Descent updates for linear regression take the form,

$$\hat{\mathbf{w}}_{k+1}^{\text{GD}} = \hat{\mathbf{w}}_k^{\text{GD}} - \eta(\mathbf{S} \hat{\mathbf{w}}_k^{\text{GD}} - \mathbf{X}^\top \mathbf{y}). \quad (15)$$

Like Iterative Newton, we can express $\hat{\mathbf{w}}_k^{\text{GD}}$ in terms of powers of \mathbf{S} and $\mathbf{X}^\top \mathbf{y}$. However, after k steps of Gradient Descent, the highest power of \mathbf{S} is only $O(k)$. This exponential separation is consistent with the exponential gap in terms of the parameter dependence in the convergence rate— $\mathcal{O}(\kappa(\mathbf{S}) \log(1/\epsilon))$ steps for Gradient Descent compared to $\mathcal{O}(\log \kappa(\mathbf{S}) + \log \log(1/\epsilon))$ steps for Iterative Newton. Therefore, a natural question is whether Transformers can actually represent as complicated of a method as Iterative Newton with only polynomially many layers?

Theorem 1 shows that this is indeed possible.

Theorem 1. *There exist Transformer weights such that on any set of in-context examples $\{\mathbf{x}_i, y_i\}_{i=1}^n$ and test point \mathbf{x}_{test} , the Transformer predicts on \mathbf{x}_{test} using $\mathbf{x}_{\text{test}}^\top \hat{\mathbf{w}}_k^{\text{Newton}}$. Here $\hat{\mathbf{w}}_k^{\text{Newton}}$ are the Iterative Newton updates given by $\hat{\mathbf{w}}_k^{\text{Newton}} = \mathbf{M}_k \mathbf{X}^\top \mathbf{y}$ where \mathbf{M}_j is updated as*

$$\mathbf{M}_j = 2\mathbf{M}_{j-1} - \mathbf{M}_{j-1} \mathbf{S} \mathbf{M}_{j-1}, 1 \leq j \leq k, \quad \mathbf{M}_0 = \alpha \mathbf{S},$$

for some $\alpha > 0$ and $\mathbf{S} = \mathbf{X}^\top \mathbf{X}$. The number of layers of the transformer is $\mathcal{O}(k)$ and the dimensionality of the hidden layers is $\mathcal{O}(d)$.

We note that our proof uses full attention instead of causal attention and ReLU activations for the self-attention layers. The definitions of these and the full proof appear in Appendix B.

6 CONCLUSION AND DISCUSSION

In this work, we studied how Transformers perform in-context learning for linear regression. In contrast with the hypothesis that Transformers learn in-context by implementing gradient descent, our experimental results show that different Transformer layers match iterations of Iterative Newton’s method *linearly* and Gradient Descent *exponentially*. This suggests that Transformers share a similar rate of convergence to Iterative Newton’s method but not Gradient Descent. Moreover, Transformers can perform well empirically on ill-conditioned linear regression, whereas first-order methods such as Gradient Descent struggle. Theoretically, we provide exact Transformer circuits that can implement k -steps of Iterative Newton’s method with $\mathcal{O}(k)$ layers to make predictions from in-context examples.

An interesting direction is to explore a wider range of higher-order methods that Transformers can implement. It also seems promising to extend our analysis to classification problems, especially given recent work showing that Transformers resemble SVMs in classification tasks (Li et al., 2023; Tarzanagh et al., 2023b;a). Finally, a natural question is to understand the differences in the model architecture that make Transformers better in-context learners than LSTMs. Based on our investigations with LSTMs, we hypothesize that Transformers can implement more powerful algorithms because of having access to a longer history of examples. Investigating the role of this additional memory in learning appears to be an intriguing direction.

REFERENCES

- Kwangjun Ahn, Xiang Cheng, Hadi Daneshmand, and Suvrit Sra. Transformers learn to implement preconditioned gradient descent for in-context learning. *ArXiv*, abs/2306.00297, 2023. 1, 3
- Ekin Akyürek, Dale Schuurmans, Jacob Andreas, Tengyu Ma, and Denny Zhou. What learning algorithm is in-context learning? investigations with linear models. *ArXiv*, abs/2211.15661, 2022. 2, 5, 19, 20, 21
- Yu Bai, Fan Chen, Haiquan Wang, Caiming Xiong, and Song Mei. Transformers as statisticians: Provable in-context learning with in-context algorithm selection. *ArXiv*, abs/2306.04637, 2023. 3
- Adi Ben-Israel. An iterative method for computing the generalized inverse of an arbitrary matrix. *Mathematics of Computation*, 19(91):452–455, 1965. ISSN 00255718, 10886842. URL <http://www.jstor.org/stable/2003676>. 5
- Stephen P Boyd and Lieven Vandenbergh. *Convex optimization*. Cambridge university press, 2004. 4
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (eds.), *Advances in Neural Information Processing Systems*, volume 33, pp. 1877–1901. Curran Associates, Inc., 2020. URL https://proceedings.neurips.cc/paper_files/paper/2020/file/1457c0d6bfc4967418bfb8ac142f64a-Paper.pdf. 1, 2
- Han Cai, Chuang Gan, and Song Han. Efficientvit: Enhanced linear attention for high-resolution low-computation visual recognition. *ArXiv*, abs/2205.14756, 2022. 4
- Ting-Yun Chang and Robin Jia. Data curation alone can stabilize in-context learning. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 8123–8144, Toronto, Canada, July 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.acl-long.452. URL <https://aclanthology.org/2023.acl-long.452>. 2
- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, Parker Schuh, Kensen Shi, Sasha Tsvyashchenko, Joshua Maynez, Abhishek Rao, Parker Barnes, Yi Tay, Noam Shazeer, Vinodkumar Prabhakaran, Emily Reif, Nan Du, Ben Hutchinson, Reiner Pope, James Bradbury, Jacob Austin, Michael Isard, Guy Gur-Ari, Pengcheng Yin, Toju Duke, Anselm Levskaya, Sanjay Ghemawat, Sunipa Dev, Henryk Michalewski, Xavier Garcia, Vedant Misra, Kevin Robinson, Liam Fedus, Denny Zhou, Daphne Ippolito, David Luan, Hyeontaek Lim, Barret Zoph, Alexander Spiridonov, Ryan Sepassi, David Dohan, Shivani Agrawal, Mark Omernick, Andrew M. Dai, Thanumalayan Sankaranarayanan Pillai, Marie Pellat, Aitor Lewkowycz, Erica Moreira, Rewon Child, Oleksandr Polozov, Katherine Lee, Zongwei Zhou, Xuezhi Wang, Brennan Saeta, Mark Diaz, Orhan Firat, Michele Catasta, Jason Wei, Kathy Meier-Hellstern, Douglas Eck, Jeff Dean, Slav Petrov, and Noah Fiedel. Palm: Scaling language modeling with pathways, 2022. 1
- Arthur Conmy, Augustine N. Mavor-Parker, Aengus Lynch, Stefan Heimersheim, and Adrià Garriga-Alonso. Towards automated circuit discovery for mechanistic interpretability, 2023. 3
- Damai Dai, Yutao Sun, Li Dong, Yaru Hao, Zhifang Sui, and Furu Wei. Why can gpt learn in-context? language models secretly perform gradient descent as meta-optimizers. *ArXiv*, abs/2212.10559, 2023. 1, 3
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of*

- the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pp. 4171–4186, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1423. URL <https://aclanthology.org/N19-1423>. 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. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=YicbFdNTTy>. 1
- Shivam Garg, Dimitris Tsipras, Percy Liang, and Gregory Valiant. What can transformers learn in-context? a case study of simple function classes. *ArXiv*, abs/2208.01066, 2022. 1, 2, 3, 6
- Chi Han, Ziqi Wang, Han Zhao, and Heng Ji. In-context learning of large language models explained as kernel regression, 2023. 2
- Michael Hassid, Hao Peng, Daniel Rotem, Jungo Kasai, Ivan Montero, Noah A. Smith, and Roy Schwartz. How much does attention actually attend? questioning the importance of attention in pretrained transformers, 2022. 3
- Sepp Hochreiter and Jürgen Schmidhuber. Long Short-Term Memory. *Neural Computation*, 9(8):1735–1780, 11 1997. ISSN 0899-7667. doi: 10.1162/neco.1997.9.8.1735. URL <https://doi.org/10.1162/neco.1997.9.8.1735>. 5
- Yingcong Li, Muhammed Emrullah Ildiz, Dimitris Papailiopoulos, and Samet Oymak. Transformers as algorithms: Generalization and stability in in-context learning. In *International Conference on Machine Learning*, 2023. 2, 9
- Jiachang Liu, Dinghan Shen, Yizhe Zhang, Bill Dolan, Lawrence Carin, and Weizhu Chen. What makes good in-context examples for GPT-3? In *Proceedings of Deep Learning Inside Out (DeLIO 2022): The 3rd Workshop on Knowledge Extraction and Integration for Deep Learning Architectures*, pp. 100–114, Dublin, Ireland and Online, May 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.deelio-1.10. URL <https://aclanthology.org/2022.deelio-1.10>. 2
- Yao Lu, Max Bartolo, Alastair Moore, Sebastian Riedel, and Pontus Stenetorp. Fantastically ordered prompts and where to find them: Overcoming few-shot prompt order sensitivity. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 8086–8098, Dublin, Ireland, May 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.acl-long.556. URL <https://aclanthology.org/2022.acl-long.556>. 2
- Sewon Min, Mike Lewis, Hannaneh Hajishirzi, and Luke Zettlemoyer. Noisy channel language model prompting for few-shot text classification. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 5316–5330, Dublin, Ireland, May 2022a. Association for Computational Linguistics. doi: 10.18653/v1/2022.acl-long.365. URL <https://aclanthology.org/2022.acl-long.365>. 2
- Sewon Min, Mike Lewis, Luke Zettlemoyer, and Hannaneh Hajishirzi. MetaICL: Learning to learn in context. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp. 2791–2809, Seattle, United States, July 2022b. Association for Computational Linguistics. doi: 10.18653/v1/2022.naacl-main.201. URL <https://aclanthology.org/2022.naacl-main.201>. 2
- Sewon Min, Xinxi Lyu, Ari Holtzman, Mikel Artetxe, Mike Lewis, Hannaneh Hajishirzi, and Luke Zettlemoyer. Rethinking the role of demonstrations: What makes in-context learning work? In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pp. 11048–11064, Abu Dhabi, United Arab Emirates, December 2022c. Association for Computational Linguistics. doi: 10.18653/v1/2022.emnlp-main.759. URL <https://aclanthology.org/2022.emnlp-main.759>. 2

