# Decomposing Prediction Mechanisms
# for In-Context Recall

**Sultan Daniels**[†]    **Dylan Davis**[†]    **Dhruv Gautam**[†]    **Wentinn Liao**[‡]
**Gireeja Ranade**[†]    **Anant Sahai**[†]
[†]University of California, Berkeley    [‡]University of Pennsylvania
{sultan_daniels, dylanjd, dhruvgautam, gireeja, asahai}@berkeley.edu
wenliao@seas.upenn.edu

## Abstract

We introduce a new family of toy problems that combine features of linear-regression-style continuous in-context learning (ICL) with discrete associative recall. We pretrain transformer models on sample traces from this toy, specifically symbolically-labeled interleaved state observations from randomly drawn linear deterministic dynamical systems. We study if the transformer models can recall the state of a sequence previously seen in its context when prompted to do so with the corresponding in-context label. Taking a closer look at this task, it becomes clear that the model must perform two functions: (1) identify which system's state should be recalled and apply that system to its last seen state, and (2) continuing to apply the correct system to predict the subsequent states. Training dynamics reveal that the first capability emerges well into a model's training. Surprisingly, the second capability, of continuing the prediction of a resumed sequence, develops much earlier.

Via out-of-distribution experiments, and a mechanistic analysis on model weights via edge pruning, we find that next-token prediction for this toy problem involves at least two separate mechanisms. One mechanism uses the discrete symbolic labels to do the associative recall required to predict the start of a resumption of a previously seen sequence. The second mechanism, which is largely agnostic to the discrete symbolic labels, performs a "Bayesian-style" prediction based on the previous token and the context. These two mechanisms have different learning dynamics.

To confirm that this multi-mechanism (manifesting as separate phase transitions) phenomenon is not just an artifact of our toy setting, we used OLMo training checkpoints on an ICL translation task to see a similar phenomenon: a decisive gap in the emergence of first-task-token performance vs second-task-token performance.

## 1   Introduction

The release of GPT-3 [7] demonstrated the power of Large Language Models' (LLMs) ability to do in-context learning (ICL). Since then, there has been significant progress in understanding ICL for language models themselves [56, 1, 83, 33, 80, 84, 81, 57, 47]. There has also been work that focuses on understanding ICL for simpler toy problems [20, 64, 17, 71, 65, 15, 54]. Toys (e.g., linear regression [20, 65, 25]) allow us to study the learned ICL behavior of deep neural networks in settings where optimal strategies are known, allowing complex prediction mechanisms to be disentangled. In this paper, we build on previous work to create a new toy problem involving interleaved vector-valued time-series.

We start with underlying time-series that come from the evolution of random deterministic linear systems and thus these time-series play the role of noise-free least-squares problems in [20] — each

consecutive time-series observation is defined by its underlying deterministic linear system (defined by an unknown matrix, just as in linear regression). This continuous-state problem has a naturally continuous error metric: mean-squared-error. As in [20], ICL here implicitly involves identifying the underlying system from observations of its evolution.

Segments of random length from different time-series are interleaved with "symbolic punctuation labels," tokens that unambiguously demarcate different segments as belonging to different time-series. These discrete symbolic labels and the fact that they can occur repeatedly introduces a dimension of recall similar to multi-query-associative-recall (MQAR) [3]. However, successful recall is not simply a matter of copying a particular surface-level value from the context. Instead, the corresponding task (predicting the next observation in this particular sequence) must be done.

The discrete symbolic nature (in the sense of [79]) of these punctuation labels means that their meanings must be learned in context to be able to complete associative recall — *they cannot be memorized as the associations change for every new instance of the problem.* Similarly, the details of each distinct time-series itself must also be learned in context. We restrict attention to noiseless time-series defined by orthogonal matrices — consequently, once we have seen enough information in the observation sequence segments for a specific time-series, in principle, perfect prediction accuracy is possible.

## 1.1 Contributions

In this work, we combine the discrete spirit of associative recall with continuous Markov time-series prediction [15, 32] to create a new toy problem (interleaved labeled time-series prediction) with natural continuous performance metrics and LLM-style pretraining with next-output prediction.

We find clear evidence of the emergence of associative recall during training in our toy problem. The compound nature of the problem allows us to see how different abilities emerge sequentially. Specifically, the ability to continue predictions for a resumed sequence after a state has been observed develops first, and the ability to identify sequences from just their corresponding symbolic label emerges later.

Given this, we explore two natural hypotheses for mechanisms by which recall is actually performed:

> **H1: Label-based recall.** The model uses in-context learning of the association of symbolic labels to time-series, and then performs inference based on recalling the queried time-series and continuing its evolution.

> **H2: Observation-based Bayesian recall.** The model ignores the symbolic labels. The noise-free nature of the toy problem means that once we see an observation, we can figure out which previously seen time-series it could have come from. Then, we can do Bayesian prediction [83, 41, 50] based on previous observations.

H1 corresponds to (what we think is) the natural human approach to the task: recall the relevant information and use it. However, **we find that H1 and H2 are both false as complete explanations.** Instead, both are true simultaneously! H1 is used for predicting the first token after a particular time-series is being resumed. To get this to be better than guessing, there is no choice other than using the information in the symbolic label. But for the second token and beyond in a resumed sequence, the information in the observation allows a variant of H2 to work.

We observe further that there is a difference between the emergence of the ability to learn to predict the first versus second token after a symbolic label, even though information-theoretically, they both require recalling the in-context-learned nature of that specific time-series. And somewhat surprisingly, the successful use of the symbolic label (which feels conceptually easier for a human) occurs *after* the model has clearly learned to do some approximate version of the more Bayesian H2 for those tokens on which H2 is a viable strategy. (We verify this using out-of-distribution experiments in Section 4.1.)

A further edge-pruning based investigation shows that the circuits for predicting the first and second tokens are completely distinct. This leads to the following conjecture:

> **C3: Transformers use multiple mechanisms for a single multi-token task.** Distinct mechanisms are used to initiate a new episode of an ICL-specified task (i.e., predict the first token), versus continuing that task (i.e., predict the second token).

Given this novel conjecture emerging from our toy problem, we seek to confirm this phenomenon in the natural language setting of LLMs. By modifying one of the classic LLM emergence experiments [80, 79] and using OLMo checkpoints [55], we confirm that this conjecture holds for an NLP task — i.e. even before the emergence of successful initiation of an ICL-specified task, models can successfully continue performing that task.

## 2 Setup

Consider predicting the continuous-state of an *unknown* linear dynamical system. We focus[1] on the orthogonally evolved system family [67], where the system is defined by $U \in \mathbb{R}^{5 \times 5}$, a random orthogonal matrix. Each $U$ is generated by the algorithm presented in [45], which ensures a uniform sampling over all $\mathbb{R}^{5 \times 5}$ orthogonal matrices. The initial state is $\mathbf{x}_0 \sim \mathcal{N}\left(0, \frac{1}{5}I\right)$, with state updates:

$$\mathbf{x}_{i+1} = U\mathbf{x}_i = U^{i+1}\mathbf{x}_0. \tag{2}$$

The system state is in-principle perfectly predictable, but only after six positions in the sequence are observed by solving for

$$U = \begin{bmatrix} \mathbf{x}_1 & \mathbf{x}_2 & \mathbf{x}_3 & \mathbf{x}_4 & \mathbf{x}_5 \end{bmatrix} \begin{bmatrix} \mathbf{x}_0 & \mathbf{x}_1 & \mathbf{x}_2 & \mathbf{x}_3 & \mathbf{x}_4 \end{bmatrix}^{-1}. \tag{3}$$

### 2.1 Optimal Pseudoinverse Predictor

Following from (3), given the state observations $\{\mathbf{x}_0, \ldots, \mathbf{x}_i\}$, an optimal predictor for this problem computes $\widehat{\mathbf{x}}_{i+1} = \widehat{U}\mathbf{x}_i$, where

$$\widehat{U} = \begin{bmatrix} \mathbf{x}_1 & \ldots & \mathbf{x}_i \end{bmatrix} \begin{bmatrix} \mathbf{x}_0 & \ldots & \mathbf{x}_{i-1} \end{bmatrix}^{\dagger}, \tag{4}$$

and $X^{\dagger}$ denotes the Moore-Penrose pseudoinverse of $X$. Essentially, this baseline only makes non-zero errors on the first, second, third, fourth, fifth, and sixth entry in any sequence — it gets everything else perfectly correct.

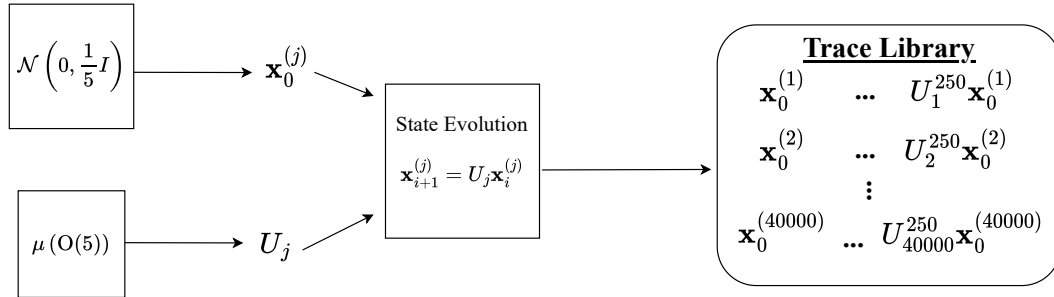### 2.2 Data generation and training



Figure 1: The generation of a train or test library of sequences.

---

[1]In the Appendix, we provide results for the identity system family which has dynamics that are even simpler than the orthogonal system family. For the identity systems, the initial state is $\mathbf{x}_0 \sim \mathcal{N}\left(0, \frac{1}{5}I\right) \in \mathbb{R}^5$. Now, the state updates as

$$\mathbf{x}_{i+1} = \mathbf{x}_i. \tag{1}$$

This trivial process of copying a constant is perfectly predictable after one realization is observed.

**Generating a library of training sequences:** Depicted visually in Fig. 1, we first compile a training library by the following method:

1. Generate 40000 orthogonal matrices iid uniformly over all orthogonal matrices $U_1, \ldots, U_{40000} \overset{iid}{\sim} \mu\left(\mathrm{O}(5)\right)$, where $\mu\left(\mathrm{O}(5)\right)$ is the Haar measure over orthogonal matrices in $\mathbb{R}^{5 \times 5}$. [45]

2. Generate 40000 iid initial states that will correspond to each training system $\mathbf{x}_0^{(1)}, \ldots, \mathbf{x}_0^{(40000)} \overset{iid}{\sim} \mathcal{N}\left(0, \frac{1}{5}I\right)$.

3. Roll out the states to get observation sequences that are each 251 entries long, and compile the sequences as our training library.

**Cutting and interleaving training sequences** To form a training trace, we interleave segments of observation sequences from the library into a context window of length 251, by this process:

1. Insert the start symbol at index 0.

2. Sample the maximum number of systems in the trace $N$ from a $\mathrm{Zipf}(1.5, 25)$ distribution depicted graphically[2] in Fig. 3a. This means that no more than 25 systems will ever appear in a training trace.

3. Choose $N$ of the 40,000 systems in the training library uniformly at random without replacement.

4. Randomly assign to each of the $N$ systems a pair of symbolic open and close labels for this training example.

5. Sample the number of cuts $C \sim \mathrm{Poisson}(2N)$ to be made in the trace. This means that there will be $C + 1$ trace segments in the trace.[3]

6. Place the $C$ cuts uniformly at random with replacement within the context window.

7. For each segment created by the cuts, in order, uniformly at random choose one of the $N$ systems with replacement.

8. At the cut at the beginning of the segment, the open label for this segment's system is inserted.[4]

9. For the system chosen, check if it has appeared in a previous segment of the trace. If not, insert the system's segment from the training trace library starting at index 0. If this system has appeared in a previous segment, insert the system's segment from the training trace library starting at the index that corresponds to the continuation of the previous segment for this system.

10. At one index before the next cut, insert the close label for this segment's system.

$$\texttt{<start>} \left[ \mathbf{x}_0^{(30)} \ U_{30}\mathbf{x}_0^{(30)} \ \cdots \ U_{30}^{24}\mathbf{x}_0^{(30)} \right] \left\{ \mathbf{x}_0^{(2)} \ U_2\mathbf{x}_0^{(2)} \ \cdots \ U_2^4\mathbf{x}_0^{(2)} \right\} \left\{ U_2^5\mathbf{x}_0^{(2)} \ \cdots \ U_2^{26}\mathbf{x}_0^{(2)} \right\} \left( \mathbf{x}_0^{(771)} \cdots \ U_{771}^{45}\mathbf{x}_0^{(771)} \right) \left[ U_{30}^{25}\mathbf{x}_0^{(30)} \ \cdots \ U_{30}^{166}\mathbf{x}_0^{(30)} \right]$$

Figure 2: Example of a 251-element-long interleaved training example.

Note that within a single training example, segments of a particular system always start with the same open token and always end with its corresponding close token. These random assignments

---

[2]The Zipf distribution was chosen for its ubiquity in nature [4] and natural language [69], along with recent work pointing to its importance in modern neural networks [46].

[3]On average, each training trace has $2N + 1$ segments to ensure that the trained model has seen ample interruptions and continuations of systems.

[4]Since the open and close labels occupy two indices in the context window, there are three special cases that can occur: (1) If two cuts are sampled to be on top of each other, then the first of the two cuts that were sampled is ignored; (2) If the two cuts are sampled to occupy adjacent indices, then only the close label for the system corresponding to first of the two cuts is inserted, effectively making that index meaningless as close labels are masked; (3) If the two cuts are sampled so that there is only one index between them, then the open label for the system corresponding to the first of the two cuts is inserted and is immediately followed by the close label for that system, effectively making both indices meaningless due to the masking of the labels. Note that the distributions shown in Fig. 3 do not account for these rare special cases.

are redrawn at the beginning of the interleaving process for each training example; therefore, *the same system can have different symbolic open and close labels when it appears in different training examples*. See Fig. 2 for a diagram of an interleaved training example.

Given this randomized procedure for generating interleaved training examples, we can analyze the training distribution to better understand how frequently a model must recall a system, or sees many systems in a trace.
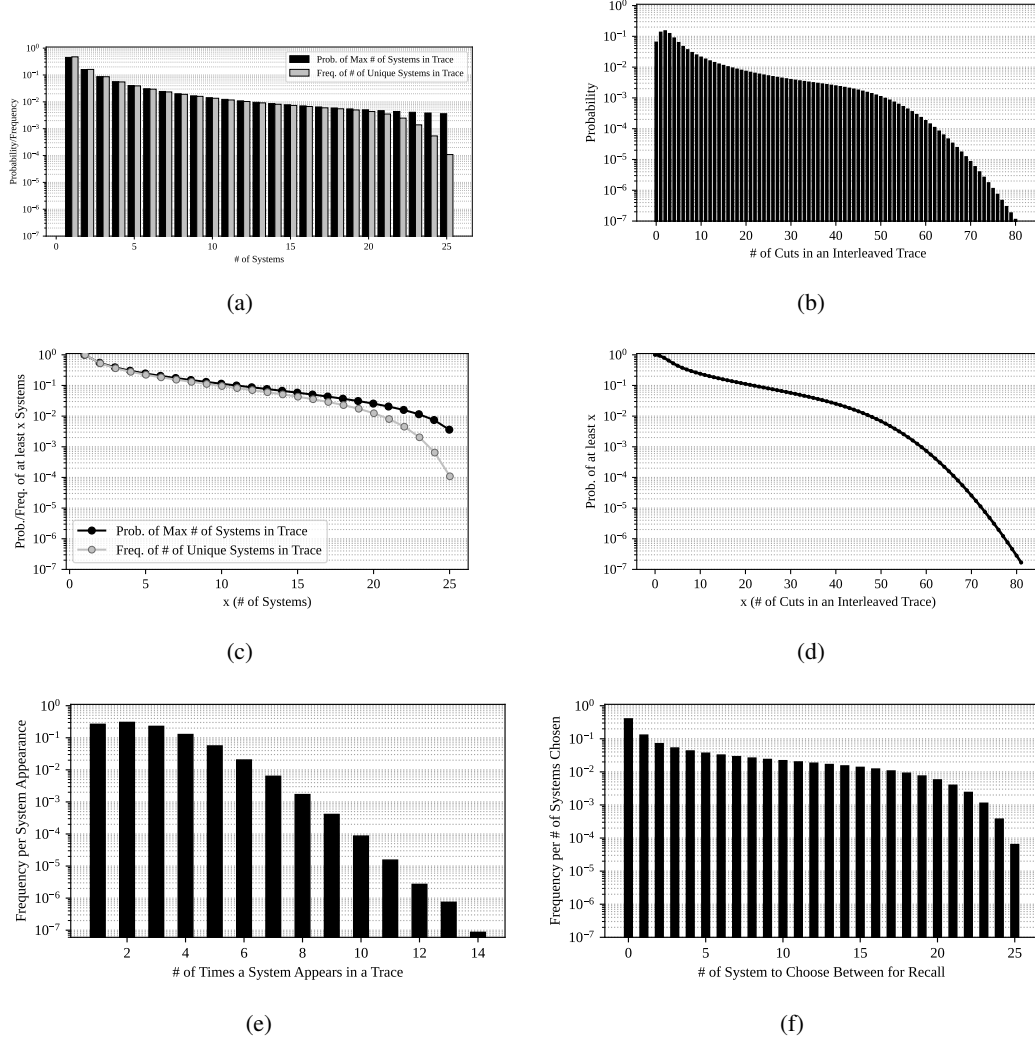


(a)

(b)

(c)

(d)

(e)

(f)

Figure 3: Distributions and complementary cumulative distribution functions (CCDFs) used in data generation — Fig. 3a is the $\mathrm{Zipf}(1.5, 25)$ distribution for the maximum number of systems per trace in black and the number of unique systems per trace for $1 \times 10^7$ traces in silver. Fig. 3c shows the CCDFs for these distributions. Fig. 3b shows the PMF and Fig. 3d shows the CCDF of the number of cuts per trace. Fig. 3e shows the frequency of the number of times a system will appear in the same trace per system appearance. Lastly, Fig. 3f shows the frequency of the number of previously seen systems a predictor must choose between to recall in a trace per system appearance. For example, if system 0 is chosen then system 1, then system 0, then the model must choose between two systems to recall.

In Fig. 3, we show relevant distributions that are derived from the randomized interleaving procedure in this section. Figs. 3a and 3c show the $\mathrm{Zipf}(1.5, 25)$ distribution for the maximum number of systems per trace in black and the number of unique systems per trace in silver. The frequency of number of unique systems per trace follows closely the $\mathrm{Zipf}(1.5, 25)$ distribution for the smaller quantities of systems, but diminishes quicker for larger numbers of systems, due to the coupon-

collecting phenomenon of picking the same system multiple times in a trace. The PMF for the number of cuts made in an interleaved trace is shown in Fig. 3b, while the CCDF for this quantity is given in Fig. 3d. Fig. 3e shows the frequency of how many times a system appears in a training trace. If the same system appears more than once, than the model must perform recall. Therefore, Fig. 3e gives an idea of how often the model must recall a system during training. Finally, Fig. 3f provides the frequency of the number of previously seen systems in the training trace that are candidates to be continued when a predictor is tasked to recall a system. The value zero on the x-axis of this figure means that the model is seeing a system for the first time in a training example and has no need to recall. Later, in Section 3.1.2, we construct tests for the associative recall ability of the trained model for different numbers of candidate systems to be continued in the trace. Fig. 3f shows that for 19 candidate systems, the largest number of candidate systems that we tested on, the model has been presented with this situation less than $1\%$ of the time during training.

**Input structure and embedding**    The input dimension of our models is 57. There are 50 dimensions for encoding paired symbolic open and close labels, a dimension for the start symbol, a dimension for the payload flag and 5 dimensions to hold the 5-dimensional observation vectors. The special symbols are one-hot encoded vectors; see Fig. 4 for an example of the open symbol. For the observation sequence between the SPLs, the 5-dimensional state vectors are inserted into the payload portion of the input vector, the payload flag is set to 1, and the rest of the input vector dimensions are zeroed out.



Figure 4: The one-hot encoding of an open symbolic label. In this example, the system corresponding to this label is assigned to be "system 0."

The entire randomized procedure of generating interleaved training traces is depicted in Fig. 5 along with the structure of the inputs into the model.
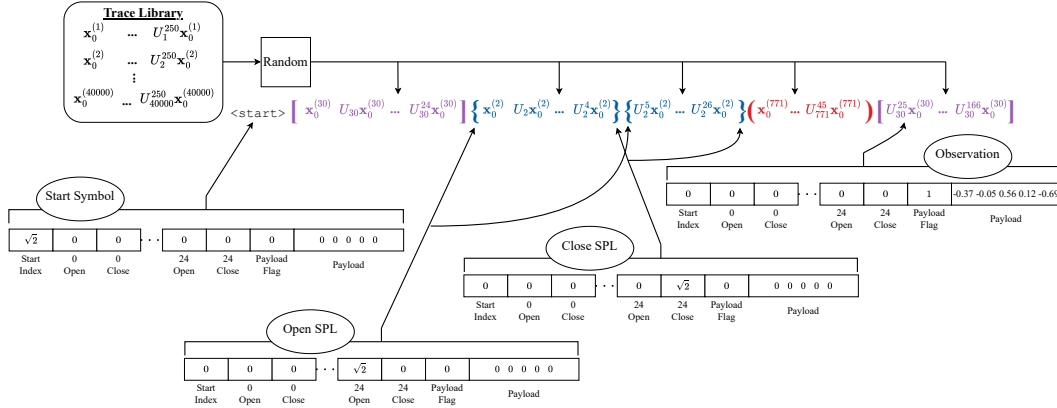


Figure 5: Generating a training example — Notice in this example the continuation from the first segment to the last (system $U_{30}$), and from the 2nd segment to the 3rd (system $U_2$). The "parentheses" (symbolic punctuating labels) are encoded as special tokens as shown.

**Model and embedding**    Building off of the codebase in [15], which was influenced by [20], we train a 2.42M parameter GPT-2 style transformer to perform this task. Our model[5] has hidden dimension 128, 12 layers, and 8 heads. Our model's input embedding is $128 \times 57$ dimensional. The model's output layer is $5 \times 128$ to ensure that the model makes 5-dimensional predictions. The input and output layers are untied [15].

---

[5]These parameter counts and model dimensions are for our "medium" model. Three other models of different sizes were also trained, and their model dimensions are given table 2 in Appendix E.

**Training and Hyperparameters**   New interleaved training examples are generated for each training iteration and our GPT-2-style model was trained for next token prediction on these training traces.[6] The loss for all SPLs on the output were zeroed out. A model that successfully recalls the state of a system seen previously in its context should make predictions after the open token that perform as it would have if the relevant sequence had continued on without interruption.

Following the choice made in [15], we trained our model with a weight decay of $1 \times 10^{-2}$. We used a batch size of 512, a learning rate of $\approx 1.58 \times 10^{-5}$, and trained on a single NVIDIA L40S GPU with 45GB of RAM. A single training run takes around 5 days. We used the AdamW Optimizer [37] and trained using mean-squared error loss.

## 3   Results

This toy problem has two qualitatively different capabilities that a model can acquire during training: in-context learning a linear system as it is seen (ICL), and recalling what has been in-context learned about a previously seen system (associative recall). Within the ICL ability, again, there are qualitatively two sub-abilities: ICL for the first system that is seen, and ICL for the subsequent new systems that are seen. The second sub-ability requires a model to learn to restart its in-context learning for a new system.

### 3.1   Test Setup

#### 3.1.1   Uninterleaved sequence test

To test the model's ICL ability for the first system that is seen (Section 3.2), we generate 100 held-out systems and 1000 different held-out initial states[7] by the same method described in Section 2.2 to form our testing library. We then evaluate the model on the uninterleaved traces from this testing library.

#### 3.1.2   Needle-in-a-haystack test

To evaluate the model's ability to restart ICL on a new system (Section 3.2.2) and recall a previously seen system (Section 3.3), we generate a series of structured "needle-in-a-haystack" test traces through interleaving the traces in the testing library generated in Section 3.1.1. A single "needle-in-a-haystack" trace is generated by the following procedure:

1. Choose $N \in [1, 19]$ to be the number of distinct systems in a test trace.
2. For each of the $N$ systems, insert a segment of 10 state observations starting from index 0 from the testing library into the "needle-in-a-haystack" test trace. Each of these segments are individually punctuated with a unique open and close symbol pair. We call this portion of the trace the "haystack".
3. Append a query open symbol to the test trace that signifies which system in the haystack will be be continued. The segment that will be continued is called the "needle".
4. Append 10 state observations from the continuation of the system in the haystack corresponding to the query open symbol. This portion is called the "test segment".

See Fig. 6 for a diagram of a test trace for $N = 2$ systems in the haystack and system $U_1$ as the needle.

For the full "needle-in-a-haystack" test dataset, we would like to ensure that we test on the same systems for the different values of $N$, while having diverse systems in our dataset so that our results are statistically meaningful. To achieve this, we test on 50 "needle-in-a-haystack" trace configurations. A trace configuration is a specific ordering of systems from the testing library in the positions of the haystack. For the first needle-in-a-haystack trace configuration that we generate, we place a segment

---

[6]The model sees newly interleaved training examples at each iteration, but the training traces that are interleaved into the training example are drawn from the fixed training library of 40,000 sequences. Therefore, the model undergoes single-epoch training where the information within a training example might have appeared in many other training examples.

[7]We generate 1000 initial states for each system to narrow down the quartiles in the squared-error curves.
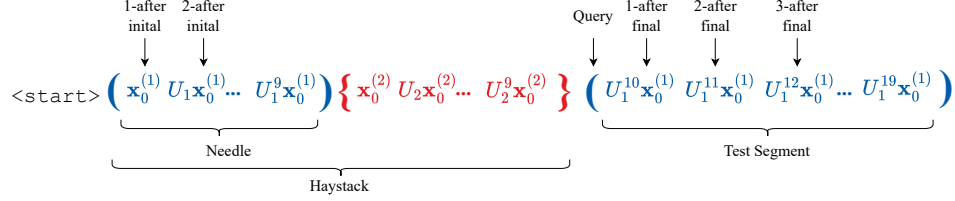
Figure 6: Needle-in-a-haystack test example. (A two system haystack.)

from "system 0" from the testing library into the first position in the haystack, and fill the rest of the haystack positions consecutively until the $N$-th position is filled with a segment from "system $N-1$". For the next trace configuration, the first position in the haystack is filled with a segment from "system 1" and the rest of the haystack positions are filled consecutively until the $N$-th position is filled with a segment from "system $N$". This pattern continues until the last trace configuration. In our case, we tested on 50 trace configurations, meaning the haystack of the last trace configuration started with "system 49" and ended with "system $48 + N$". Each trace configuration is populated with 1000 different initial states for each system. For the results in the main paper, the test segment is a continuation of the segment in the first position of the haystack. For results where segments in other haystack positions are continued in the test segment see Section H in the Appendix.

$$
\begin{matrix}
\{0 & \ldots & N-1\} \\
\{1 & \ldots & N\} \\
\vdots & \vdots & \vdots \\
\{49 & \ldots & 48+N\}
\end{matrix}
$$

Figure 7: When testing on 50 needle-in-a-haystack trace configurations, the order of system indices from the testing library in a haystack of size $N$ for each needle-in-a-haystack test trace is given above. For each system, 1000 sequences are interleaved to build a testing dataset of shape $50 \times 1000 \times (12N+1) \times 5$. The shape of the second to last axis is due to the start token and haystack segments being 10 context indices long plus two indices for the symbolic open and close labels. The last axis is 5-dimensional since every system has 5-dimensional observations.

## 3.2 In-context learning system dynamics

We now present results on how well the transformer-based model can in-context learn a linear system's dynamics. Section 3.2.1 shows the results for learning about the first system that is seen, while Section 3.2.2 shows the results for learning about subsequent systems that are seen and the ability to restart ICL.

### 3.2.1 Can a model learn the first system in-context? Yes

We first confirm that our trained model is able to in-context learn the first system seen in its context. We evaluate the model on the uninterleaved traces from the testing library specified in Section 3.1.1. In Fig. 8a, we plot the median squared-error over these test traces vs the context index, and the color of each curve represents how far along the model is in training. The dotted line in this figure is the median squared-error of the pseudoinverse predictor in (4) over the same test traces. In Fig. 8b we plot the median squared-error over these test traces as training proceeds (measured by the number of training examples seen so far), the color of each solid curve represents the context index, the blue and green dotted horizontal lines are the optimal pseudoinverse predictor's median squared error for the specific early context indices.