- E.H Moore. On the reciprocal of the general algebraic matrix. *Bulletin of American Mathematical Society*, 26:394–395, 1920. 4
- Neel Nanda, Lawrence Chan, Tom Lieberum, Jess Smith, and Jacob Steinhardt. Progress measures for grokking via mechanistic interpretability, 2023. 3
- Tai Nguyen and Eric Wong. In-context example selection with influences. *arXiv preprint arXiv:2302.11042*, 2023. 2
- OpenAI. Gpt-4 technical report, 2023. URL <http://arxiv.org/abs/2303.08774v3>. 1
- Victor Y. Pan and Robert S. Schreiber. An improved newton iteration for the generalized inverse of a matrix, with applications. *SIAM J. Sci. Comput.*, 12:1109–1130, 1991. 5
- Alethea Power, Yuri Burda, Harri Edwards, Igor Babuschkin, and Vedant Misra. Grokking: Generalization beyond overfitting on small algorithmic datasets, 2022. 3
- Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language models are unsupervised multitask learners. 2019. 3
- Jack W. Rae, Sebastian Borgeaud, Trevor Cai, Katie Millican, Jordan Hoffmann, Francis Song, John Aslanides, Sarah Henderson, Roman Ring, Susannah Young, Eliza Rutherford, Tom Hennigan, Jacob Menick, Albin Cassirer, Richard Powell, George van den Driessche, Lisa Anne Hendricks, Maribeth Rauh, Po-Sen Huang, Amelia Glaese, Johannes Welbl, Sumanth Dathathri, Saffron Huang, Jonathan Uesato, John Mellor, Irina Higgins, Antonia Creswell, Nat McAleese, Amy Wu, Erich Elsen, Siddhant Jayakumar, Elena Buchatskaya, David Budden, Esme Sutherland, Karen Simonyan, Michela Paganini, Laurent Sifre, Lena Martens, Xiang Lorraine Li, Adhiguna Kuncoro, Aida Nematzadeh, Elena Gribovskaya, Domenic Donato, Angeliki Lazaridou, Arthur Mensch, Jean-Baptiste Lespiau, Maria Tsimpoukelli, Nikolai Grigorev, Doug Fritz, Thibault Sottiaux, Mantas Pajarskas, Toby Pohlen, Zhitao Gong, Daniel Toyama, Cyprien de Masson d’Autume, Yujia Li, Tayfun Terzi, Vladimir Mikulik, Igor Babuschkin, Aidan Clark, Diego de Las Casas, Aurelia Guy, Chris Jones, James Bradbury, Matthew Johnson, Blake Hechtman, Laura Weidinger, Iason Gabriel, William Isaac, Ed Lockhart, Simon Osindero, Laura Rimell, Chris Dyer, Oriol Vinyals, Kareem Ayoub, Jeff Stanway, Lorraine Bennett, Demis Hassabis, Koray Kavukcuoglu, and Geoffrey Irving. Scaling language models: Methods, analysis & insights from training gopher, 2022. 1
- Allan Raventós, Mansheej Paul, Feng Chen, and Surya Ganguli. Pretraining task diversity and the emergence of non-bayesian in-context learning for regression, 2023. 2
- Ohad Rubin, Jonathan Herzig, and Jonathan Berant. Learning to retrieve prompts for in-context learning. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp. 2655–2671, Seattle, United States, July 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.naacl-main.191. URL <https://aclanthology.org/2022.naacl-main.191>. 2
- Günther Schulz. Iterative berechnung der reziproken matrix. *Zeitschrift für Angewandte Mathematik und Mechanik (Journal of Applied Mathematics and Mechanics)*, 13:57–59, 1933. 5
- Vatsal Sharan, Aaron Sidford, and Gregory Valiant. Memory-sample tradeoffs for linear regression with small error. In *Proceedings of the 51st Annual ACM SIGACT Symposium on Theory of Computing*, pp. 890–901, 2019. 8
- Kai Shen, Junliang Guo, Xu Tan, Siliang Tang, Rui Wang, and Jiang Bian. A study on relu and softmax in transformer, 2023. 4
- Shaden Smith, Mostofa Patwary, Brandon Norick, Patrick LeGresley, Samyam Rajbhandari, Jared Casper, Zhun Liu, Shrimai Prabhumoye, George Zerveas, Vijay Korthikanti, Elton Zhang, Rewon Child, Reza Yazdani Aminabadi, Julie Bernauer, Xia Song, Mohammad Shoeybi, Yuxiong He, Michael Houston, Saurabh Tiwary, and Bryan Catanzaro. Using deepspeed and megatron to train megatron-turing nlg 530b, a large-scale generative language model, 2022. 1

- Torsten Soderstrom and G. W. Stewart. On the numerical properties of an iterative method for computing the moore- penrose generalized inverse. *SIAM Journal on Numerical Analysis*, 11(1): 61–74, 1974. ISSN 00361429. URL <http://www.jstor.org/stable/2156431>. 5
- Hongjin Su, Jungo Kasai, Chen Henry Wu, Weijia Shi, Tianlu Wang, Jiayi Xin, Rui Zhang, Mari Ostendorf, Luke Zettlemoyer, Noah A. Smith, and Tao Yu. Selective annotation makes language models better few-shot learners. In *The Eleventh International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=qY1h1v7gwg>. 2
- Davoud Ataee Tarzanagh, Yingcong Li, Christos Thrampoulidis, and Samet Oymak. Transformers as support vector machines. *ArXiv*, abs/2308.16898, 2023a. 2, 9
- Davoud Ataee Tarzanagh, Yingcong Li, Xuechen Zhang, and Samet Oymak. Max-margin token selection in attention mechanism, 2023b. 2, 9
- Romal Thoppilan, Daniel De Freitas, Jamie Hall, Noam Shazeer, Apoorv Kulshreshtha, Heng-Tze Cheng, Alicia Jin, Taylor Bos, Leslie Baker, Yu Du, YaGuang Li, Hongrae Lee, Huaixiu Steven Zheng, Amin Ghafouri, Marcelo Menegali, Yanping Huang, Maxim Krikun, Dmitry Lepikhin, James Qin, Dehao Chen, Yuanzhong Xu, Zhifeng Chen, Adam Roberts, Maarten Bosma, Vincent Zhao, Yanqi Zhou, Chung-Ching Chang, Igor Krivokon, Will Rusch, Marc Pickett, Pranesh Srinivasan, Laichee Man, Kathleen Meier-Hellstern, Meredith Ringel Morris, Tulsee Doshi, Renelito Delos Santos, Toju Duke, Johnny Soraker, Ben Zevenbergen, Vinodkumar Prabhakaran, Mark Diaz, Ben Hutchinson, Kristen Olson, Alejandra Molina, Erin Hoffman-John, Josh Lee, Lora Aroyo, Ravi Rajakumar, Alena Butryna, Matthew Lamm, Viktoriya Kuzmina, Joe Fenton, Aaron Cohen, Rachel Bernstein, Ray Kurzweil, Blaise Agüera-Arcas, Claire Cui, Marian Croak, Ed Chi, and Quoc Le. Lamda: Language models for dialog applications, 2022. 1
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Ł ukasz Kaiser, and Illia Polosukhin. Attention is all you need. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017. URL https://proceedings.neurips.cc/paper_files/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf. 1, 4
- Johannes von Oswald, Eyvind Niklasson, E. Randazzo, João Sacramento, Alexander Mordvintsev, Andrey Zhmoginov, and Max Vladymyrov. Transformers learn in-context by gradient descent. In *International Conference on Machine Learning*, 2022. 1, 3, 6
- Johannes von Oswald, Eyvind Niklasson, Maximilian Schlegel, Seijin Kobayashi, Nicolas Zucchet, Nino Scherrer, Nolan Miller, Mark Sandler, Blaise Agüera y Arcas, Max Vladymyrov, Razvan Pascanu, and Joao Sacramento. Uncovering mesa-optimization algorithms in transformers. *ArXiv*, abs/2309.05858, 2023. 1, 3
- Kevin Wang, Alexandre Variengien, Arthur Conmy, Buck Shlegeris, and Jacob Steinhardt. Interpretability in the wild: a circuit for indirect object identification in gpt-2 small, 2022. 3
- Jerry Wei, Jason Wei, Yi Tay, Dustin Tran, Albert Webson, Yifeng Lu, Xinyun Chen, Hanxiao Liu, Da Huang, Denny Zhou, and Tengyu Ma. Larger language models do in-context learning differently, 2023. 2
- Mitchell Wortsman, Jaehoon Lee, Justin Gilmer, and Simon Kornblith. Replacing softmax with relu in vision transformers, 2023. 4
- Kang Min Yoo, Junyeob Kim, Hyuhng Joon Kim, Hyunsoo Cho, Hwiyeol Jo, Sang-Woo Lee, Sang-goo Lee, and Taeuk Kim. Ground-truth labels matter: A deeper look into input-label demonstrations. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pp. 2422–2437, Abu Dhabi, United Arab Emirates, December 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.emnlp-main.155. URL <https://aclanthology.org/2022.emnlp-main.155>. 2
- Ruiqi Zhang, Spencer Frei, and Peter L. Bartlett. Trained transformers learn linear models in-context. *ArXiv*, abs/2306.09927, 2023. 2

Zihao Zhao, Eric Wallace, Shi Feng, Dan Klein, and Sameer Singh. Calibrate before use: Improving few-shot performance of language models. In *International Conference on Machine Learning*, pp. 12697–12706. PMLR, 2021. 2

APPENDIX

A ADDITIONAL EXPERIMENTAL RESULTS

A.1 ADDITIONAL RESULTS ON ISOTROPIC DATA

A.1.1 HEATMAPS

We present heatmaps with all values of similarities.