Notice in Fig. 8a that the early model checkpoints saturate out and cannot continue to make better predictions with more context. Nevertheless, after seeing $6.25 \times 10^7$ training examples, the model's prediction error on indices 2 through 7 closely match those of the pseudoinverse baseline from (4). This is also seen in Fig. 8b where, as training proceeds, the dark curves representing the model's prediction errors on early indices converge to the prediction error of the pseudoinverse predictor for the corresponding index given by the blue and green curves. The best performance of the model is at
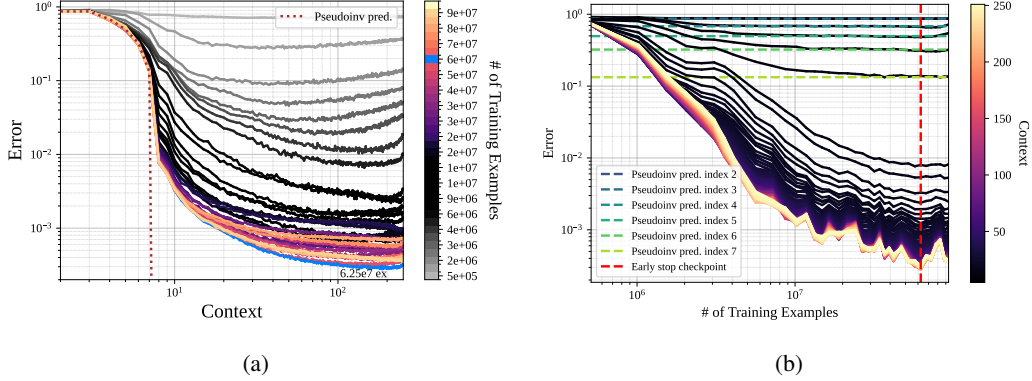
Figure 8: Performance on a long uninterleaved trace — 8a and 8b depict the in-context learning performance of the model on long uninterleaved test examples. The median squared-error is plotted against the context index in Fig. 8a, and against the number of training examples seen in Fig. 8b. The red line in Fig. 8b at $6.25 \times 10^7$ training examples, and the blue curve in Fig. 8a denote the checkpoint that we use for early stopping. in Fig. 8a the optimal pseudoinverse predictor is the brown dashed curve, while in Fig. 8b it is denoted by the blue and green horizontal lines for each of the early context indices. in Fig. 8b, notice that the model gradually learns to make optimal predictions as opposed to the sudden emergence of assoiative recall ability seen in Section 3.3, since the prediction errors for early indices slowly converge towards the pseudoinverse predictor's performance.

the end of the context window with a median squared-error of $\approx 2 \times 10^{-4}$. According to [34], this is near the practical precision threshold for transformer models.

Notice in Fig. 8b the model is gradually learning to make better predictions as training continues. This ICL ability for the first system seen is the same ability that is studied by [15, 67] and using this evaluation metric, we see that sudden emergence is not present. Nonetheless, using the same evaluation metric, the ability to restart ICL for a new system (Section 3.2.2) and to recall a previously seen system (Section 3.3) both exhibit emergence.

Additionally, we notice in Fig. 8b that the squared-error bottoms out after $6.25 \times 10^7$ training examples, showing that the model suffers from overfitting late in training. Having seen this, we set our early stopping checkpoint at $6.25 \times 10^7$ training examples, as denoted by the red vertical line in Fig. 8b which corresponds to the blue curve in Fig. 8a.

In summary, the model is able to use context to make better predictions of state observations from held-out test systems. The model's ability gradually develops during training, as opposed to how associative recall develops suddenly later in training as seen in Section 3.3.

### 3.2.2 Emergence of the ability to restart predictions on new systems

We now show the training dynamics for the ability to restart in-context learning for a new system that was not the first system seen in-context. We find that learning this restart ability is not gradual, as it was for learning the first system seen (Section 3.2.1). Instead, we see that the model begins to transition from poor restart performance to good restart performance early in training as compared to when associative recall emerges in training which we will see in Section 3.3. We study the model's performance on the haystack segments of "needle-in-a-haystack" test examples (see the diagram in Fig. 6).

In Fig. 9, the median squared-error vs the number of training examples seen is plotted for steps 1 through 8 into the first system segment in the haystack and the third system segment in the haystack. Specifically, in Fig. 9a we see that at the beginning of training the model has not learned to restart its predictions for the third system segment, as its median squared-error in segment 3 is well above its counterpart predictions in segment 1 for all steps into each segment except for step 1. This shows that early on in training, there is clearly substantial interference from earlier segments when the model tries to learn to predict the behavior for a new sequence (explicitly labeled as such) that it is seeing in the third position in the haystack. As training proceeds, we see that the median squared-error for

each step in segment 3 converges to the value of its segment 1 counterpart. Fig. 9b shows that the model transitions towards restarting ICL correctly when training has processed $\approx 2 \times 10^6$ training examples.
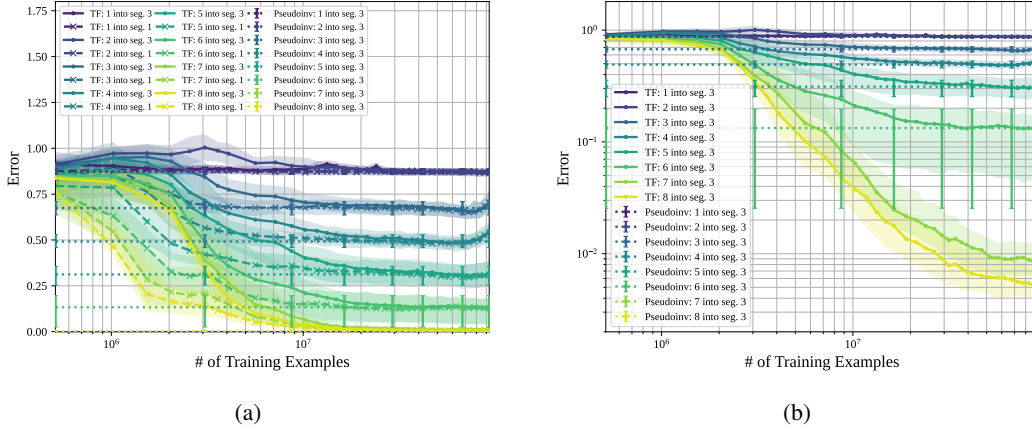


(a)                                                                                              (b)

Figure 9: Performance on new subsequent segments. 9a is the squared-error of predictions on steps 1 through 8 into the first and third system segments, where each segment is seen for the first time in context. 9b is the squared-error for steps 1 through 8 into the third system segments on log scale. The error bars show the 25th and 75th percentiles across trace configurations of the model's prediction error across the medians over the 1000 initial states in each trace configuration. The horizontal dotted lines are the median squared-error of the optimal pseudoinverse predictor.



Figure 10: Restarting for a new system at the early-stopping checkpoint after seeing $6.25 \times 10^7$ training examples.

In Fig. 10, we show the median squared-error of the model's predictions for up to 8 steps into each new segment at the early-stopping checkpoint of $6.25 \times 10^7$ training examples. This log-scale plot shows that even at the end of training, the model's ability to restart ICL for new systems slowly degrades when predicting indices 6, 7 and 8 as more new systems are presented in the context window. This is seen as the upward trend in the green and yellow curves in Fig. 10.

### 3.3 Emergence of associative recall

We now study the prediction performance of the model when queried for recall. In particular, how it depends on the index into the test segment. Furthermore, we study how the ability to predict different indices in the test segment develops differently during training. Particularly, the ability to predict the first index into the test segment develops much later in training and much more abruptly than the ability to predict the subsequent indices into the test segment.

10

To uncover this, we test the model on the "needle-in-a-haystack" test traces described in Section 3.2 with different values of $N$, for the number of systems in the haystack. Fig. 11 shows results for $N = 1, 2$, and 5. Results for more $N$ values are given in Appendix F. In Fig. 11 the solid curves marked with dots show the median squared-error of the model's predictions on the first, second, third, seventh, and eighth indices into the test segment. These curves are labelled as "after final" in the legend, because their indices are after the final open symbol. To contrast, the dashed curves marked with crosses show the median squared-error of the model's predictions on the same indices into the first segment in the haystack. These curves are labelled as "after initial" in the legend, since their indices occur directly after the initial open symbol. These indices are depicted in Fig. 6.



| (a) 1 system haystack. | (b) 2 system haystack. | (c) 5 system haystack. |

Figure 11: Training dynamics for recall — The $25^{th}$, $50^{th}$, and $75^{th}$ quantiles of the squared-error of the model's predictions vs the number of training examples seen during training so far are plotted on log-log plots for $N = 1$ in Fig. 11a, $N = 2$ in Fig. 11b, and $N = 5$ in Fig. 11c. All haystack segments are of length 10 (excluding delimiting tokens). The test set consisted of 1,000 "needle-in-a-haystack" traces from each of 50 systems. The dashed curves marked with crosses show the performance for indices 1, 2, 3, 7, and 8 steps after the initial open symbol, while the solid curves marked with dots show the performance for the same indices after the final open symbol. Notice in all the above figures that the solid black curve, showing the model's ability to recall the correct system from just seeing its corresponding open symbol, very sharply improves after $10^7$ training examples.

To perform perfectly on the associative recall task, the model must implicitly learn and remember the correct $5 \times 5$ orthogonal matrix corresponding to a particular system. In all subfigures in Fig. 11, the solid curves largely decrease as training proceeds.[8] This shows that training is improving the model's ability to recall. More specifically, the solid black curve for the pure recall task of predicting the first observation in the test segment shows emergence-style behavior: first a steady high squared-loss until it begins to drop at some point in training. Interestingly, this ability begins to emerge *earlier* ($\approx 1.5 \times 10^7$ training examples seen) for the simple case of $N = 1$ (Fig. 11a), than for case when $N = 5$ where the transition happens closer to $2.5 \times 10^7$ (Fig. 11c). One system in the haystack is an arguably-trivial recall test, since the model must recall the only system that it has seen so far. While recall performance for $N = 5$ is showing seemingly no improvement in Fig. 11c, the recall for $N = 1$ in Fig. 11a has gotten substantially better. The whole time, we are also seeing improvements in recognizing the underlying state evolution of the time-series as evidenced by better progress in the other blue, red, yellow, and green solid curves — all of which involve predicting the correct orthogonal transformation of the observed state.

In Fig. 12, we plot the quartiles of the squared-error of the model's predictions as a function of the number of systems in the haystack. Observe that the 1-after final open label performance remains steady with more systems, while the predictions for the later indices get progressively worse. Information theoretically, making a prediction for 2-after the final open label is just as easy as making a prediction for 1-after the final open label, if not easier, since the open label provides the information required to make an optimal prediction. In light of this, Fig. 12 suggests that the model has not learned to use symbolic labels well enough to maintain steady performance on the indices two or more after the final open symbol as the number of systems in the haystack increases.

---

[8]For the solid blue and red curves showing the model's recall ability on the second and third indices after the query symbol, we notice double-descent behavior that is more pronounced for larger haystacks.
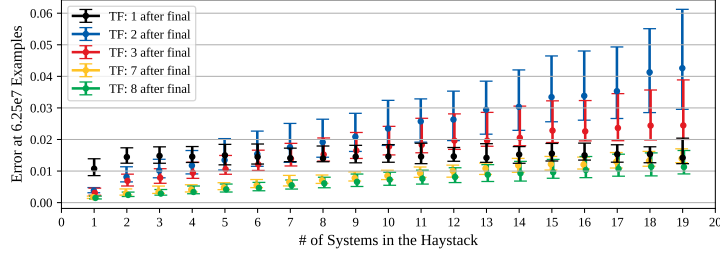
Figure 12: The $25^{th}$, $50^{th}$, and $75^{th}$ quartiles of the squared-error after $6.25 \times 10^7$ training examples as the number of systems in the haystack $N$ increases. Notice that predicting 1-after the open symbol is largely unaffected by the value of $N$, as the black markers stay steady around $0.015$. Although the information encoded in the open symbolic label keeps the difficulty of the task constant as $N$ increases, nevertheless, the performance for 2,3,7, and 8-after the open symbol predictions steadily degrades as $N$ increases.

# 4    Multiple Mechanisms for Prediction

As discussed in Section 1.1, we explore whether the trained model performs the associative recall task through a **label-based recall (H1)** mechanism or a **observation-based Bayesian recall (H2)** mechanism. If the mechanism for recall was label-based, the model would in-context learn the association of symbolic labels to their corresponding systems, and then perform inference based on recalling the queried system and continuing its evolution. If the mechanism for recall was observation-based and Bayesian, the model would ignore the symbolic labels and instead would use the state observation to figure out which system the observation could have come from. The model then performs Bayesian prediction based on previous observations to make future predictions.

## 4.1    Out-of-Distribution Inference-Time Experiments

To decide if H1 or H2 is a valid hypothesis for the associative recall mechanism, we conducted three out-of-distribution experiments at inference-time: misdirecting the model towards the incorrect sequence in the haystack (Section 4.1.1), synchronizing sequences in the haystack from different systems so that they would all have the same state after the final open symbol (Section 4.1.2), and misdirecting the model towards an unseen sequence not present in the haystack (Section 4.1.3). Through these three out-of-distribution experiments, we found that neither H1 nor H2 fully describe the model's mechanism for associative recall. In fact, the evidence indicates that model uses H1 for predicting 1-after the final open symbol, and a suboptimal version of H2 for predicting two or more steps after the final open symbol.

To further explore the mechanism for restarting ICL on a new system, we conduct a fourth out-of-distribution experiment: misdirecting the model towards a seen sequence in the haystack (Section 4.1.4). This experiment's results provide evidence that the model ignores the open symbol when restarting ICL on a new system.

### 4.1.1    Experiment 1: Misdirection towards the incorrect sequence in the haystack

For this out-of-distribution experiment, we test the model on "needle-in-a-haystack" test traces generated as specified in Section 3.1.2, except we swap the final open symbol with another open symbol that corresponds to a segment in the haystack that is not the "needle" as seen in Fig. 13. If H1 were true, the model would be using label-based recall and we would expect it to make predictions on the test segment for the wrong system. If H2 were true, the swapping of the label would not affect the prediction performance of the model, since the model would be using the seen states to make its predictions rather than the symbolic labels.
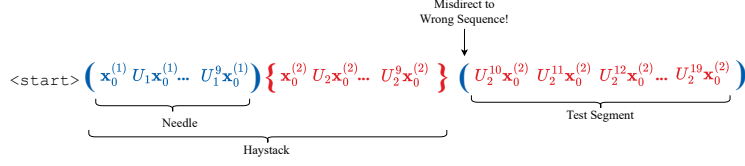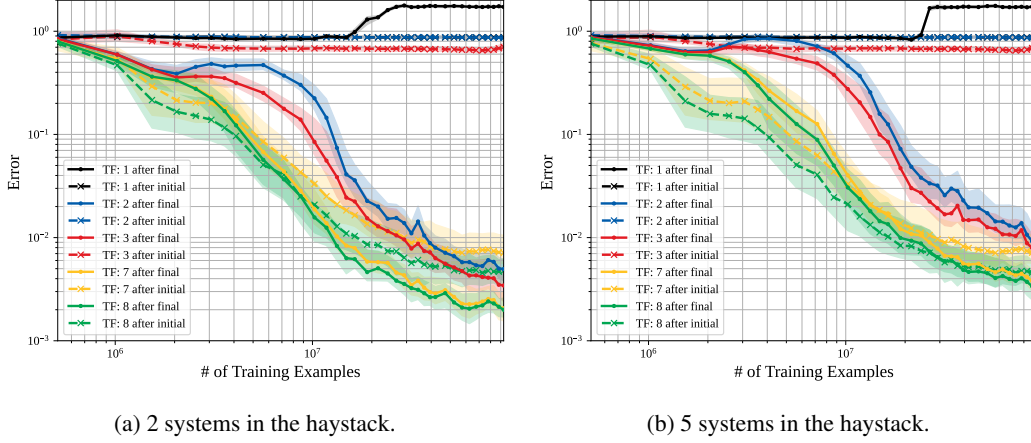
12

$$\texttt{<start>}\ \left(\ \mathbf{x}_0^{(1)}\ U_1\mathbf{x}_0^{(1)}\cdots\ U_1^9\mathbf{x}_0^{(1)}\ \right)\left\{\ \mathbf{x}_0^{(2)}\ U_2\mathbf{x}_0^{(2)}\cdots\ U_2^9\mathbf{x}_0^{(2)}\ \right\}\ \left(\ U_2^{10}\mathbf{x}_0^{(2)}\ U_2^{11}\mathbf{x}_0^{(2)}\ U_2^{12}\mathbf{x}_0^{(2)}\cdots\ U_2^{19}\mathbf{x}_0^{(2)}\ \right)$$

Figure 13: Misdirecting the model towards the incorrect sequence in the haystack.



(a) 2 systems in the haystack.

(b) 5 systems in the haystack.

Figure 14: Misdirection towards incorrect sequence — The median squared-error of the model's predictions on the test segment after observing the incorrect final open label vs the number of training examples seen so far. The solid black curve sharply increases in these figures where it would have sharply decreased for a normal test trace as seen in Fig. 11, suggesting H1 is true for predicting the first index of the test segment. In contrast, we find that the model ignores the misdirection when predicting two or more indices into the test segment, since the solid blue, red, yellow, and green curves are almost identical to the corresponding curve in Fig. 11. This suggests that the model is using a Bayesian approach when predicting these indices.

**Misdirecting the model towards the incorrect sequence in the haystack falsifies pure label-based recall and provides strong evidence that an observation-based Bayesian recall mechanism is present.** The results of this experiment are shown in Fig. 14 for $N = 2$ and $N = 5$, where the median squared-error of the model's predictions are plotted against the number of training examples seen, mirroring the format of the Fig. 11. The first thing to notice in Fig. 14 is that the solid black curve now sharply rises late in training, as opposed to its sharp fall for a normal "needle-in-a-haystack" test trace as shown in Fig. 11. Presumably, the model is using the label-based recall to predict the first index into the test segment.

In contrast, the solid curves in Fig. 14 for 2, 3, 7, and 8 after the final open symbol look almost identical to the corresponding solid curves in Fig. 11. See the blue curve in Fig. 14a and 11b at $2 \times 10^7$ training examples. Their median squared-error both sit at around $2 \times 10^{-2}$. The same correspondence can be seen for the red, yellow, and green curves as well. This shows that the model's predictions for 2, 3, 7, and 8 after the final open symbol are not very affected by the open symbol. Instead, they must be using the state observations after the final open symbol to decide which system to use for making predictions. This strengthens the case that the model uses an observation-based Bayesian recall mechanism for predicting the indices after the first index into the test segment.

### 4.1.2 Experiment 2: Synchronizing "rotations" in the haystack

In Section 4.1.1, the misdirection towards the incorrect sequence experiment showed that the model can make accurate predictions for indices 2 and further into the test segment without using the final open symbol. However, in that experiment the observation that is 1 index after the final open symbol can determine which system should be appplied for predicting the subsequent indices, since if a predictor has observed $\mathbf{x}_9^{(1)}$ from system 1 and $\mathbf{x}_9^{(2)}$ from system 2, when it observes $\mathbf{x}_{10}$ from either

13

system 1 or system 2, it can check whether $\mathbf{x}_{10} = U_1\mathbf{x}_9^{(1)}$ or $\mathbf{x}_{10} = U_2\mathbf{x}_9^{(2)}$. Now, there is still the question of whether the model can use the final open symbol to predict the later indices in a situation where the 1 after final observation does not provide the necessary information. To answer this question, we conduct a synchronizing rotations experiment, where all of the sequences in the haystack from different systems all have the same state at the $10^{\text{th}}$ index, which corresponds to 1 index after the final open symbol. To do this, we first generate a single vector $x_{10} \sim \mathcal{N}\left(0, \frac{1}{5}I\right)$ for all systems and generate the haystack by "rewinding" our systems back to their initial state $x_0$ by $x_{i-1} = U^T x_i$ as is shown in Fig. 15. The ambiguity of which system the first observation in the test segment comes from, means that the model must use the symbolic label to make an accurate prediction on the second index into the test segment.



Figure 15: Synchronizing previous systems in the haystack, so the first observation in the test segment no longer reveals the system that is being continued.



(a) 2 systems in the haystack.

(b) 5 systems in the haystack.

Figure 16: Synchronizing rotations — The median squared-error of the model's predictions on the test segment when the first index into the test segment is a valid continuation for all haystack segments vs how many training examples have been processed during training. Synchronizing the rotations of the haystack segments means that after the first observation in the test segment, a predictor cannot determine which system is being continued. The solid black curve still shows emergence at the same point in training as it did in Fig. 11, supporting the validity of H1 for predicting the first index into the test segment. Otherwise, the model performance is significantly degraded for the rest of the indices into the test segment. For example, the solid blue curve in these plots for prediction errors on the second index into the test segment is near 1.0 in Fig. 16a where it is near $1.0 \times 10^{-2}$ in Fig. 11b. This indicates that the model is unable to use the symbolic label well enough to make accurate predictions.

**Synchronizing the sequences in the haystack so the observation at the first index in the test segment is not informative of the continuing system shows that the model has not learned to use the symbolic label to make accurate predictions for two or more indices into the test segment.** In Fig. 16, which plots the median squared-error of the predictions on the synchronizing test traces vs. the number of training examples seen, we see that the solid black curve still sharply decreases late in training[9], as it does in Fig. 11, showing that the model is able to use the final open symbol to predict $\mathbf{x}_{10}$. On the other hand, the solid blue curves are almost horizontal throughout all of training, and the solid red, yellow, and green curves are have a significantly higher squared-error than their

---

[9]Due to the synchronization, the first observation in the test segment is a valid continuation for all systems in the haystack and therefore is always predictable.

14

counterpart curves in Fig. 11. This means that the model is unable to make accurate predictions on the subsequent indices in the test segment after $\mathbf{x}_{10}$, although the final open symbol provides all of the necessary information to do so. This strengthens the case for H2 being the right hypothesis for predicting the indices after $\mathbf{x}_{10}$, although the model is clearly not a Bayes optimal predictor since seeing $\mathbf{x}_{10}$ and the next observation $U\mathbf{x}_{10}$ would disambiguate which system $U$ is being continued. Despite the system being disambiguated after seeing $U\mathbf{x}_{10}$, the red solid curves in Fig. 16 are still much worse than their counterparts in Fig. 11.

### 4.1.3 Experiment 3: Misdirection towards an unseen sequence

To further study whether the model uses the symbolic labels at all to continue its predictions on a recalled sequence, we misdirect the model to restart ICL for a sequence that has not appeared in the haystack so far. Fig. 17 illustrates how we swap out the final open symbol of a normal "needle-in-a-haystack" test trace with an open symbol that does not correspond to any of the systems in the haystack. Given the results from the previous two sections, we expect the model to predict zero for the first index as that is optimal for restarting a new sequence, but we expect the model to make accurate predictions on the subsequent indices that are continuing a segment that is in the haystack. This is almost what happens, but late in training, the model suddenly transitions to (at least partially) restarting ICL on a new sequence.



Figure 17: Misdirecting the model with an unseen symbolic label indicating a new sequence.



(a) 2 systems in the haystack.

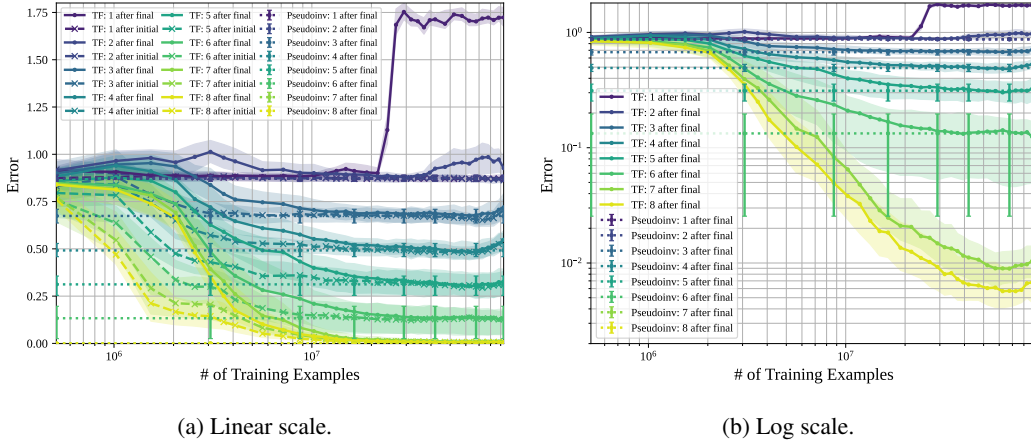(b) 5 systems in the haystack.

Figure 18: Misdirection towards an unseen system — The median squared-error of the model's predictions on the test segment when the final open symbol does not correspond to any haystack systems vs the number of training examples seen so far. The solid blue and red curves match their counterparts in Fig. 11 until $\approx 3 \times 10^7$ training examples, at which point they sharply increase. This means the model continues predicting the existing system in early training then suddenly starts treating it more like an unseen system late in training. We note that this transition happens shortly after the emergence of associative recall, suggesting the model has learned that unseen labels also require ICL to be restarted.

**Misdirecting with an unseen symbolic label highlights a new abrupt phase transition in model behavior late in training from disregarding the symbolic label and continuing to predict the observed test segment to restarting ICL for a new system.** Fig. 18 shows the median squared-error of the model predictions after being presented with a final open symbol that does not correspond to any system in the haystack while the test segment is actually a continuation of a haystack segment vs

15

how far along the model is in training. In the same figure, the first after final predictions match the expected behavior of restarting predictions as seen by the solid black curve being a horizontal line near 1. For the rest of the indices, the model ignores the symbolic label through the initial stages of training and continues to correctly predict the test segment, since the solid curves in Fig. 18 match their counterparts in Fig. 11 up to around $3 \times 10^7$ training examples. Once associative recall fully emerges later in training, and the model learns how to use the symbolic labels, the model transitions to treating a continuation of the old sequence as a brand new sequence corresponding to the new label, since the blue and red solid curves abruptly increase after $3 \times 10^7$ training examples in Fig. 18. This shows that the model predictions for two or more indices into the test segment *are* affected by the final open symbol, although Section 4.1.2 showed that the model is unable to use it to make accurate predictions for recall on those indices.

### 4.1.4   Experiment 4: Misdirection towards a seen sequence in the haystack

The first three out-of-distribution experiments studied the model's mechanisms for performing associative recall, but this section will study the mechanism for restarting its prediction on a new system. Again, we devise a misdirection experiment. This time we provide a final open symbol that points toward the "needle" in the haystack, but the test segment is a sequence from a system that is not in the haystack (Fig. 19).



Figure 19: Misdirecting the model with a previously seen symbolic label.



(a) Linear scale.

(b) Log scale.

Figure 20: Misdirection towards a seen system — The median squared-error of the model's predictions on the third segment that is an unseen sequence while the final open symbol corresponds to the first segment vs the number of training examples seen so far. The solid curves other than the 1-after final curve match the solid curves in Fig. 9. This supports the hypothesis that the model uses the state observations to continue its prediction on a newly seen segment, rather than the open symbolic label.

**Misdirecting with a previously seen symbolic label does not stop the model from restarting its predictions on the new sequence for two or more indices into the test segment.** In Fig. 20, the median squared-error on this misdirection experiment is plotted against the number of training examples seen so far during training. We find that when given the correct unseen symbolic label (Fig. 9a), and when shown a symbolic label misdirecting to a sequence that has already been observed in context (Fig. 20), the model performs equivalently on predicting two or more indices into the newly seen segment, as the solid curves match except for the 1-after final curve that is using the symbolic open label. The model recognizes that the sequence it is seeing does not correspond to a sequence it has already seen. This supports a hypothesis that that model does not use the symbolic labels to

16

restart ICL on a new system. Interestingly, misdirecting the model with a previously unseen symbolic label indicating a new system (Section 4.1.3) shows that these unseen symbolic labels do affect the model's predictions in the test segment. This evidence shows that the mechanism for restarting ICL may not be purely label-based nor purely observation-based.

### 4.1.5 Summary of out-of-distribution experiment results

After conducting the four out-of-distribution experiments, we find that neither H1 nor H2 can fully explain the model's mechanism for predicting interleaved time-series. The misdirection to the incorrect sequence experiment (Section 4.1.1) supports H1 for predicting the first index into the test segment, but it also shows that the model can make accurately continue its predictions in the test segment without using the final open symbolic label. Therefore, H1 is insufficient. The synchronizing rotations in the haystack experiment (Section 4.1.2) further showed that the model is unable to use the final open symbolic label well enough to accurately continue its predictions in the test segment, providing strong evidence that the model performs a suboptimal version of H2 for this subtask. This evidence is further strengthed by the misdirection towards a seen sequence in Section 4.1.4. Finally, the misdirection to an unseen sequence experiment (Section 4.1.3) showed that even for the later indices, the final open symbolic label can signal the model to restart its predictions for a new segment, invalidating H2 as the sole mechanism for continuing predictions. It is clear from these results that the conjecture C3, where the transformer model uses multiple mechanisms for the single task of interleaved time-series prediction holds true.

### 4.2 Transformer circuit analysis

| Circuit | # Edges | 1-After squared-error | 2-After squared-error |
|---|---|---|---|
| Orthogonal Sys Full Model | 32936 | 0.002 | 0.004 |
| Orthogonal Sys 1-after Circuit | 200 | 0.01 | 0.64 |
| Orthogonal Sys 2-after Circuit | 40 | 0.32 | 0.02 |

Table 1: Edge pruning finds sparse transformer circuits with high evaluation accuracy in our GPT2-style model. We prune late checkpoints of a model using interleaved traces and data consisting of 5 systems in the haystack. We report the number of edges in the final circuit and the mean squared-error of the circuits' predictions and the ground truth state observations for both the 1-after and 2-after tasks. We run an edge thresholding binary search to reach the target edge sparsity and set all pruned edges to have weights of 0, then run inference. We visualize the 1-after circuit in Fig. 24.