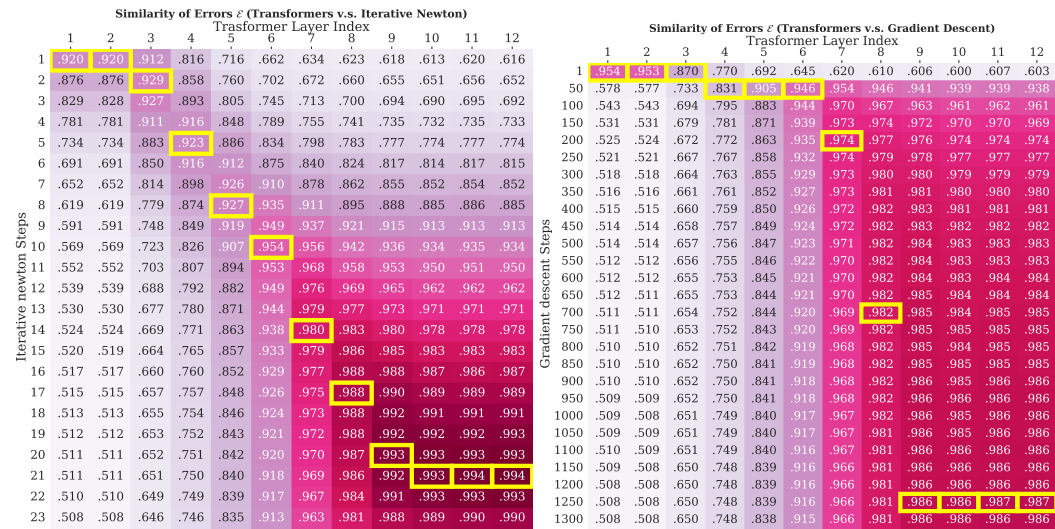


Figure 6: **Similarity of Errors.** The best matching hyper-parameters are highlighted in yellow.

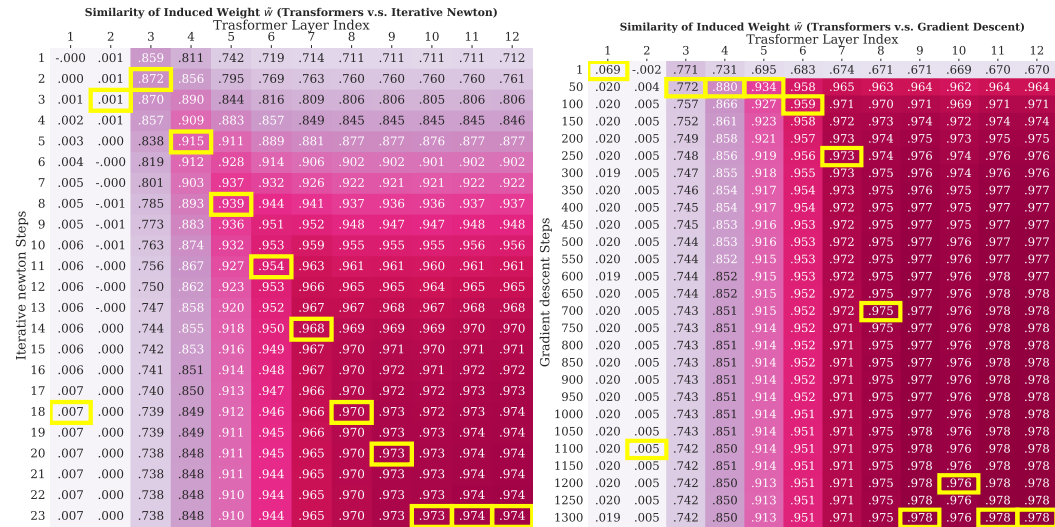


Figure 7: **Similarity of Induced Weight Vectors.** The best matching hyper-parameters are highlighted in yellow.

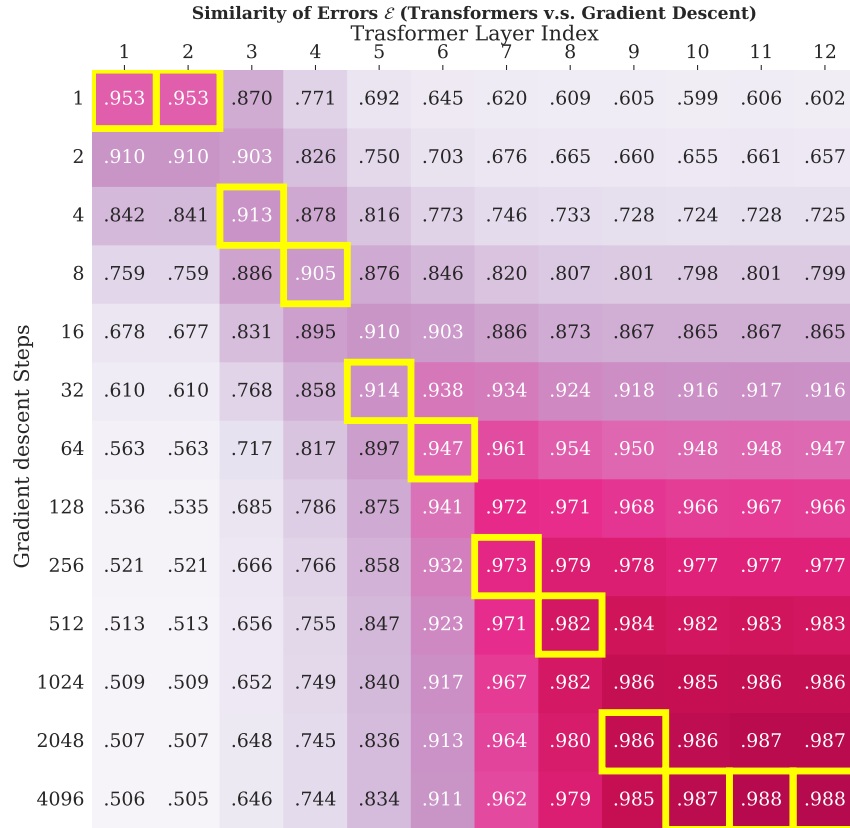


Figure 8: **Similarity of Errors of Gradient Descent in Log Scale.** The best matching hyper-parameters are highlighted in yellow. Putting the number of steps of Gradient Descent in log scale further verifies the claim that Transformer’s rate of coverage is exponentially faster than that of Gradient Descent.

A.1.2 ADDITIONAL RESULTS ON COMPARISON OVER TRANSFORMER LAYERS

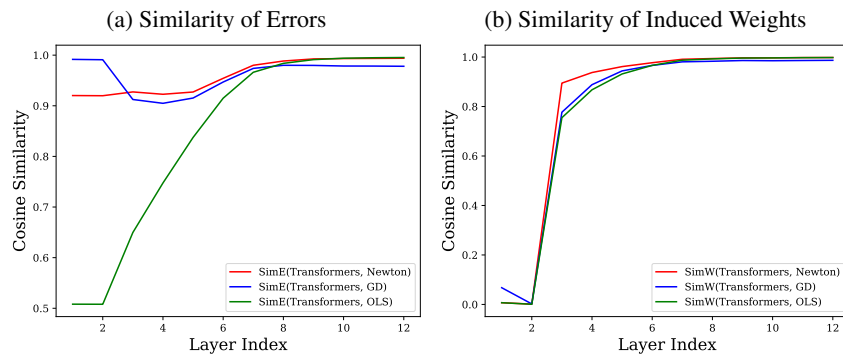


Figure 9: Similarities between Transformer and candidate algorithms. Transformers resemble *Iterative Newton’s Method* the most.

A.1.3 ADDITIONAL RESULTS ON SIMILARITY OF INDUCED WEIGHTS

We present more details line plots for how the similarity of weights changes as the models see more in-context observations $\{\mathbf{x}_i, y_i\}_{i=1}^n$, i.e., as n increases. We fix the number of Transformers layers ℓ and compare with other algorithms with their best-match hyperparameters to ℓ in Figure 10.

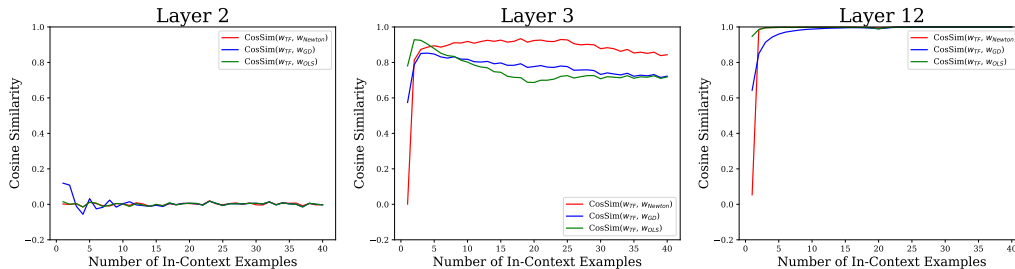


Figure 10: *Similarity of induced weights* over varying number of in-context examples, on three layer indices of Transformers, indexed as 2, 3 and 12. We find that initially at layer 2, the Transformers model hasn't learned so it has zero similarity to all candidate algorithms. As we progress to the next layer number 3, we find that Transformers start to learn, and when provided few examples, Transformers are more similar to OLS but soon become most similar to the Iterative Newton's Method. Layer 12 shows that Transformers in the later layers converge to the OLS solution when provided more than 1 example. We also find there is a dip around $n = d$ for similarity between Transformers and OLS but not for Transformers and Newton, and this is probably because OLS has a more prominent double-descent phenomenon than Transformers and Newton.

A.2 EXPERIMENTS ON ILL-CONDITIONED PROBLEMS

In this section, we repeat the same experiments as we did on isotropic data in the main text and in Appendix A.1, and we change the covariance matrix to be ill-conditioned such that $\kappa(\Sigma) = 100$.

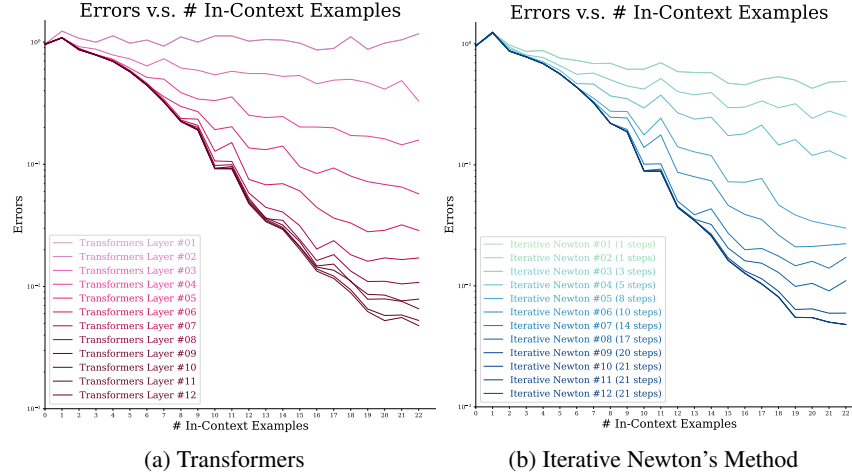


Figure 11: **Progression of Algorithms on Ill-Conditioned Data.** Transformer’s performance still improves over the layer index ℓ ; Iterative Newton’s Method’s performance improves over the number of iterations t and we plot the best-matching t to Transformer’s ℓ following definition 4.

We also present the heatmaps to find the best-matching hyper-parameters and conclude that Transformers are similar to Newton’s method than GD in ill-conditioned data.

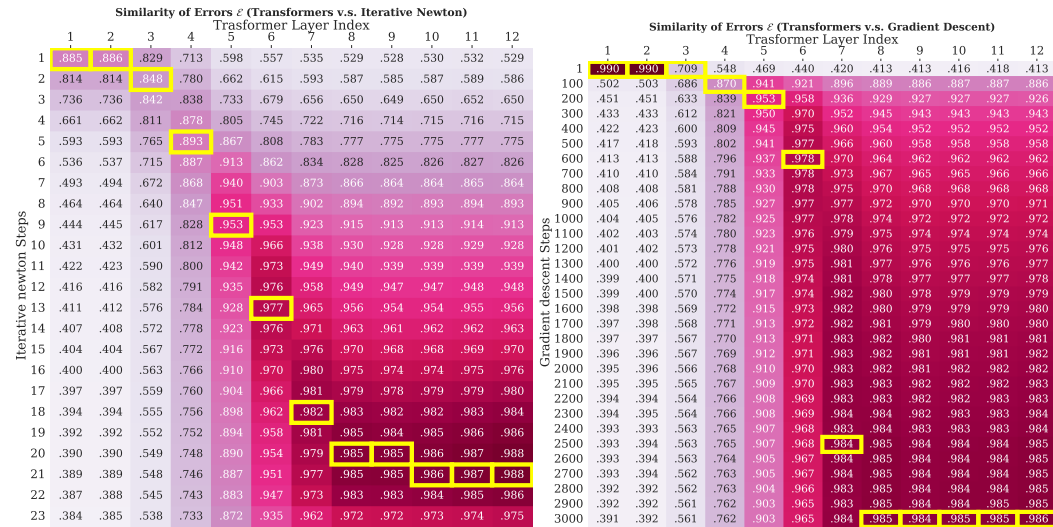


Figure 12: **Similarity of Errors on Ill-Conditioned Data.** The best matching hyper-parameters are highlighted in yellow.

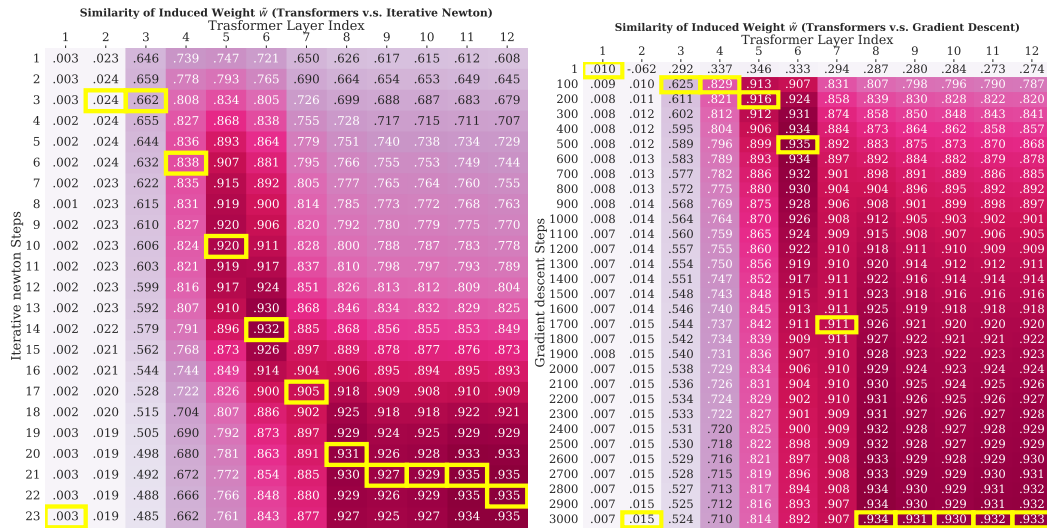


Figure 13: **Similarity of Induced Weights on Ill-Conditioned Data.** The best matching hyper-parameters are highlighted in yellow.

A.3 EXPERIMENTS ON ALTERNATIVE CONFIGURATIONS

In this section, we present experimental results from an alternative model configurations than the main text. We show in the main text that Transformers learn higher-order optimization methods in-context where the experiments are using a GPT-2 model with 12 layers and 8 heads per layer. In this section, we present experiments with a GPT-2 model with 12 layers but only 1 head per layer.

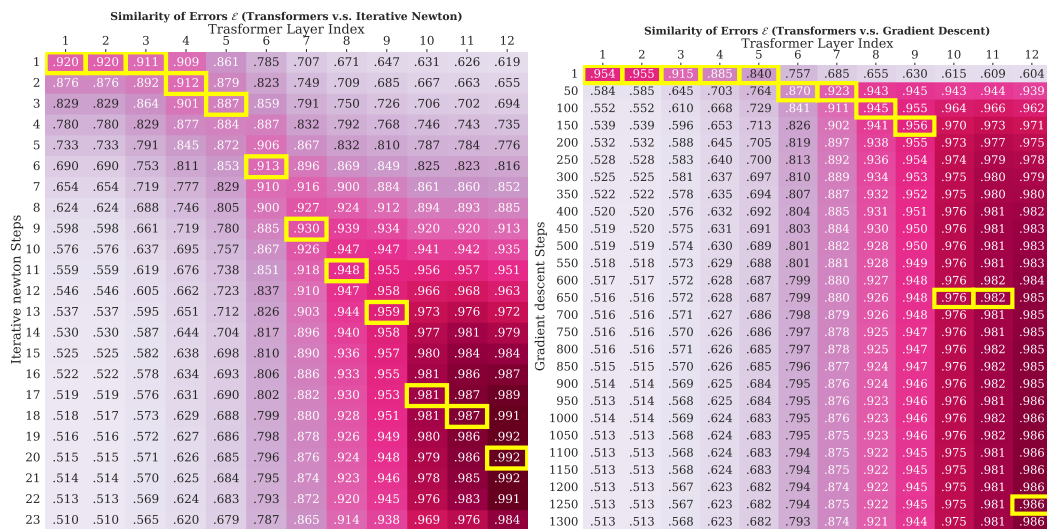


Figure 14: **Similarity of Errors on an alternative Transformers Configuration.** The best matching hyper-parameters are highlighted in yellow.