Now that it is clear the model uses multiple mechanisms for the single task of interleaved time-series prediction, we study if these different mechanisms have different computation graphs in the weights of the learned transformer model. To do this we use Edge Pruning, a transformer circuit-discovery method that optimizes over continuous masks over a disentangled transformer to find a sparse representation of a task [6]. We run a modified version of Edge Pruning to distinguish the circuits being used by our model for the 1-after and 2-after tasks. As our model is solely trained with squared-error, we remove the KL objective in the original Edge Pruning method's loss function in [6] and optimize on a scaled up squared-error added to the original edge loss. The loss function that we optimize to find a sparse computation graph is $\mathcal{L}' = k \cdot \mathcal{L}_{\text{squared}-\text{error}} + \mathcal{L}_{\text{edge},s}$. Our dataset for the edge pruning method consisted of one "needle-in-a-haystack" trace configuration generated in the same way as in Section 3.1.2 for 5 systems in the haystack.[10]

We report the size of the circuits and the squared-error of the ground truth with the predictions outputted from both the final checkpoint of the trained model and the pruned circuits in Table 1.[11] Since the pruning method optimizes continuous gates for each node in the computation graph, these continuous gates must be quantized to 0 or 1 to get a pruned circuit.[6] Our results in Table 1 show the accuracy of the pruned model *after* this quantizing has been done, which differs from how [6]

---

[10]Further experiments should test more "needle-in-a-haystack" configurations, and would have a larger testing library of traces to create full train and test splits for training and evaluating the pruning method.

[11]The trained model that was used for these results is from an earlier training run than the "orthogonal medium" model that is throughout the rest of this paper. This earlier training run was trained on $5 \times 5$ orthogonal matrices that were sampled from a non-uniform distribution over all $5 \times 5$ orthogonal matrices [45].

reports their results. Post-quantization performance ensures that edges that we believe to be irrelevant truly have no contribution to the output. Importantly, we find high accuracy and *0% edge overlap* between the 1-after and 2-after circuits, indicating that *our model mostly leverages mechanistically different learnt mechanisms for consecutive tokens*. See Appendix B for further details on the pruned circuits.

## 5 Do Pretrained LLMs Also Display Multi-Mechanism Tendencies? Yes!

To see whether our conjecture C3 (multiple mechanisms for a single multi-token task) holds for natural language problems solvable by prompting LLMs, we leverage OLMo-2 7B checkpoints [55] and a basic English to Spanish translation task that is inspired[12] by the IPA translation task used in previous works benchmarking and studying emergent behaviors [78, 5]. In Fig. 21 (right), we see a similar phase transition in the first token prediction task, a parallel of the 1-after dynamics of the associative recall setup (Section 3.3). Meanwhile, the second-token performance is both better and more gradual in its improvement across training. This matches what we saw in our toy problem in Section 3.3.

One natural question is whether what we are seeing in Fig. 21 is the emergence of true ICL recall or just the underlying ability to start a translation itself. This can be probed by replacing the purely symbolic task labels "X:" and "Y:" in the few-shot examples with semantically informative "Spanish:" and "English:" labels. This replacement switches the problem from pure ICL for task recognition (*in-context associative recall*) to leveraging a learnt label (*in-weights associative recall*). The resulting performance is seen in Fig. 21 (left). Notice the marked improvement in the first-token performance that erases the entire gap to the second-token performance. This establishes that the model knows how to start a translation, it just can't in-context-learn well enough to know that is what it is supposed to do before the phase transition around 50k training steps.



Figure 21: Comparative example of in-weights associative recall (left) and in-context associative recall (right) in a 2-shot prompting setting. Each point is a separate OLMo-2 7B model at different training steps. We report 95% credible intervals using Jeffreys prior ($Beta(0.5, 0.5)$) based on 100 samples per evaluation point.

## 6 Pretraining Loss

In [16], it was shown that many emergent abilities emerge for models of different sizes once a model's pretraining loss performance on a broad corpus of held-out data in the style of their pretraining data reaches a certain threshold. We want to see if this holds true in our toy setting of interleaved time-series prediction.

We evaluate four different model sizes (see Table 2 in Appendix E for model parameter details) and show how their pretraining loss dynamics differ throughout training. In Fig. 22, we see the

---

[12]We use Spanish instead of the International Phonetic Language (IPA) as IPA has tokens that are not compatible with the OLMo 2 tokenizer. Examples of the English to Spanish task are shown in Appendix C. We also change the in-context labels to have no semantic meaning in light of [79].

Figure 22: Pretraining loss — The squared-error of each transformer model's predictions on traces interleaved in the style of the training data averaged over each timestep of the trace. For both the train and test data, averaging was done over 40,000 different interleaved traces, each of length 251. The horizontal black dotted line is the averaged squared-error of a predictor that computes an estimate of the underlying system dynamics by using the Moore-Penrose pseudoinverse of the observed data.

performance of model training checkpoints on data that is in the style of what the model saw during training as specified in Section 2.2. The dotted curves are the models' performance on freshly drawn interleavings of traces from a held-out test library, while the crossed curves are on freshly drawn interleavings of traces from the training library.

In Fig. 22, it is clear that larger models decrease their pretraining loss earlier in training than smaller models. Furthermore, classic overfitting behavior is evident, especially in the larger models. We see that the pretraining error on the held-out dataset deteriorates late in training, while the error on the training data continues to decrease. For the "big" model, its error on the training data even goes lower than the fundamental lower bound imposed by the error of the perfect pseudoinverse predictor.

## 6.1 Training dynamics with respect to the pretraining loss

Taking inspiration from [16], in Fig. 23 we look at when the ability to resume a recalled sequence, the ability to continue predicting a recalled sequence, and the ability to restart ICL on a previously unseen system emerge with respect to the pretraining loss as opposed to the number of training examples. This is to see if these abilities emerge at a specific pretraining loss threshold that is independent of the model size. This pretraining loss corresponds to the held-out loss specified in Section 6. The red vertical line in the plots is the pretraining loss achieved by the pseudoinverse predictor specified in (4) averaged over all context indices. The blue horizontal line in Figs. 23e and 23f is the median error of the pseudoinverse predictor at the specific index corresponding to the curve plotted.

Although figures in Appendices E.1 and F show that larger models develop these emergent abilities before smaller models with respect to the number of training examples, in Fig. 23b we see that all the models exhibit the "1-after final" phase transition at a pretraining loss of around $6.5 \times 10^{-2}$. Furthermore, in Fig. 23d, the "2-after final" phase transition occurs for all models begin at a pretraining loss of a bit more than $1 \times 10^{-1}$. In Fig. 23a the different models are again tightly linked, except for the "tiny" model whose emergence for "1-after final" occurs at a pretraining loss that is around $0.02$ larger than the rest. Overall, these results seems to qualitatively match the conclusions in [16] that pretraining loss is a good indicator for when emergent abilities emerge.

(a) Recall 1 step after final — 1 system haystack.

(b) Recall 1 step after final — 7 system haystack.

(c) Recall 2 steps after final — 1 system haystack.

(d) Recall 2 steps after final — 7 system haystack.

(e) Restart 4 steps into segment 2.

(f) Restart 4 steps into segment 8.

Figure 23: Recall and restart performance vs pretraining loss on held-out data. The red vertical line is the fundamental lower bound pretraining loss achieved by the pseudoinverse predictor. The blue horizontal line in Figs. 23e and 23f is the median error of the pseudoinverse predictor at the specific index plotted in each figure.

# 7  Discussion

At this point, the community understands that ICL is rich and nuanced [33, 76, 47, 59, 29]. We contribute a new dimension of nuance by empirically pointing out that *a single ICL-driven task can be performed, on different tokens, using multiple mechanisms that emerge separately*. When tasks have tight local coherency, there can often be approximate local underspecification — there are multiple ways of knowing what the model is supposed to be doing here. The intrinsic Bayesian orientation of autoregressive next-token prediction [83] means that this local underspecification can get picked up on, while the average-loss-over-tokens driving the training gradients means that an approximately correct mechanism that works most of the time will be rewarded and improved even if it can't solve the task completely. The very success of this approximately correct mechanism during training will further reduce the overall gradient pressure [60] for alternative and potentially better mechanisms, potentially forcing them to develop more slowly [70]. However, structurally, there are certain aspects of tasks that are likely to require the use of the better mechanisms — and it seems that starting an episode of a task might be one of them. These better mechanisms therefore can emerge later in training — but their emergence does not mean that the better information they are acting on will automatically be incorporated by the mechanisms already favored for other parts of the task.

This contrasts with situations where the earlier emerging mechanism foregrounds information that can be used by the later mechanism [71].

# References

[1] Ekin Akyürek, Bailin Wang, Yoon Kim, and Jacob Andreas. In-context language learning: Architectures and algorithms. *arXiv preprint arXiv:2401.12973*, 2024.

[2] Suraj Anand, Michael A. Lepori, Jack Merullo, and Ellie Pavlick. Dual process learning: Controlling use of in-context vs. in-weights strategies with weight forgetting, 2025.

[3] Simran Arora, Sabri Eyuboglu, Aman Timalsina, Isys Johnson, Michael Poli, James Zou, Atri Rudra, and Christopher Ré. Zoology: Measuring and improving recall in efficient language models. *arXiv preprint arXiv:2312.04927*, 2023.

[4] Per Bak. *How nature works : the science of self-organized criticality*. Copernicus, New York, NY, USA, 1996.

[5] BIG bench authors. Beyond the imitation game: Quantifying and extrapolating the capabilities of language models. *Transactions on Machine Learning Research*, 2023.

[6] Adithya Bhaskar, Alexander Wettig, Dan Friedman, and Danqi Chen. Finding transformer circuits with edge pruning. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2024.

[7] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.

[8] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared 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 M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher 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, 2020.

[9] Bryan Chan, Xinyi Chen, András György, and Dale Schuurmans. Toward understanding in-context vs. in-weight learning, 2025.

[10] Stephanie C. Y. Chan, Adam Santoro, Andrew K. Lampinen, Jane X. Wang, Aaditya Singh, Pierre H. Richemond, Jay McClelland, and Felix Hill. Data distributional properties drive emergent in-context learning in transformers, 2022.

[11] Angelica Chen, Ravid Shwartz-Ziv, Kyunghyun Cho, Matthew L Leavitt, and Naomi Saphra. Sudden drops in the loss: Syntax acquisition, phase transitions, and simplicity bias in MLMs. In *The Twelfth International Conference on Learning Representations*, 2024.

[12] Hang Chen, Xinyu Yang, Jiaying Zhu, and Wenya Wang. Quantifying semantic emergence in language models, 2024.

[13] Zhongtian Chen, Edmund Lau, Jake Mendel, Susan Wei, and Daniel Murfet. Dynamical versus bayesian phase transitions in a toy model of superposition, 2023.

[14] Xander Davies, Lauro Langosco, and David Krueger. Unifying grokking and double descent. *arXiv preprint arXiv:2303.06173*, 2023.

[15] Zhe Du, Haldun Balim, Samet Oymak, and Necmiye Ozay. Can transformers learn optimal filtering for unknown systems? *IEEE Control Systems Letters*, 7:3525–3530, 2023.

[16] Zhengxiao Du, Aohan Zeng, Yuxiao Dong, and Jie Tang. Understanding emergent abilities of language models from the loss perspective, 2025.

[17] Benjamin L Edelman, Ezra Edelman, Surbhi Goel, Eran Malach, and Nikolaos Tsilivis. The evolution of statistical induction heads: In-context learning markov chains. *arXiv preprint arXiv:2402.11004*, 2024.

[18] Nelson Elhage, Neel Nanda, Catherine Olsson, Tom Henighan, Nicholas Joseph, Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, Nova DasSarma, Dawn Drain, Deep Ganguli, Zac Hatfield-Dodds, Danny Hernandez, Andy Jones, Jackson Kernion, Liane Lovitt, Kamal Ndousse, Dario Amodei, Tom Brown, Jack Clark, Jared Kaplan, Sam McCandlish, and Chris Olah. A mathematical framework for transformer circuits. *Transformer Circuits Thread*, 2021. https://transformer-circuits.pub/2021/framework/index.html.

[19] Deep Ganguli, Danny Hernandez, Liane Lovitt, Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, Nova Dassarma, Dawn Drain, Nelson Elhage, et al. Predictability and surprise in large generative models. In *Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency*, pages 1747–1764, 2022.

[20] Shivam Garg, Dimitris Tsipras, Percy S Liang, and Gregory Valiant. What can transformers learn in-context? a case study of simple function classes. *Advances in Neural Information Processing Systems*, 35:30583–30598, 2022.

[21] Dirk Groeneveld, Iz Beltagy, Pete Walsh, Akshita Bhagia, Rodney Kinney, Oyvind Tafjord, Ananya Harsh Jha, Hamish Ivison, Ian Magnusson, Yizhong Wang, et al. Olmo: Accelerating the science of language models. *arXiv preprint arXiv:2402.00838*, 2024.

[22] Andrey Gromov. Grokking modular arithmetic. *arXiv preprint arXiv:2301.02679*, 2023.

[23] Jesse Hoogland, George Wang, Matthew Farrugia-Roberts, Liam Carroll, Susan Wei, and Daniel Murfet. Loss landscape degeneracy drives stagewise development in transformers, 2025.

[24] John J Hopfield. Neural networks and physical systems with emergent collective computational abilities. *Proceedings of the national academy of sciences*, 79(8):2554–2558, 1982.

[25] Ruomin Huang and Rong Ge. Task descriptors help transformers learn linear models in-context. In *The Thirteenth International Conference on Learning Representations*, 2025.

[26] Ahmed Imtiaz Humayun, Randall Balestriero, and Richard Baraniuk. Deep networks always grok and here is why. In *Forty-first International Conference on Machine Learning*, 2024.

[27] G. Kamradt. Needle in a haystack — pressure testing llms. https://github.com/gkamradt/LLMTest_NeedleInAHaystack, 2023. GitHub repository.

[28] Tanishq Kumar, Blake Bordelon, Samuel J Gershman, and Cengiz Pehlevan. Grokking as the transition from lazy to rich training dynamics. In *The Twelfth International Conference on Learning Representations*, 2024.

[29] Andrew Kyle Lampinen, Stephanie C. Y. Chan, Aaditya K. Singh, and Murray Shanahan. The broader spectrum of in-context learning, 2024.

[30] Edmund Lau, Zach Furman, George Wang, Daniel Murfet, and Susan Wei. The local learning coefficient: A singularity-aware complexity measure, 2024.

[31] Shuaipeng Li, Penghao Zhao, Hailin Zhang, Xingwu Sun, Hao Wu, Dian Jiao, Weiyan Wang, Chengjun Liu, Zheng Fang, Jinbao Xue, et al. Surge phenomenon in optimal learning rate and batch size scaling. *arXiv preprint arXiv:2405.14578*, 2024.

[32] 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*, pages 19565–19594. PMLR, 2023.

[33] Ziqian Lin and Kangwook Lee. Dual operating modes of in-context learning, 2024.

[34] Jerry Weihong Liu, Jessica Grogan, Owen M Dugan, Simran Arora, Atri Rudra, and Christopher Re. Can transformers solve least squares to high precision? In *ICML 2024 Workshop on In-Context Learning*, 2024.

[35] Nelson F. Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and Percy Liang. Lost in the middle: How language models use long contexts. *Transactions of the Association for Computational Linguistics*, 12:157–173, 2024.

[36] Ziming Liu, Eric J Michaud, and Max Tegmark. Omnigrok: Grokking beyond algorithmic data. In *The Eleventh International Conference on Learning Representations*, 2022.

[37] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization, 2019.

[38] Sheng Lu, Irina Bigoulaeva, Rachneet Sachdeva, Harish Tayyar Madabushi, and Iryna Gurevych. Are emergent abilities in large language models just in-context learning? *arXiv preprint arXiv:2309.01809*, 2023.

[39] Ekdeep Singh Lubana, Kyogo Kawaguchi, Robert P. Dick, and Hidenori Tanaka. A percolation model of emergence: Analyzing transformers trained on a formal language, 2024.

[40] Kaifeng Lyu, Jikai Jin, Zhiyuan Li, Simon S Du, Jason D Lee, and Wei Hu. Dichotomy of early and late phase implicit biases can provably induce grokking. *arXiv preprint arXiv:2311.18817*, 2023.

[41] David JC MacKay. A practical bayesian framework for backpropagation networks. *Neural computation*, 4(3):448–472, 1992.

[42] Nischal Mainali and Lucas Teixeira. Exact learning dynamics of in-context learning in linear transformers and its application to non-linear transformers, 2025.

[43] Neil Rohit Mallinar, Daniel Beaglehole, Libin Zhu, Adityanarayanan Radhakrishnan, Parthe Pandit, and Mikhail Belkin. Emergence in non-neural models: grokking modular arithmetic via average gradient outer product. In *NeurIPS 2024 Workshop on Mathematics of Modern Machine Learning*, 2024.

[44] Stephan Mandt, Matthew D. Hoffman, and David M. Blei. Stochastic gradient descent as approximate bayesian inference. *Journal of machine learning research*, 18, 2017.

[45] Francesco Mezzadri. How to generate random matrices from the classical compact groups. *arXiv preprint math-ph/0609050*, 2006.

[46] Eric Michaud, Ziming Liu, Uzay Girit, and Max Tegmark. The quantization model of neural scaling. *Advances in Neural Information Processing Systems*, 36, 2023.

[47] 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?, 2022.

[48] Chris Mingard, Guillermo Valle-Pérez, Joar Skalse, and Ard A. Louis. Is sgd a bayesian sampler? well, almost. *Journal of Machine Learning Research*, 22(79):1–64, 2021.

[49] Mohamad Amin Mohamadi, Zhiyuan Li, Lei Wu, and Danica Sutherland. Grokking modular arithmetic can be explained by margin maximization. In *NeurIPS 2023 Workshop on Mathematics of Modern Machine Learning*, 2023.

[50] Samuel Müller, Noah Hollmann, Sebastian Pineda Arango, Josif Grabocka, and Frank Hutter. Transformers can do bayesian inference, 2024.

[51] Yoonsoo Nam, Nayara Fonseca, Seok Hyeong Lee, Chris Mingard, and Ard A. Louis. An exactly solvable model for emergence and scaling laws in the multitask sparse parity problem, 2024.

[52] Neel Nanda, Lawrence Chan, Tom Lieberum, Jess Smith, and Jacob Steinhardt. Progress measures for grokking via mechanistic interpretability. In *The Eleventh International Conference on Learning Representations*, 2023.

[53] Alex Nguyen and Gautam Reddy. Differential learning kinetics govern the transition from memorization to generalization during in-context learning, 2024.

[54] Eshaan Nichani, Jason D. Lee, and Alberto Bietti. Understanding factual recall in transformers via associative memories, 2024.

[55] Team OLMo, Pete Walsh, Luca Soldaini, Dirk Groeneveld, Kyle Lo, Shane Arora, Akshita Bhagia, Yuling Gu, Shengyi Huang, Matt Jordan, Nathan Lambert, Dustin Schwenk, Oyvind Tafjord, Taira Anderson, David Atkinson, Faeze Brahman, Christopher Clark, Pradeep Dasigi, Nouha Dziri, Michal Guerquin, Hamish Ivison, Pang Wei Koh, Jiacheng Liu, Saumya Malik, William Merrill, Lester James V. Miranda, Jacob Morrison, Tyler Murray, Crystal Nam, Valentina Pyatkin, Aman Rangapur, Michael Schmitz, Sam Skjonsberg, David Wadden, Christopher Wilhelm, Michael Wilson, Luke Zettlemoyer, Ali Farhadi, Noah A. Smith, and Hannaneh Hajishirzi. 2 olmo 2 furious, 2024.

[56] Catherine Olsson, Nelson Elhage, Neel Nanda, Nicholas Joseph, Nova DasSarma, Tom Henighan, Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, et al. In-context learning and induction heads. *arXiv preprint arXiv:2209.11895*, 2022.

[57] Jane Pan, Tianyu Gao, Howard Chen, and Danqi Chen. What in-context learning "learns" in-context: Disentangling task recognition and task learning, 2023.

[58] Madhur Panwar, Kabir Ahuja, and Navin Goyal. In-context learning through the bayesian prism, 2024.

[59] Core Francisco Park, Ekdeep Singh Lubana, Itamar Pres, and Hidenori Tanaka. Competition dynamics shape algorithmic phases of in-context learning, 2025.

[60] Mohammad Pezeshki, Oumar Kaba, Yoshua Bengio, Aaron C Courville, Doina Precup, and Guillaume Lajoie. Gradient starvation: A learning proclivity in neural networks. *Advances in Neural Information Processing Systems*, 34:1256–1272, 2021.

[61] Mohammad Pezeshki, Amartya Mitra, Yoshua Bengio, and Guillaume Lajoie. Multi-scale feature learning dynamics: Insights for double descent. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvari, Gang Niu, and Sivan Sabato, editors, *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pages 17669–17690. PMLR, 17–23 Jul 2022.

[62] Alethea Power, Yuri Burda, Harri Edwards, Igor Babuschkin, and Vedant Misra. Grokking: Generalization beyond overfitting on small algorithmic datasets. *arXiv preprint arXiv:2201.02177*, 2022.

[63] Lucas Prieto, Melih Barsbey, Pedro AM Mediano, and Tolga Birdal. Grokking at the edge of numerical stability. *arXiv preprint arXiv:2501.04697*, 2025.

[64] Nived Rajaraman, Marco Bondaschi, Ashok Vardhan Makkuva, Kannan Ramchandran, and Michael Gastpar. Transformers on markov data: Constant depth suffices. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024.

[65] Allan Raventós, Mansheej Paul, Feng Chen, and Surya Ganguli. Pretraining task diversity and the emergence of non-Bayesian in-context learning for regression. *Advances in Neural Information Processing Systems*, 36, 2024.

[66] Gautam Reddy. The mechanistic basis of data dependence and abrupt learning in an in-context classification task, 2023.

[67] Michael Eli Sander, Raja Giryes, Taiji Suzuki, Mathieu Blondel, and Gabriel Peyré. How do transformers perform in-context autoregressive learning ? In Ruslan Salakhutdinov, Zico Kolter, Katherine Heller, Adrian Weller, Nuria Oliver, Jonathan Scarlett, and Felix Berkenkamp, editors, *Proceedings of the 41st International Conference on Machine Learning*, volume 235 of *Proceedings of Machine Learning Research*, pages 43235–43254. PMLR, 21–27 Jul 2024.

[68] Rylan Schaeffer, Brando Miranda, and Sanmi Koyejo. Are emergent abilities of large language models a mirage? *Advances in Neural Information Processing Systems*, 36, 2024.

[69] Hinrich Schutze and Christopher Manning. I preliminaries. In *Foundations of Statistical Natural Language Processing*. MIT Press, United States, 1999.

[70] Harshay Shah, Kaustav Tamuly, Aditi Raghunathan, Prateek Jain, and Praneeth Netrapalli. The pitfalls of simplicity bias in neural networks. *Advances in Neural Information Processing Systems*, 33:9573–9585, 2020.

[71] Aaditya K. Singh, Ted Moskovitz, Sara Dragutinovic, Felix Hill, Stephanie C. Y. Chan, and Andrew M. Saxe. Strategy coopetition explains the emergence and transience of in-context learning, 2025.

[72] Aaditya K. Singh, Ted Moskovitz, Felix Hill, Stephanie C. Y. Chan, and Andrew M. Saxe. What needs to go right for an induction head? a mechanistic study of in-context learning circuits and their formation, 2024.

[73] Daniel Soudry, Elad Hoffer, Mor Shpigel Nacson, Suriya Gunasekar, and Nathan Srebro. The implicit bias of gradient descent on separable data. *Journal of Machine Learning Research*, 19:1–57, 2018.

[74] Johannes Von Oswald, Eyvind Niklasson, Ettore Randazzo, João Sacramento, Alexander Mordvintsev, Andrey Zhmoginov, and Max Vladymyrov. Transformers learn in-context by gradient descent. In *International Conference on Machine Learning*, pages 35151–35174. PMLR, 2023.

[75] Lean Wang, Lei Li, Damai Dai, Deli Chen, Hao Zhou, Fandong Meng, Jie Zhou, and Xu Sun. Label words are anchors: An information flow perspective for understanding in-context learning, 2023.

[76] Xiaolei Wang, Xinyu Tang, Wayne Xin Zhao, and Ji-Rong Wen. Investigating the pre-training dynamics of in-context learning: Task recognition vs. task learning, 2024.

[77] Sumio Watanabe. *Algebraic geometry and statistical learning theory / Sumio Watanabe*. Cambridge monographs on applied and computational mathematics ; 25. Cambridge University Press, Cambridge ;, 2009.

[78] Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, et al. Emergent abilities of large language models. *Transactions on Machine Learning Research*, 2022.

[79] Jerry Wei, Le Hou, Andrew Lampinen, Xiangning Chen, Da Huang, Yi Tay, Xinyun Chen, Yifeng Lu, Denny Zhou, Tengyu Ma, et al. Symbol tuning improves in-context learning in language models. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 968–979, 2023.

[80] 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.

[81] Noam Wies, Yoav Levine, and Amnon Shashua. The learnability of in-context learning, 2023.

[82] Jingfeng Wu, Difan Zou, Zixiang Chen, Vladimir Braverman, Quanquan Gu, and Peter L. Bartlett. How many pretraining tasks are needed for in-context learning of linear regression?, 2024.

[83] Sang Michael Xie, Aditi Raghunathan, Percy Liang, and Tengyu Ma. An explanation of in-context learning as implicit bayesian inference. *arXiv preprint arXiv:2111.02080*, 2021.

[84] Kayo Yin and Jacob Steinhardt. Which attention heads matter for in-context learning?, 2025.

[85] Wei Zhang and Bowen Zhou. Learning to update auto-associative memory in recurrent neural networks for improving sequence memorization. *arXiv preprint arXiv:1709.06493*, 2017.

[86] Yedi Zhang, Aaditya K. Singh, Peter E. Latham, and Andrew Saxe. Training dynamics of in-context learning in linear attention, 2025.

[87] Ziqian Zhong, Ziming Liu, Max Tegmark, and Jacob Andreas. The clock and the pizza: Two stories in mechanistic explanation of neural networks. *Advances in Neural Information Processing Systems*, 36, 2024.

# A Extended Related Work

**Emergence**   Benchmark performance of large language models (LLMs) has been observed to improve abruptly at certain scales in a seemingly unpredictable manner [78, 19], leading to the conclusion that LLMs exhibit *emergent* abilities, where different abilities may emerge at different scales or points during training. A natural question was posed by [68]: is emergence a mirage due to the discrete nature of the evaluation metrics used? Would alternative continuous metrics (e.g. logits instead of generated responses) show continuous improvement and thus lead to more predictable performance? Or are these phase transitions real? At this point, the reality of phase-transitions in abilities during training is well established, with arguably the strongest evidence of this coming from the "grokking" literature [62, 52, 22, 87, 43, 26, 49, 73, 36, 63, 61, 14] where there is a clear emergence of substantially improved generalization performance after what looks like a long period without improvement during training. Similar behavior has been observed in stylized regression-style examples that are very different from LLMs and ICL [40, 28, 51]. The toy problem in our paper is both regression-style and very much related to ICL and LLMs.

What exactly drives emergence is also an ongoing topic of investigation. While earlier work talked about model sizes and total compute, the story now is more nuanced. Powerful evidence connects the emergence of abilities to the pretraining losses attained [16], other information-oriented metrics [12], and the idea that more complex or specialized abilities can only emerge after a model acquires a prerequisite abilities during training [11, 39].

**In-context learning**   Our understanding of in-context learning (ICL) has greatly benefited from experiments in simple domains, with many works studying ICL through linear dynamics [20, 25, 82]. To add a proxy for the order-dependence of language, works have used sequences from Markov chains to study pretraining, revealing emergent behaviors related to the formation of induction heads [17, 67, 15, 32].

Recent works have analyzed the role of labels in ICL, giving additional perspective into information flow and the inherent limitations of ICL [47, 75, 25]. This is closely linked to the literature on in-weights learning (IWL) [9, 2], which reveals structural contrasts in the learning mechanisms models can leverage.

Mechanistic explorations of ICL have advanced through both top-down classification of higher-order attention head behaviors [18, 56, 84] and bottom-up analyses tracing the dynamics of these features [72]. From analyzing the learning dynamics of ICL in various setups [42, 86], works have found that ICL enables a mixture of emergent behaviors [66, 10, 53, 38] and competes with other learning mechanisms [59]. Contemporaneous work has also explored the transient dynamics of ICL [71], finding that a form of IWL both cooperates and competes with ICL in attention-only models.

With learning dynamics displaying various surprising phenomena in different setups, works have also tackled discerning whether there are distinct modes of ICL (e.g. task learning (TL) vs. task recognition (TR)[13]) through various lenses like ICL risk (premature convergence), competition dynamics, and differing data distributions [57, 76, 33, 81]. Other work has argued for broader definitions than the "dual mode hypothesis", pushing for ICL to be understood as a spectrum of behaviors or mixture of algorithms [29, 59]. In a sense, our findings here add another nuance to this story since we provide evidence that different modes of ICL can be active at the same time for the same task, just on different tokens, and that these modes can emerge separately during training.

**Associative recall**   With the resurgence of state-space models, [3] highlighted associative recall [24, 85] as the key capability where attention-based models hold a distinct advantage. This was formalized by introducing a discrete toy problem they called Multi-Query Associative Recall (MQAR). However, in basic MQAR, the solution involves the exact copying of the right token from context rather than recalling something conceptual that has to be applied locally. Meanwhile, in continuous Gauss-Markov models, next-observation prediction can still be performed effectively using state-space models, as least-squares algorithms can be adapted into a streaming form via standard Woodbury-

---

[13]From [76]: "TR refers to the ability of an LLM to recognize the target task from demonstrations and only utilize its own knowledge obtained from pre-training to solve the task, while TL refers to the ability of an LLM to solve the target task solely based on demonstrations." The NLP example we choose in Section 5 of this paper falls into the "Task Recognition" subtype.

matrix tricks. By using MQAR-style discrete symbolic labels to segment time-series drawn from Gauss-Markov models, we effectively create a naturally continuous toy problem where the information to be recalled is not an exact copying.