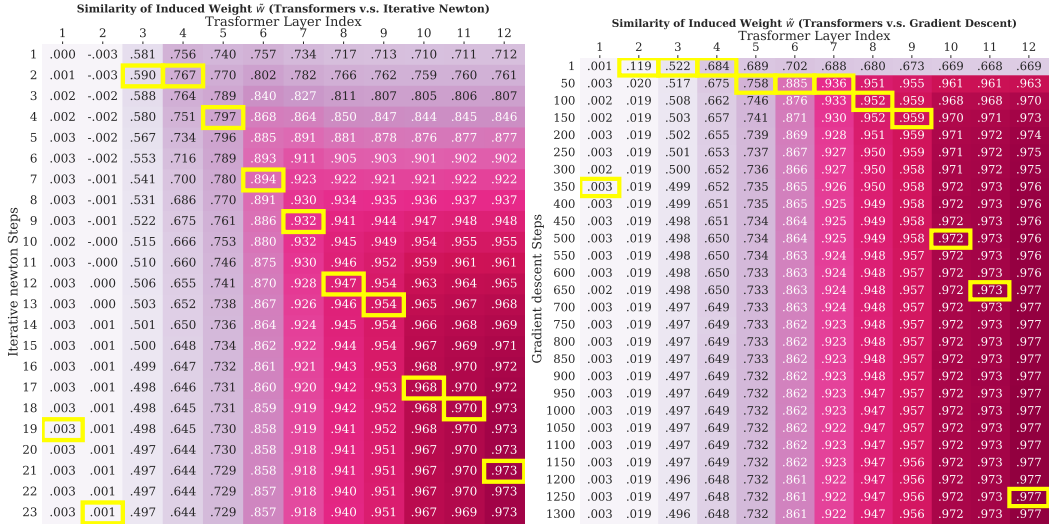


Figure 15: **Similarity of Induced Weights on an alternative Transformers Configuration.** The best matching hyper-parameters are highlighted in yellow.

We conclude that our experimental results are not restricted to a specific model configurations, smaller models such as GPT-2 with 12 layers and 1 head each layer also suffice in implementing the Iterative Newton’s method, and more similar than gradient descents, in terms of rate of convergence.

A.4 DEFINITIONS FOR EVALUATING FORGETTING

We measure the phenomenon of model forgetting by reusing an in-context example within $\{\mathbf{x}_i, y_i\}_{i=1}^n$ as the test example \mathbf{x}_{test} . In experiments of Figure 5, we fix $n = 20$ and reuse $\mathbf{x}_{\text{test}} = \mathbf{x}_i$. We denote the “Time Stamp Gap” as the distance the reused example index i from the current time stamp $n = 20$. We measure the forgetting of index i as

$$\text{Forgetting}(\mathcal{A}, i) = \mathbb{E}_{\{\mathbf{x}_i, y_i\}_{i=1}^n \sim P_D} \text{MSE}\left(\mathcal{A}(\mathbf{x}_i \mid \{\mathbf{x}_i, y_i\}_{i=1}^n), y_i\right) \tag{16}$$

Note: the further away i is from n , the more possible algorithm \mathcal{A} forgets.

B DETAILED PROOFS FOR SECTION 5

In this section, we work on full attention layers with normalized ReLU activation $\sigma(\cdot) = \frac{1}{n} \text{ReLU}(\cdot)$ given n examples.

Definition 5. A full attention layer with M heads and ReLU activation is also denoted as Attn on any input sequence $\mathbf{H} = [\mathbf{h}_1, \dots, \mathbf{h}_N] \in \mathbb{R}^{D \times N}$, where D is the dimension of hidden states and N is the sequence length. In the vector form,

$$\tilde{\mathbf{h}}_t = [\text{Attn}(\mathbf{H})]_t = \mathbf{h}_t + \frac{1}{n} \sum_{m=1}^M \sum_{j=1}^n \text{ReLU}(\langle \mathbf{Q}_m \mathbf{h}_t, \mathbf{K}_m \mathbf{h}_j \rangle) \cdot \mathbf{V}_m \mathbf{h}_j \tag{17}$$

Remark 1. This is slightly different from the *causal* attention layer (see definition 1) in that at each time stamp t , the attention layer in definition 5 has full information of all hidden states $j \in [n]$, unlike causal attention layer which requires $j \in [t]$.

B.1 HELPER RESULTS

We begin by constructing a useful component for our proof, and state some existing constructions from Akyurek et al. (2022).

Lemma 1. Given hidden states $\{\mathbf{h}_1, \dots, \mathbf{h}_n\}$, there exists query, key and value matrices $\mathbf{Q}, \mathbf{K}, \mathbf{V}$ respectively such that one attention layer can compute $\sum_{j=1}^n \mathbf{h}_j$.

Proof. We can pad each hidden state by 1 and 0's such that $\mathbf{h}'_t \leftarrow \begin{bmatrix} \mathbf{h}_t \\ 1 \\ \mathbf{0}_d \end{bmatrix} \in \mathbb{R}^{2d+1}$. We con-

struct two heads where $\mathbf{Q}_1 = \mathbf{K}_1 = \mathbf{Q}_2 = \begin{bmatrix} \mathbf{O}_{d \times d} & \mathbf{O}_{d \times 1} & \mathbf{O}_{d \times d} \\ \mathbf{O}_{1 \times d} & 1 & \mathbf{O}_{1 \times d} \\ \mathbf{O}_{d \times d} & \mathbf{O}_{d \times 1} & \mathbf{O}_{d \times d} \end{bmatrix}$ and $\mathbf{K}_2 = -\mathbf{K}_1$. Then

$$\begin{bmatrix} \mathbf{O}_{d \times d} & \mathbf{O}_{d \times 1} & \mathbf{O}_{d \times d} \\ \mathbf{O}_{1 \times d} & 1 & \mathbf{O}_{1 \times d} \\ \mathbf{O}_{d \times d} & \mathbf{O}_{d \times 1} & \mathbf{O}_{d \times d} \end{bmatrix} \mathbf{h}'_t = \begin{bmatrix} \mathbf{0}_d \\ 1 \\ \mathbf{0}_d \end{bmatrix}.$$

Let $\mathbf{V}_1 = \mathbf{V}_2 = \begin{bmatrix} \mathbf{O}_{(d+1) \times d} & \mathbf{O}_{(d+1) \times (d+1)} \\ n\mathbf{I}_{d \times d} & \mathbf{O}_{d \times (d+1)} \end{bmatrix}$ so that $\mathbf{V}_m \begin{bmatrix} \mathbf{h}_j \\ 1 \\ \mathbf{0}_d \end{bmatrix} = \begin{bmatrix} \mathbf{0}_{d+1} \\ n\mathbf{h}_j \end{bmatrix}$.

We apply one attention layer to these 1-padded hidden states and we have

$$\begin{aligned} \tilde{\mathbf{h}}_t &= \mathbf{h}'_t + \frac{1}{n} \sum_{m=1}^2 \sum_{j=1}^n \text{ReLU}(\langle \mathbf{Q}_m \mathbf{h}'_t, \mathbf{K}_m \mathbf{h}'_j \rangle) \cdot \mathbf{V}_m \mathbf{h}'_j \\ &= \mathbf{h}'_t + \frac{1}{n} \sum_{j=1}^n [\text{ReLU}(1) + \text{ReLU}(-1)] \cdot \begin{bmatrix} \mathbf{0}_{d+1} \\ n\mathbf{h}_j \end{bmatrix} \\ &= \begin{bmatrix} \mathbf{h}_t \\ 1 \\ \mathbf{0}_d \end{bmatrix} + \begin{bmatrix} \mathbf{0}_{d+1} \\ \sum_{j=1}^n \mathbf{h}_j \end{bmatrix} = \begin{bmatrix} \mathbf{h}_t \\ 1 \\ \sum_{j=1}^n \mathbf{h}_j \end{bmatrix} \end{aligned} \quad (18)$$

□

Proposition 1 (Akyürek et al., 2022). Each of `mov`, `aff`, `mul`, `div` can be implemented by a single transformer layer. These four operations are mappings $\mathbb{R}^{D \times N} \rightarrow \mathbb{R}^{D \times N}$, expressed as follows,

`mov`($\mathbf{H}; s, t, i, j, i', j'$): selects the entries of the s -th column of \mathbf{H} between rows i and j , and copies them into the t -th column ($t \geq s$) of \mathbf{H} between rows i' and j' .

`mul`($\mathbf{H}; a, b, c, (i, j), (i', j'), (i'', j'')$): in each column \mathbf{h} of \mathbf{H} , interprets the entries between i and j as an $a \times b$ matrix \mathbf{A}_1 , and the entries between i' and j' as a $b \times c$ matrix \mathbf{A}_2 , multiplies these matrices together, and stores the result between rows i'' and j'' , yielding a matrix in which each column has the form $[\mathbf{h}_{:i'-1}, \mathbf{A}_1 \mathbf{A}_2, \mathbf{h}_{j''}]^\top$. This allows the layer to implement inner products.

`div`($\mathbf{H}; (i, j), i', (i'', j'')$): in each column \mathbf{h} of \mathbf{H} , divides the entries between i and j by the absolute value of the entry at i' , and stores the result between rows i'' and j'' , yielding a matrix in which every column has the form $[\mathbf{h}_{:i'-1}, \mathbf{h}_{i:j}/|\mathbf{h}_{i'}|, \mathbf{h}_{j''}]^\top$.