**Are deep neural networks Bayesian predictors?**    The Bayesian perspective on machine-learning generally is longstanding and consequently has been used to better understand the ICL capabilities of deep neural networks. [83] frames in-context learning as implicit Bayesian inference. Using a simple, synthetic task setup, this work showed theoretically that if a model makes its predictions according to its pretraining distribution, it will asymptotically approach the performance of the Bayes optimal predictor as the number of in-context examples for this task increase to infinity. [58] builds off of these ideas, and empirically investigates how close a transformer model's predictions are to the Bayes optimal predictor. They find that high-capacity transformers successfully perform Bayes-optimal prediction for Gaussian mixtures and other more complex mixtures. Interestingly, they hypothesize that low-capacity transformers attempt to approximate Bayesian inference through gradient descent. The facts that gradient descent reaches the global optimum for convex problems and [74] shows that multiple self-attention layers in a model can simulate multiple steps of gradient descent lead the authors in [58] to believe that higher-capacity transformers perform Bayesian inference because they can take more steps of gradient descent.

Many works have shown that gradient descent can approximate a Bayes-optimal predictor. For example, [44] uses stochastic gradient descent with a constant learning rate to sample from an approximate posterior, by connecting stochastic optimization with the stochastic gradients used in Markov chain Monte Carlo. Additionally, [48] uses Gaussian processes to estimate a posterior and empirically computes the posterior probability of a deep neural network trained with stochastic gradient descent. Through comparing the two posterior distributions, they conclude that the Bayes posterior is a first-order correlate with the neural network's, while second-order differences can come from hyperparameter choices when training the model.

Relating this work to the emergence that we observe for transformer models on our in-context recall task, singular learning theory predicts that Bayesian learners have phase transitions in their generalization error [77]. Grounding their empirical investigation in this theory, [23, 13] study the loss landscape of transformer and autoencoder models, and are able to predict the occurrence of a phase transition by computing the local learning coefficient [30].

# B   Edge Pruning Circuit Visualization

Section 4.2 describes how optimizing over continuous masks on disentagled transformer nodes leads to a sparse computation graph after the continuous masks are quantized to 0 or 1. To perform this quantization, the quantization threshold must be set. We use binary search to find the smallest threshold (within 1e-5 precision) such that the pruned model's edge sparsity is close to the target edge sparsity of 0.98. Importantly, this does not force a viable path from the start and end of the residual stream, it only captures significant contributions over edges with respect to the loss. In Fig. 24, the pruned circuit is visually presented in a directed computation graph. Nodes in this computation graph are specific QKV matrices of attention heads, entire attention heads, MLPs at a specific layer, or the output of the residual stream. The output of the residual stream is labeled as `resid_post`, MLPs are labeled as `m{layer_num}`, attention heads are labeled as `a{layer_num}.h{head_num}`, and `.q`, `.k`, and `.v` stand for QKV matrices respectively.



Figure 24: 1-after final open symbol circuit in orthogonal systems. The output of the residual stream is labeled as `resid_post`, MLPs are labeled as `m{layer_num}`, attention heads are labeled as `a{layer_num}.h{head_num}`, and `.q`, `.k`, and `.v` stand for QKV matrices respectively.

## C   NLP Prompts

In this section, we provide sample promplts for the NLP experiments from Section 5. We use 100 English prompts from the IPA transliterate benchmark task in BIG-bench [5], and pass these English prompts through Google Translate to obtain their Spanish translations.

Below are examples of three different task instances for the English to Spanish for the in-weights associative recall task.

```
"input":  Spanish:  Inició el ascenso del Imperio Británico en la India.
English:
"target":  It started the rise of the British Empire in India.
"input":  Spanish:  Hicieron acusaciones sobre la plataforma.  English:
"target":  They made accusations about the platform.
"input":  Spanish:  Che fue otro dictador que ascendió al poder mediante un
levantamiento militar.  English:
"target":  Che was another dictator who rose to power by military uprising.
```

Below are examples of three different task instances for the English to Spanish for the in-context associative recall task. Notice that they are the same as the in-weights associative recall examples, except the semantically meaningful labels "English" and "Spanish" are swapped for the labels "X" and "Y".

```
"input":  X: Inició el ascenso del Imperio Británico en la India.  Y:
"target":  It started the rise of the British Empire in India.
"input":  X: Hicieron acusaciones sobre la plataforma.  Y:
"target":  They made accusations about the platform.
"input":  X: Che fue otro dictador que ascendió al poder mediante un
levantamiento militar.  Y:
"target":  Che was another dictator who rose to power by military uprising.
```

## D   Additional "Out-of-distribution" Experiment in OLMo

In this section we explore a counterpart to the misdirecting-to-an-unseen-token task (Fig. 17 and results shown in Fig. 18) using training checkpoints of the OLMo 2 model. Here, we consider the task where the symbolic label that is used for directing the associative recall (e.g. the open SPL in the Markovian problem) is replaced with a different label that is so far unseen. To adapt this task to a natural language setup, we keep our translation task with two few-shot examples exactly the same as in Section C, but alter the final testing prompt to not use the 'Y:' label tokens and to instead use a previously unseen 'Z:' token.

As such, at inference time, the first token prompt functionally looks like this.

```
X: Inició el ascenso del Imperio Británico en la India.
Y: It started the rise of the British Empire in India.
X: Hicieron acusaciones sobre la plataforma.
Y: They made accusations about the platform.
X: Che fue otro dictador que ascendió al poder mediante un levantamiento
militar.
Z:
```

A comparison of the two plots in Fig. 25 shows the results. Note that the figure on the left in Fig. 25 is identical to the figure on the right in Fig. 21 and has been copied here for ease of visual comparisons. The key behavior to notice is that there is no qualitative difference between the two blue curves in the left and right subfigures of Fig. 25. What this means is that the second-token-prediction task is functionally unaffected by the "misdirection" by the unseen token 'Z:'. This gives further evidence that the second token prediction does not rely on the in-context label in order to complete the task once it has started. This matches what we saw in our toy model as well, lending further support for our multi-mechanism conjecture to hold in pretrained LLMs.

Meanwhile, for the first token task (black curve in Fig. 25), in a sense, there is objectively no right answer for what to do when presented with 'Z:'. Nonetheless, we measure accuracy against a correct ground-truth English translation. It seems that the model accuracy oscillates wildly as it trains. This oscillation during training is an unexpected behavior for single-epoch training and deserves more investigation in future work.



Figure 25: Comparison of the original in-context associative recall task (left) vs the in-context associative recall task with an irrelevant token (right). Importantly, we see that the irrelevant token does not effect the second token accuracy when compared to the original task.

# E   The Effects of Model Size

| Model Name | $n_{params}$ | $n_{layers}$ | $d_{model}$ | $n_{heads}$ | $d_{head}$ | **Learning Rate** | **Batch Size** |
|---|---|---|---|---|---|---|---|
| Orthogonal Tiny | 212K | 3 | 72 | 6 | 12 | $1.7 \times 10^{-4}$ | 2048 |
| Orthogonal Small | 701K | 6 | 96 | 6 | 16 | $4.5 \times 10^{-5}$ | 1024 |
| Orthogonal Medium | 2.42M | 12 | 128 | 8 | 16 | $1.6 \times 10^{-5}$ | 512 |
| Orthogonal Big | 10.7M | 24 | 192 | 12 | 16 | $1.5 \times 10^{-5}$ | 640 |
| Identity Tiny | 212K | 3 | 72 | 6 | 12 | $6.3 \times 10^{-5}$ | 8192 |
| Identity Small | 701K | 6 | 96 | 6 | 16 | $3.2 \times 10^{-5}$ | 4096 |
| Identity Medium | 2.42M | 12 | 128 | 8 | 16 | $1.6 \times 10^{-5}$ | 1024 |
| Identity Big | 10.7M | 24 | 192 | 12 | 16 | $1.3 \times 10^{-5}$ | 512 |

Table 2: Model size and training hyperparameters.

In order to test the effect of model size on our emergence results, we trained models across 4 different model sizes as shown in Table 2. We originally tuned the learning rate for the medium model with a batch size of 512 on a single GPU. Following the model scaling that was done[14] in [8], when decreasing the size of our model from "medium" to "small" we halved the number of layers, multiplied the model dimension by $0.75$, maintained the same head dimension, and doubled the learning rate. To go from "small" to "tiny", we used the same process except we chose the head dimension to be $12$ to maintain an integer value for the number of heads. If we would have maintained the head dimension of $16$, the tiny model would have had $4.5$ heads. For this reason, $12$ was chosen as the head dimension since it is the largest integer less than $16$ that leads to an integer when dividing $72$. To go from "medium" to "big", we doubled the number of layers, multiplied the model dimension by $1.5$, maintained the same head dimension, and multiplied the learning rate by $\frac{5}{6}$. This scaling was used for the identity models and batch size was not taken into account. It was later brought to our attention that the learning rate should also scale with the batch size. For our later orthogonal runs, we additionally adopted the square-root learning-rate scaling as indicated in [31]. Specifically, we took the learning rate we had scaled by model size and multiplied it by $\sqrt{\frac{batch\_size}{512}}$. However, we have not verified if this scaling is the best way to proceed with our training and further

---

[14]Alternatively, one could follow the model scaling done in [21]. There, they also halve the number of layers when decreasing the model size, but decide to halve the model dimension and number of attention heads as well.

testing is still required. Furthermore, we trained the identity models on Nvidia GH200 with 80GB of RAM whereas the orthogonal models were trained on one or two L40S GPUs with 48GB of RAM each. This is the reason for the differing batch sizes across different model types.

### E.1 When does the associative recall ability emerge for predicting the first test segment observation?

Evidence that larger models see earlier emergence is shown in Figs. 26, 27 and 28, where the number of training examples until the phase transition is shown for each model size for the identity and orthogonal systems. This was done by recording the checkpoint before and after the mean-squared error for the pure recall task dropped below the cutoff values of $0.4$ for the identity systems and $0.5$ for the orthogonal systems. These cutoff values were chosen by visual inspection.



Figure 26: Emergence of associative recall on different haystack lengths across model sizes - The plot above shows that associative recall emerges much earlier for larger model sizes both when trained on identity systems as well as when trained on orthogonal systems. We also see that associative recall emerges earlier when there is only one system in the haystack but levels out as the model learns to generalize associative recall regardless of the haystack size.

(a) Orthogonal: 1 system haystack.    (b) Orthogonal: 2 system haystack.    (c) Orthogonal: 3 system haystack.

(d) Orthogonal: 17 system haystack.    (e) Orthogonal: 18 system haystack.    (f) Orthogonal: 19 system haystack.

Figure 27: Emergence of associative recall in varying model sizes across haystack lengths for Orthogonal systems



(a) Identity: 1 system haystack.    (b) Identity: 2 system haystack.    (c) Identity: 3 system haystack.
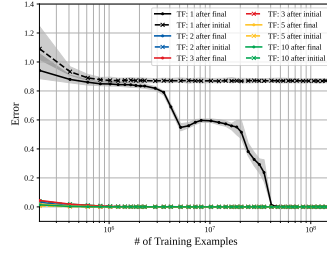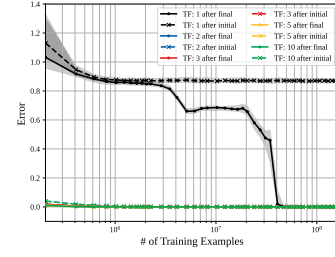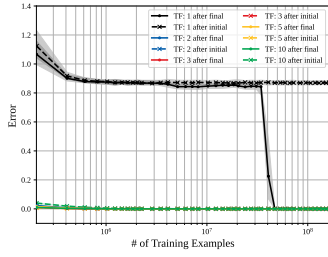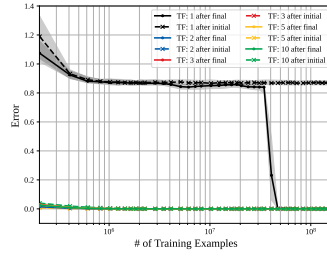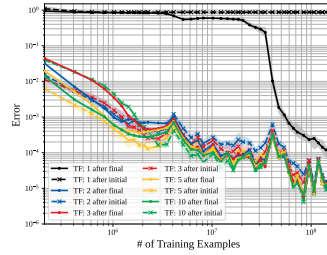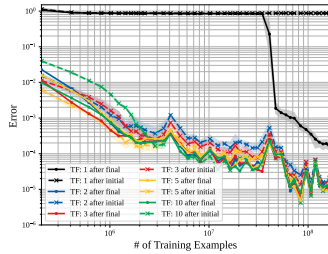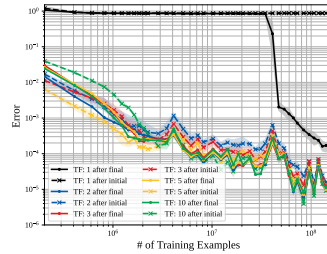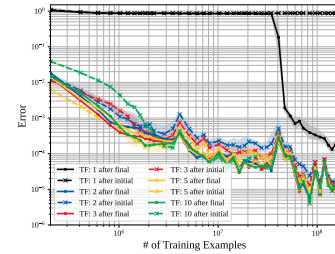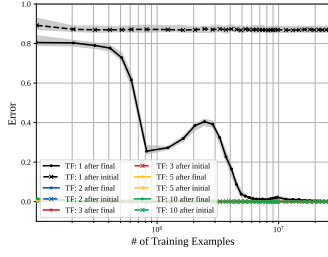
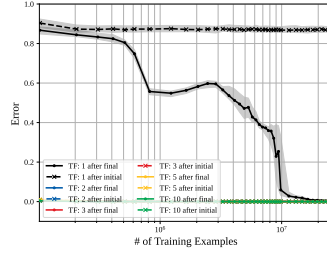(d) Identity: 17 system haystack.    (e) Identity: 18 system haystack.    (f) Identity: 19 system haystack.