`aff`($\mathbf{H}; (i, j), (i', j'), (i'', j''), \mathbf{W}_1, \mathbf{W}_2, \mathbf{b}$): in each column \mathbf{h} of \mathbf{H} , applies an affine transformation to the entries between i and j and i' and j' , then stores the result between rows i'' and j'' , yielding a matrix in which every column has the form $[\mathbf{h}_{:i'-1}, \mathbf{W}_1 \mathbf{h}_{i:j} + \mathbf{W}_2 \mathbf{h}_{i':j'} + \mathbf{b}, \mathbf{h}_{j''}]^\top$. This allows the layer to implement subtraction by setting $\mathbf{W}_1 = \mathbf{I}$ and $\mathbf{W}_2 = -\mathbf{I}$.

B.2 PROOF OF THEOREM 1

Theorem 1. There exist Transformer weights such that on any set of in-context examples $\{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^n$ and test point \mathbf{x}_{test} , the Transformer predicts on \mathbf{x}_{test} using $\mathbf{x}_{\text{test}}^\top \hat{\mathbf{w}}_k^{\text{Newton}}$. Here $\hat{\mathbf{w}}_k^{\text{Newton}}$ are the Iterative Newton updates given by $\hat{\mathbf{w}}_k^{\text{Newton}} = \mathbf{M}_k \mathbf{X}^\top \mathbf{y}$ where \mathbf{M}_j is updated as

$$\mathbf{M}_j = 2\mathbf{M}_{j-1} - \mathbf{M}_{j-1} \mathbf{S} \mathbf{M}_{j-1}, 1 \leq j \leq k, \quad \mathbf{M}_0 = \alpha \mathbf{S},$$

for some $\alpha > 0$ and $\mathbf{S} = \mathbf{X}^\top \mathbf{X}$. The number of layers of the transformer is $\mathcal{O}(k)$ and the dimensionality of the hidden layers is $\mathcal{O}(d)$.

Proof. We break the proof into parts.

Transformers Implement Initialization $T^{(0)} = \alpha S$. Given input sequence $\mathbf{H} := \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$, with $\mathbf{x}_i \in \mathbb{R}^d$, we first apply the `mov` operations given by Proposition 1 (similar to Akyürek et al. (2022)), we show only non-zero rows when applying these operations):

$$\begin{bmatrix} \mathbf{x}_1 & \cdots & \mathbf{x}_n \end{bmatrix} \xrightarrow{\text{mov}} \begin{bmatrix} \mathbf{x}_1 & \cdots & \mathbf{x}_n \\ \mathbf{x}_1 & \cdots & \mathbf{x}_n \end{bmatrix} \quad (19)$$

We call each column after `mov` as \mathbf{h}_j . With an full attention layer, one can construct two heads with query and value matrices of the form $\mathbf{Q}_1^\top \mathbf{K}_1 = -\mathbf{Q}_2^\top \mathbf{K}_2 = \begin{bmatrix} \mathbf{I}_{d \times d} & \mathbf{O}_{d \times d} \\ \mathbf{O}_{d \times d} & \mathbf{O}_{d \times d} \end{bmatrix}$ such that for any $t \in [n]$, we have

$$\sum_{m=1}^2 \text{ReLU}(\langle \mathbf{Q}_m \mathbf{h}_t, \mathbf{K}_m \mathbf{h}_j \rangle) = \text{ReLU}(\mathbf{x}_t^\top \mathbf{x}_j) + \text{ReLU}(-\mathbf{x}_t^\top \mathbf{x}_j) = \langle \mathbf{x}_t, \mathbf{x}_j \rangle \quad (20)$$

Let all value matrices $\mathbf{V}_m = n\alpha \begin{bmatrix} \mathbf{I}_{d \times d} & \mathbf{O}_{d \times d} \\ \mathbf{O}_{d \times d} & \mathbf{O}_{d \times d} \end{bmatrix}$ for some $\alpha \in \mathbb{R}$. Combining the skip connections, we have

$$\tilde{\mathbf{h}}_t = \begin{bmatrix} \mathbf{x}_t \\ \mathbf{x}_t \end{bmatrix} + \frac{1}{n} \sum_{j=1}^n \langle \mathbf{x}_t, \mathbf{x}_j \rangle n\alpha \begin{bmatrix} \mathbf{x}_j \\ \mathbf{0} \end{bmatrix} = \begin{bmatrix} \mathbf{x}_t \\ \mathbf{x}_t \end{bmatrix} + \begin{bmatrix} \alpha \left(\sum_{j=1}^n \mathbf{x}_j \mathbf{x}_j^\top \right) \mathbf{x}_t \\ \mathbf{0} \end{bmatrix} = \begin{bmatrix} \mathbf{x}_t + \alpha \mathbf{S} \mathbf{x}_t \\ \mathbf{x}_t \end{bmatrix} \quad (21)$$

Now we can use the `aff` operator to make subtractions and then

$$\begin{bmatrix} \mathbf{x}_t + \alpha \mathbf{S} \mathbf{x}_t \\ \mathbf{x}_t \end{bmatrix} \xrightarrow{\text{aff}} \begin{bmatrix} (\mathbf{x}_t + \alpha \mathbf{S} \mathbf{x}_t) - \mathbf{x}_t \\ \mathbf{x}_t \end{bmatrix} = \begin{bmatrix} \alpha \mathbf{S} \mathbf{x}_t \\ \mathbf{x}_t \end{bmatrix} \quad (22)$$

We call this transformed hidden states as $\mathbf{H}^{(0)}$ and denote $T^{(0)} = \alpha S$:

$$\mathbf{H}^{(0)} = \begin{bmatrix} \mathbf{h}_1^{(0)} & \cdots & \mathbf{h}_n^{(0)} \end{bmatrix} = \begin{bmatrix} T^{(0)} \mathbf{x}_1 & \cdots & T^{(0)} \mathbf{x}_n \\ \mathbf{x}_1 & \cdots & \mathbf{x}_n \end{bmatrix} \quad (23)$$

Notice that S is symmetric and thereafter $T^{(0)}$ is also symmetric.

Transformers implement Newton Iteration. Let the input prompt be the same as Equation (23),

$$\mathbf{H}^{(0)} = \begin{bmatrix} \mathbf{h}_1^{(0)} & \cdots & \mathbf{h}_n^{(0)} \end{bmatrix} = \begin{bmatrix} T^{(0)} \mathbf{x}_1 & \cdots & T^{(0)} \mathbf{x}_n \\ \mathbf{x}_1 & \cdots & \mathbf{x}_n \end{bmatrix} \quad (24)$$

We claim that the ℓ 's hidden states can be of the similar form

$$\mathbf{H}^{(\ell)} = \begin{bmatrix} \mathbf{h}_1^{(\ell)} & \cdots & \mathbf{h}_n^{(\ell)} \end{bmatrix} = \begin{bmatrix} T^{(\ell)} \mathbf{x}_1 & \cdots & T^{(\ell)} \mathbf{x}_n \\ \mathbf{x}_1 & \cdots & \mathbf{x}_n \end{bmatrix} \quad (25)$$

We prove by induction that assuming our claim is true for ℓ , we work on $\ell + 1$:

Let $\mathbf{Q}_m = \underbrace{\tilde{\mathbf{Q}}_m \begin{bmatrix} \mathbf{O}_d & -\frac{n}{2} \mathbf{I}_d \\ \mathbf{O}_d & \mathbf{O}_d \end{bmatrix}}_G$, $\mathbf{K}_m = \underbrace{\tilde{\mathbf{K}}_m \begin{bmatrix} \mathbf{I}_d & \mathbf{O}_d \\ \mathbf{O}_d & \mathbf{O}_d \end{bmatrix}}_J$ where $\tilde{\mathbf{Q}}_1^\top \tilde{\mathbf{K}}_1 := \mathbf{I}$, $\tilde{\mathbf{Q}}_2^\top \tilde{\mathbf{K}}_2 := -\mathbf{I}$ and $\mathbf{V}_1 = \mathbf{V}_2 = \underbrace{\begin{bmatrix} \mathbf{I}_d & \mathbf{O}_d \\ \mathbf{O}_d & \mathbf{O}_d \end{bmatrix}}_J$. A 2-head self-attention layer, with ReLU attentions, can be written has

$$\mathbf{h}_t^{(\ell+1)} = [\text{Attn}(\mathbf{H}^{(\ell)})]_t = \mathbf{h}_t^{(\ell)} + \frac{1}{n} \sum_{m=1}^2 \sum_{j=1}^n \text{ReLU}(\langle \mathbf{Q}_m \mathbf{h}_t^{(\ell)}, \mathbf{K}_m \mathbf{h}_j^{(\ell)} \rangle) \cdot \mathbf{V}_m \mathbf{h}_j^{(\ell)} \quad (26)$$

where

$$\begin{aligned}
& \sum_{m=1}^2 \text{ReLU} \left(\left\langle \mathbf{Q}_m \mathbf{h}_t^{(\ell)}, \mathbf{K}_m \mathbf{h}_j^{(\ell)} \right\rangle \right) \cdot \mathbf{V}_m \mathbf{h}_j^{(\ell)} \\
&= \left[\text{ReLU} \left((\mathbf{G} \mathbf{h}_t^{(\ell)})^\top \underbrace{\tilde{\mathbf{Q}}_1^\top \tilde{\mathbf{K}}_1}_{\mathbf{I}} (\mathbf{J} \mathbf{h}_j^{(\ell)}) \right) + \text{ReLU} \left((\mathbf{G} \mathbf{h}_t^{(\ell)})^\top \underbrace{\tilde{\mathbf{Q}}_2^\top \tilde{\mathbf{K}}_2}_{-\mathbf{I}} (\mathbf{J} \mathbf{h}_j^{(\ell)}) \right) \right] \cdot (\mathbf{J} \mathbf{h}_j^{(\ell)}) \\
&= \left[\text{ReLU} \left((\mathbf{G} \mathbf{h}_t^{(\ell)})^\top (\mathbf{J} \mathbf{h}_j^{(\ell)}) \right) + \text{ReLU} \left(-(\mathbf{G} \mathbf{h}_t^{(\ell)})^\top (\mathbf{J} \mathbf{h}_j^{(\ell)}) \right) \right] \cdot (\mathbf{J} \mathbf{h}_j^{(\ell)}) \\
&= (\mathbf{G} \mathbf{h}_t^{(\ell)})^\top (\mathbf{J} \mathbf{h}_j^{(\ell)}) (\mathbf{J} \mathbf{h}_j^{(\ell)}) \\
&= (\mathbf{J} \mathbf{h}_j^{(\ell)}) (\mathbf{J} \mathbf{h}_j^{(\ell)})^\top (\mathbf{G} \mathbf{h}_t^{(\ell)})
\end{aligned} \tag{27}$$

Plug in our assumptions that $\mathbf{h}_j^{(\ell)} = \begin{bmatrix} \mathbf{T}^{(\ell)} \mathbf{x}_j \\ \mathbf{x}_j \end{bmatrix}$, we have $\mathbf{J} \mathbf{h}_j^{(\ell)} = \begin{bmatrix} \mathbf{T}^{(\ell)} \mathbf{x}_j \\ \mathbf{0}_d \end{bmatrix}$ and $\mathbf{G} \mathbf{h}_t^{(\ell)} = \begin{bmatrix} -\frac{n}{2} \mathbf{x}_t \\ \mathbf{0}_d \end{bmatrix}$, we have

$$\begin{aligned}
\mathbf{h}_t^{(\ell+1)} &= \begin{bmatrix} \mathbf{T}^{(\ell)} \mathbf{x}_t \\ \mathbf{x}_t \end{bmatrix} + \frac{1}{n} \sum_{j=1}^n \begin{bmatrix} \mathbf{T}^{(\ell)} \mathbf{x}_j \\ \mathbf{0}_d \end{bmatrix} \begin{bmatrix} \mathbf{T}^{(\ell)} \mathbf{x}_j \\ \mathbf{0}_d \end{bmatrix}^\top \begin{bmatrix} -\frac{n}{2} \mathbf{x}_t \\ \mathbf{0}_d \end{bmatrix} \\
&= \begin{bmatrix} \mathbf{T}^{(\ell)} \mathbf{x}_t - \frac{1}{2} \sum_{j=1}^n (\mathbf{T}^{(\ell)} \mathbf{x}_j) (\mathbf{T}^{(\ell)} \mathbf{x}_j)^\top \mathbf{x}_t \\ \mathbf{x}_t \end{bmatrix} \\
&= \begin{bmatrix} \mathbf{T}^{(\ell)} \mathbf{x}_t - \frac{1}{2} \mathbf{T}^{(\ell)} \left(\sum_{j=1}^n \mathbf{x}_j \mathbf{x}_j^\top \right) \mathbf{T}^{(\ell)\top} \mathbf{x}_t \\ \mathbf{x}_t \end{bmatrix} \\
&= \begin{bmatrix} \left(\mathbf{T}^{(\ell)} - \frac{1}{2} \mathbf{T}^{(\ell)} \mathbf{S} \mathbf{T}^{(\ell)\top} \right) \mathbf{x}_t \\ \mathbf{x}_t \end{bmatrix}
\end{aligned} \tag{28}$$

Now we pass over an MLP layer with

$$\mathbf{h}_t^{(\ell+1)} \leftarrow \mathbf{h}_t^{(\ell+1)} + \begin{bmatrix} \mathbf{I}_d & \mathbf{O}_d \\ \mathbf{O}_d & \mathbf{O}_d \end{bmatrix} \mathbf{h}_t^{(\ell+1)} = \begin{bmatrix} \left(2\mathbf{T}^{(\ell)} - \mathbf{T}^{(\ell)} \mathbf{S} \mathbf{T}^{(\ell)\top} \right) \mathbf{x}_t \\ \mathbf{x}_t \end{bmatrix} \tag{29}$$

Now we denote the iteration

$$\mathbf{T}^{(\ell+1)} = 2\mathbf{T}^{(\ell)} - \mathbf{T}^{(\ell)} \mathbf{S} \mathbf{T}^{(\ell)\top} \tag{30}$$

We find that $\mathbf{T}^{(\ell+1)\top} = \mathbf{T}^{(\ell+1)}$ since $\mathbf{T}^{(\ell)}$ and \mathbf{S} are both symmetric. It reduces to

$$\mathbf{T}^{(\ell+1)} = 2\mathbf{T}^{(\ell)} - \mathbf{T}^{(\ell)} \mathbf{S} \mathbf{T}^{(\ell)} \tag{31}$$

This is exactly the same as the Newton iteration.

Transformers can implement $\hat{\mathbf{w}}_\ell^{\text{TF}} = \mathbf{T}^{(\ell)} \mathbf{X}^\top \mathbf{y}$. Going back to the empirical prompt format $\{\mathbf{x}_1, y_1, \dots, \mathbf{x}_n, y_n\}$. We can let parameters be zero for positions of y 's and only rely on the skip connection up to layer ℓ , and the $\mathbf{H}^{(\ell)}$ is then $\begin{bmatrix} \mathbf{T}^{(\ell)} \mathbf{x}_j & \mathbf{0} \\ \mathbf{x}_j & \mathbf{0} \\ 0 & y_j \end{bmatrix}_{j=1}^n$. We again apply operations from

Proposition 1:

$$\begin{bmatrix} \mathbf{T}^{(\ell)} \mathbf{x}_j & \mathbf{0} \\ \mathbf{x}_j & \mathbf{0} \\ 0 & y_j \end{bmatrix}_{j=1}^n \xrightarrow{\text{mov}} \begin{bmatrix} \mathbf{T}^{(\ell)} \mathbf{x}_j & \mathbf{T}^{(\ell)} \mathbf{x}_j \\ \mathbf{x}_j & \mathbf{0} \\ 0 & y_j \end{bmatrix}_{j=1}^n \xrightarrow{\text{mul}} \begin{bmatrix} \mathbf{T}^{(\ell)} \mathbf{x}_j & \mathbf{T}^{(\ell)} \mathbf{x}_j \\ \mathbf{x}_j & \mathbf{0} \\ \mathbf{0} & y_j \\ \mathbf{0} & \mathbf{T}^{(\ell)} y_j \mathbf{x}_j \end{bmatrix}_{j=1}^n \tag{32}$$

Now we apply Lemma 1 over all even columns in Equation (32) and we have

$$\text{Output} = \sum_{j=1}^n \begin{bmatrix} \mathbf{T}^{(\ell)} \mathbf{x}_j \\ \mathbf{0} \\ y_j \\ \mathbf{T}^{(\ell)} y_j \mathbf{x}_j \end{bmatrix} = \begin{bmatrix} \boldsymbol{\xi} \\ \mathbf{T}^{(\ell)} \sum_{j=1}^n y_j \mathbf{x}_j \end{bmatrix} = \begin{bmatrix} \boldsymbol{\xi} \\ \mathbf{T}^{(\ell)} \mathbf{X}^\top \mathbf{y} \end{bmatrix} \tag{33}$$

where ξ denotes irrelevant quantities. Note that the resulting $\mathbf{T}^{(\ell)} \mathbf{X}^\top \mathbf{y}$ is also the same as Iterative Newton’s predictor $\hat{\mathbf{w}}_k = \mathbf{M}_k \mathbf{X}^\top \mathbf{y}$ after k iterations. We denote $\hat{\mathbf{w}}_\ell^{\text{TF}} = \mathbf{T}^{(\ell)} \mathbf{X}^\top \mathbf{y}$.

Transformers can make predictions on \mathbf{x}_{test} by $\langle \hat{\mathbf{w}}_\ell^{\text{TF}}, \mathbf{x}_{\text{test}} \rangle$.

Now we can make predictions on text query \mathbf{x}_{test} :

$$\begin{bmatrix} \xi & \mathbf{x}_{\text{test}} \\ \hat{\mathbf{w}}_\ell^{\text{TF}} & \mathbf{x}_{\text{test}} \end{bmatrix} \xrightarrow{\text{mov}} \begin{bmatrix} \xi & \mathbf{x}_{\text{test}} \\ \hat{\mathbf{w}}_\ell^{\text{TF}} & \mathbf{x}_{\text{test}} \\ \mathbf{0} & \hat{\mathbf{w}}_\ell^{\text{TF}} \end{bmatrix} \xrightarrow{\text{mul}} \begin{bmatrix} \xi & \mathbf{x}_{\text{test}} \\ \hat{\mathbf{w}}_\ell^{\text{TF}} & \mathbf{x}_{\text{test}} \\ \mathbf{0} & \hat{\mathbf{w}}_\ell^{\text{TF}} \\ 0 & \langle \hat{\mathbf{w}}_\ell^{\text{TF}}, \mathbf{x}_{\text{test}} \rangle \end{bmatrix} \quad (34)$$

Finally, we can have an readout layer $\beta_{\text{ReadOut}} = \{\mathbf{u}, v\}$ applied (see definition 3) with $\mathbf{u} = [\mathbf{0}_{3d} \ 1]^\top$ and $v = 0$ to extract the prediction $\langle \hat{\mathbf{w}}_\ell^{\text{TF}}, \mathbf{x}_{\text{test}} \rangle$ at the last location, given by \mathbf{x}_{test} . This is exactly how Iterative Newton makes predictions.

To Perform k steps of Newton’s iterations, Transformers need $\mathcal{O}(k)$ layers.

Let’s count the layers:

- **Initialization:** `mov` needs $\mathcal{O}(1)$ layer; gathering $\alpha \mathbf{S}$ needs $\mathcal{O}(1)$ layer; and `aff` needs $\mathcal{O}(1)$ layer. In total, Transformers need $\mathcal{O}(1)$ layers for initialization.
- **Newton Iteration:** each exact Newton’s iteration requires $\mathcal{O}(1)$ layer. Implementing k iterations requires $\mathcal{O}(k)$ layers.
- **Implementing $\hat{\mathbf{w}}_\ell^{\text{TF}}$:** We need one operation of `mov` and `mul` each, requiring $\mathcal{O}(1)$ layer each. Apply Lemma 1 for summation also requires $\mathcal{O}(1)$ layer.
- **Making prediction on test query:** We need one operation of `mov` and `mul` each, requiring $\mathcal{O}(1)$ layer each.

Hence, in total, Transformers can implement k -step Iterative Newton and make predictions accordingly using $\mathcal{O}(k)$ layers. □

B.3 ITERATIVE NEWTON AS A SUM OF MOMENTS METHOD

Recall that Iterative Newton’s method finds \mathbf{S}^\dagger as follows

$$\mathbf{M}_0 = \frac{2}{\underbrace{\|\mathbf{S}\mathbf{S}^\top\|_2}_\alpha} \mathbf{S}^\top, \quad \mathbf{M}_k = 2\mathbf{M}_{k-1} - \mathbf{M}_{k-1} \mathbf{S} \mathbf{M}_{k-1}, \forall k \geq 1. \quad (35)$$

We can expand the iterative equation to moments of \mathbf{S} as follows.

$$\mathbf{M}_1 = 2\mathbf{M}_0 - \mathbf{M}_0 \mathbf{S} \mathbf{M}_0 = 2\alpha \mathbf{S}^\top - 4\alpha^2 \mathbf{S}^\top \mathbf{S} \mathbf{S}^\top = 2\alpha \mathbf{S} - 4\alpha^2 \mathbf{S}^3. \quad (36)$$

Let’s do this one more time for \mathbf{M}_2 .

$$\begin{aligned} \mathbf{M}_2 &= 2\mathbf{M}_1 - \mathbf{M}_1 \mathbf{S} \mathbf{M}_1 \\ &= 2(2\alpha \mathbf{S} - 4\alpha^2 \mathbf{S}^3) - (2\alpha \mathbf{S} - 4\alpha^2 \mathbf{S}^3) \mathbf{S} (2\alpha \mathbf{S} - 4\alpha^2 \mathbf{S}^3) \\ &= 4\alpha \mathbf{S} - 8\alpha^2 \mathbf{S}^3 - 4\alpha^2 \mathbf{S}^3 + 16\alpha^3 \mathbf{S}^5 - 16\alpha^4 \mathbf{S}^7 \\ &= 4\alpha \mathbf{S} - 12\alpha^2 \mathbf{S}^3 + 16\alpha^3 \mathbf{S}^5 - 16\alpha^4 \mathbf{S}^7. \end{aligned} \quad (37)$$

We can see that \mathbf{M}_k are summations of moments of \mathbf{S} , with respect to some pre-defined coefficients from the Newton’s algorithm. Hence Iterative Newton is a special of an algorithm which computes an approximation of the inverse using higher-order moments of the matrix,

$$\mathbf{M}_k = \sum_{s=1}^{2^{k+1}-1} \beta_s \mathbf{S}^s \quad (38)$$

with coefficients $\beta_s \in \mathbb{R}$.

We note that Transformer circuits can represent other sum of moments other than Newton’s method. We can introduce different coefficients β_i than in the proof of Theorem 1 by scaling the value matrices or through the MLP layers.

B.4 ESTIMATED WEIGHT VECTORS LIE IN THE SPAN OF PREVIOUS EXAMPLES

What properties can we infer and verify for the weight vectors which arise from Newton’s method? A straightforward one arises from interpreting any sum of moments method as a kernel method.

We can expand S^s as follows

$$S^s = \left(\sum_{i=1}^t \mathbf{x}_i \mathbf{x}_i^\top \right)^s = \sum_{i=1}^t \left(\sum_{j_1, \dots, j_{s-1}} \langle \mathbf{x}_i, \mathbf{x}_{j_1} \rangle \prod_{v=1}^{s-2} \langle \mathbf{x}_{j_v}, \mathbf{x}_{j_{v+1}} \rangle \right) \mathbf{x}_i \mathbf{x}_{j_{s-1}}^\top. \quad (39)$$

Then we have

$$\begin{aligned} \hat{\mathbf{w}}_t &= M_t \mathbf{X}^\top \mathbf{y} = \sum_{s=1}^{2^{t+1}-1} \beta_s S^s \mathbf{X}^\top \mathbf{y} \\ &= \sum_{s=1}^{2^{t+1}-1} \beta_s \left\{ \sum_{i=1}^t \left(\sum_{j_1, \dots, j_{s-1}} \langle \mathbf{x}_i, \mathbf{x}_{j_1} \rangle \prod_{v=1}^{s-2} \langle \mathbf{x}_{j_v}, \mathbf{x}_{j_{v+1}} \rangle \right) \mathbf{x}_i \mathbf{x}_{j_{s-1}}^\top \right\} \left\{ \sum_{i=1}^t y_i \mathbf{x}_i \right\} \\ &= \sum_{s=1}^{2^{t+1}-1} \beta_s \left(\sum_{i=1}^t \left(\sum_{j_1, \dots, j_s} y_{j_s} \langle \mathbf{x}_i, \mathbf{x}_{j_1} \rangle \prod_{v=1}^{s-1} \langle \mathbf{x}_{j_v}, \mathbf{x}_{j_{v+1}} \rangle \right) \mathbf{x}_i \right) \\ &= \sum_{i=1}^t \underbrace{\left(\sum_{s=1}^{2^{t+1}-1} \sum_{j_1, \dots, j_s} \beta_s y_{j_s} \langle \mathbf{x}_i, \mathbf{x}_{j_1} \rangle \prod_{v=1}^{s-1} \langle \mathbf{x}_{j_v}, \mathbf{x}_{j_{v+1}} \rangle \right)}_{\phi_t(i | \mathbf{X}, \mathbf{y}, \boldsymbol{\beta})} \mathbf{x}_i \\ &= \sum_{i=1}^t \phi_t(i | \mathbf{X}, \mathbf{y}, \boldsymbol{\beta}) \mathbf{x}_i \end{aligned} \quad (40)$$

where \mathbf{X} is the data matrix, $\boldsymbol{\beta}$ are coefficients of moments given by the sum of moments method and $\phi_t(\cdot)$ is some function which assigns some weight to the i -th datapoint, based on all other datapoints. Therefore if the Transformer implements a sum of moments method (such as Newton’s method), then its induced weight vector $\tilde{\mathbf{w}}_t$ (Transformers | $\{\mathbf{x}_i, y_i\}_{i=1}^t$) after seeing in-context examples $\{\mathbf{x}_i, y_i\}_{i=1}^t$ should lie in the span of the examples $\{\mathbf{x}_i\}_{i=1}^t$:

$$\tilde{\mathbf{w}}_t(\text{Transformers} | \{\mathbf{x}_i, y_i\}_{i=1}^t) \stackrel{?}{=} \text{Span}\{\mathbf{x}_1, \dots, \mathbf{x}_t\} = \sum_{i=1}^t a_i \mathbf{x}_i \quad \text{for coefficients } a_i. \quad (41)$$

We test this hypothesis. Given a sequence of in-context examples $\{\mathbf{x}_i, y_i\}_{i=1}^t$, we fit coefficients $\{a_i\}_{i=1}^t$ in Equation (41) to minimize MSE loss:

$$\{\hat{a}_i\}_{i=1}^t = \arg \min_{a_1, a_2, \dots, a_t \in \mathbb{R}} \left\| \tilde{\mathbf{w}}_t(\text{Transformers} | \{\mathbf{x}_i, y_i\}_{i=1}^t) - \sum_{i=1}^t a_i \mathbf{x}_i \right\|_2^2. \quad (42)$$

We then measure the quality of this fit across different number of in-context examples t , and visualize the residual error in Figure 16. We find that even when $t < d$, Transformers’ induced weights still lie close to the span of the observed examples \mathbf{x}_i ’s. This provides an additional validation of our proposed mechanism.

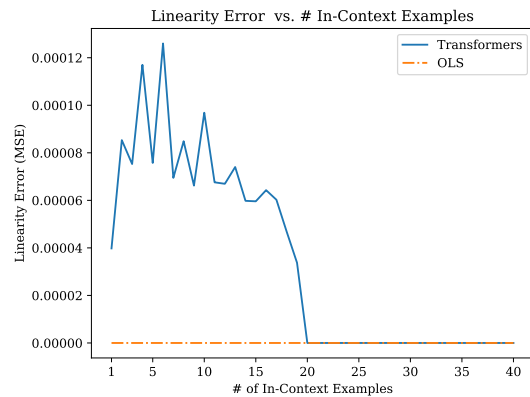


Figure 16: Verification of hypothesis that the Transformers induced weight vector w lies in the span of observed examples $\{x_i\}$.