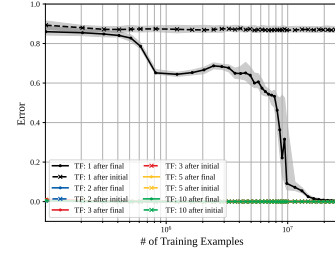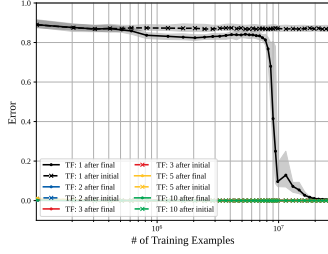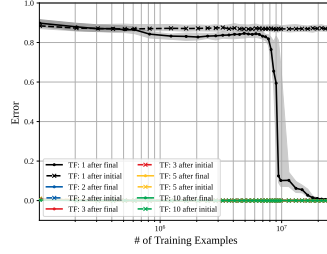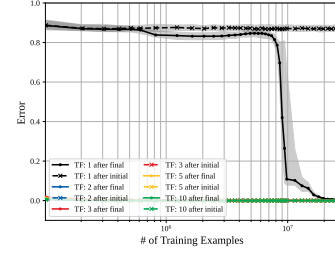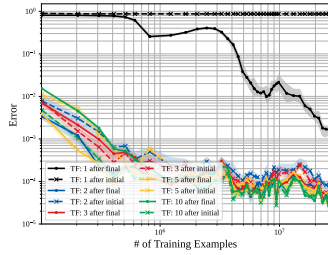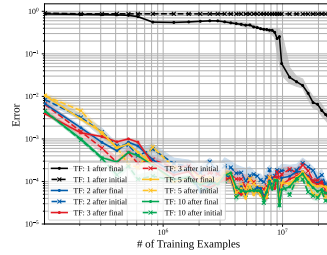Figure 28: Emergence of associative recall in varying model sizes across haystack lengths for Identity systems

# F   More training dynamics plots

Following the summary results from above, we provide full training dynamics plots for our different model sizes and show the results when tested on haystack lengths of 1, 2, 3, 17, 18 and 19. For your convenience here are the different figures: Orthogonal tiny (linear scale) Fig. 29, Orthogonal tiny (log scale) Fig. 30, Orthogonal small (linear scale) Fig. 31, Orthogonal small (log scale) Fig. 32, Orthogonal medium (linear scale) Fig. 33, Orthogonal medium (log scale) Fig. 34, Orthogonal big (linear scale) Fig. 35, Orthogonal big (log scale) Fig. 36, Identity tiny (linear scale) Fig. 37, Identity tiny (log scale) Fig. 38, Identity small (linear scale) Fig. 39, Identity small (log scale) Fig. 40, Identity medium (linear scale) Fig. 41, Identity medium (log scale) Fig. 42, Identity big (linear scale) Fig. 43, Identity big (log scale) Fig. 44



(a) Orthogonal: 1 system haystack.

(b) Orthogonal: 2 system haystack.

(c) Orthogonal: 3 system haystack.

(d) Orthogonal: 17 system haystack.

(e) Orthogonal: 18 system haystack.

(f) Orthogonal: 19 system haystack.

Figure 29: Performance of tiny orthogonal model (212K params) across training — linear-scale.



(a) Orthogonal: 1 system haystack.

(b) Orthogonal: 2 system haystack.

(c) Orthogonal: 3 system haystack.

(d) Orthogonal: 17 system haystack.

(e) Orthogonal: 18 system haystack.

(f) Orthogonal: 19 system haystack.

Figure 30: Performance of tiny orthogonal model (212K params) across training — log-scale.

(a) Orthogonal: 1 system haystack. (b) Orthogonal: 2 system haystack. (c) Orthogonal: 3 system haystack.

(d) Orthogonal: 17 system haystack. (e) Orthogonal: 18 system haystack. (f) Orthogonal: 19 system haystack.

Figure 31: Performance of small orthogonal model (701K params) across training — linear-scale.



(a) Orthogonal: 1 system haystack. (b) Orthogonal: 2 system haystack. (c) Orthogonal: 3 system haystack.

(d) Orthogonal: 17 system haystack. (e) Orthogonal: 18 system haystack. (f) Orthogonal: 19 system haystack.

Figure 32: Performance of small orthogonal model (701K params) across training — log-scale.

(a) Orthogonal: 1 system haystack.  (b) Orthogonal: 2 system haystack.  (c) Orthogonal: 3 system haystack.

(d) Orthogonal: 17 system haystack.  (e) Orthogonal: 18 system haystack.  (f) Orthogonal: 19 system haystack.

Figure 33: Performance of medium orthogonal model (2.42M params) across training — linear-scale.
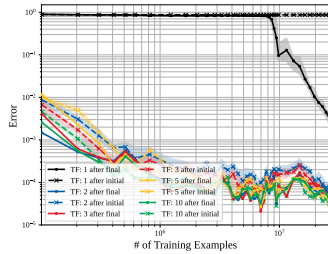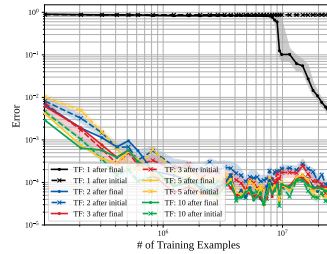


(a) Orthogonal: 1 system haystack.  (b) Orthogonal: 2 system haystack.  (c) Orthogonal: 3 system haystack.

(d) Orthogonal: 17 system haystack.  (e) Orthogonal: 18 system haystack.  (f) Orthogonal: 19 system haystack.

Figure 34: Performance of medium orthogonal model (2.42M params) across training — log-scale.

(a) Orthogonal: 1 system haystack.   (b) Orthogonal: 2 system haystack.   (c) Orthogonal: 3 system haystack.

(d) Orthogonal: 17 system haystack.   (e) Orthogonal: 18 system haystack.   (f) Orthogonal: 19 system haystack.

Figure 35: Performance of big orthogonal model (10.7M params) across training — linear-scale.



(a) Orthogonal: 1 system haystack.   (b) Orthogonal: 2 system haystack.   (c) Orthogonal: 3 system haystack.

(d) Orthogonal: 17 system haystack.   (e) Orthogonal: 18 system haystack.   (f) Orthogonal: 19 system haystack.

Figure 36: Performance of big orthogonal model (10.7M params) across training — log-scale.
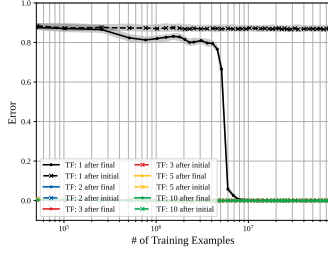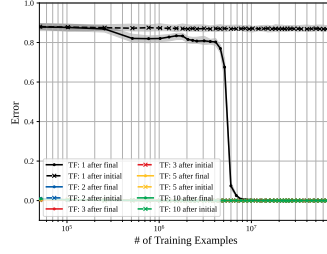
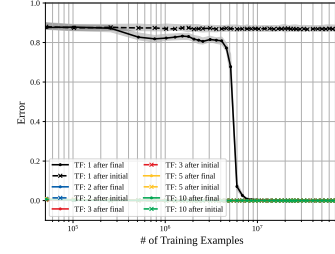(a) Identity: 1 system haystack.
(b) Identity: 2 system haystack.
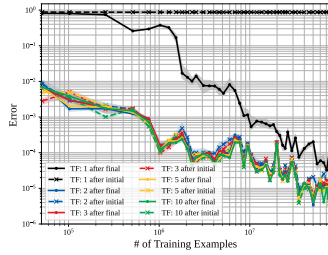(c) Identity: 3 system haystack.

(d) Identity: 17 system haystack.
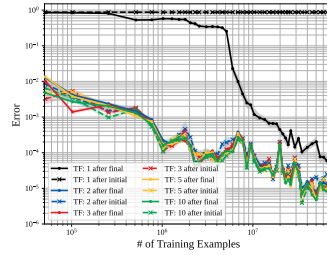(e) Identity: 18 system haystack.
(f) Identity: 19 system haystack.

Figure 37: Performance of tiny identity model (212K params) across training — linear-scale.



(a) Identity: 1 system haystack.
(b) Identity: 2 system haystack.
(c) Identity: 3 system haystack.

(d) Identity: 17 system haystack.
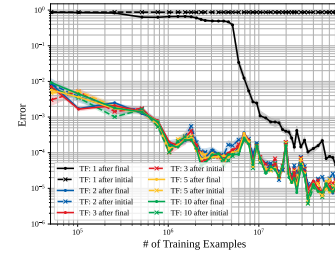(e) Identity: 18 system haystack.
(f) Identity: 19 system haystack.

Figure 38: Performance of tiny identity model (212K params) across training — log-scale.
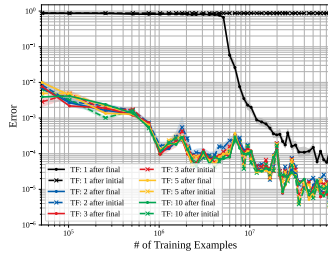
(a) Identity: 1 system haystack.

(b) Identity: 2 system haystack.

(c) Identity: 3 system haystack.

(d) Identity: 17 system haystack.

(e) Identity: 18 system haystack.

(f) Identity: 19 system haystack.

Figure 39: Performance of small identity model (701K params) across training — linear-scale.
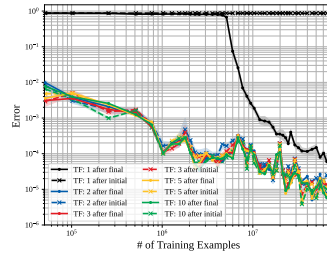


(a) Identity: 1 system haystack.
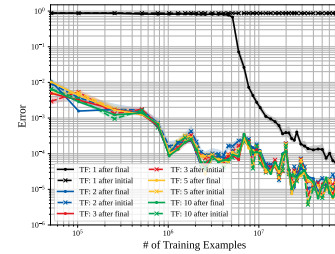
(b) Identity: 2 system haystack.

(c) Identity: 3 system haystack.

(d) Identity: 17 system haystack.

(e) Identity: 18 system haystack.

(f) Identity: 19 system haystack.

Figure 40: Performance of small identity model (701K params) across training — log-scale.

(a) Identity: 1 system haystack.  (b) Identity: 2 system haystack.  (c) Identity: 3 system haystack.

(d) Identity: 17 system haystack.  (e) Identity: 18 system haystack.  (f) Identity: 19 system haystack.

Figure 41: Performance of medium identity model (2.42M params) across training — linear-scale.



(a) Identity: 1 system haystack.  (b) Identity: 2 system haystack.  (c) Identity: 3 system haystack.

(d) Identity: 17 system haystack.  (e) Identity: 18 system haystack.  (f) Identity: 19 system haystack.

Figure 42: Performance of medium identity model (2.42M params) across training — log-scale.

(a) Identity: 1 system haystack.

(b) Identity: 2 system haystack.

(c) Identity: 3 system haystack.

(d) Identity: 17 system haystack.

(e) Identity: 18 system haystack.

(f) Identity: 19 system haystack.

Figure 43: Performance of big identity model (10.7M params) across training — linear-scale.



(a) Identity: 1 system haystack.

(b) Identity: 2 system haystack.

(c) Identity: 3 system haystack.

(d) Identity: 17 system haystack.

(e) Identity: 18 system haystack.

(f) Identity: 19 system haystack.

Figure 44: Performance of big identity model (10.7M params) across training — log-scale.

40

# G    When does the model learn to restart ICL for a new system?



(a) $2.05 \times 10^6$ training examples.

(b) $2.56 \times 10^6$ training examples.

(c) $3.07 \times 10^6$ training examples.

(d) $3.58 \times 10^6$ training examples.

(e) $4.10 \times 10^6$ training examples.

(f) $7.17 \times 10^6$ training examples.

(g) $1.33 \times 10^7$ training examples.

(h) $3.17 \times 10^7$ training examples.

(i) $3.43 \times 10^7$ training examples.

(j) $6.25 \times 10^7$ training examples.

Figure 45: Restarting for a new system — Prediction error vs. the position of a previously unseen system segment in the haystack. Position 1 is the beginning and position 20 is the end. The colored curves correspond to different indices into each system segment.

Continuing on from Section 3.2.2, this supplemental appendix takes a more comprehensive look at how the ability emerges for the model to restart its ICL predictions for a new system.

Fig. 45 depicts the prediction error of different model checkpoints on initial segments from previously unseen systems when these segments are at different positions in the haystack. The first segment in

41

the haystack is denoted position 1 and the last segment is denoted position 20. Furthermore, we look at the prediction error of the model for 1 through 8 indices after the initial open symbol initiating the segment. Subfigure (a) clearly shows that at this early checkpoint $2.05 \times 10^6$, while the model has learned to in-context learn the dynamics for the first segment, there is substantial interference with later segments where performance is poor. Rapidly across the checkpoints illustrated in (b),(c),(d),(e) the model learns how to interpret the symbolic tokens well enough to properly restart its predictions for each new segment and by checkpoint $4.1 \times 10^6$, it is quite good — except at predicting the fundamentally unpredictable second position in a new segment.

Subfigures (f),(g),(h),(i) of Fig. 45 show the more gradual improvements in this ability during much later checkpoints — where the focus is largely on improvements to the performance in predicting the second position in a new segment. Finally, Fig. 45j shows the model's good ability to restart ICL at the end of its training.

# H Final performance with position in segment and needle position within haystack



Figure 46: Comparison of recall performance for orthogonal and identity systems. Figs. 46a and 46b show predictions made in the final test segment (black), from a continuing segment (red), and from the initial segment (yellow). Figs. 46c and 46d show median-squared error vs. the position of the needle in the haystack. For these plots, the x-axis value of $-2$ corresponds to a final segment that is the continuation of the segment before with no interruption by open and close symbols. The x-axis value of 0 is the closest needle position to the final segment while the value 18 is the furthest away. The identity model is after $\approx 1.8 \times 10^7$ training examples and the orthogonal model after $\approx 6.25 \times 10^7$ training examples. "Zero" in the legend corresponds to a predictor that always predicts 0.

Using a haystack of length $N = 19$ for the results in this section, in Fig. 46a, the black dots show the performance of the orthogonal 128M parameter model on the final test segment of a needle-in-haystack sequence. The yellow points represent how well it does on the first system it saw — for which the first point is completely unpredictable and the rest do a bit better. To properly understand recall, we add one more possibility represented by the red dots. This represents the alternative of
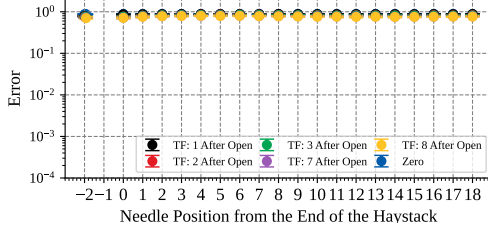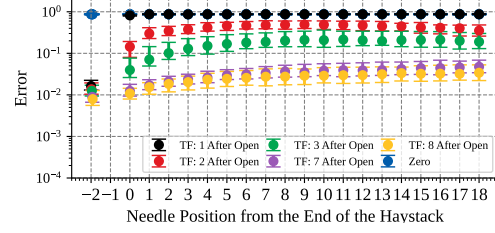
Figure 47: The effect of the needle position throughout training — The median-squared error vs. the position of the needle in the haystack for different training checkpoints of the orthogonal model. For these plots, the x-axis value of $-2$ corresponds to a test segment that is the continuation of the segment before with no interruption by open and close symbols. The x-axis value of 0 is the closest needle position to the test segment while the value 18 is the furthest away. "Zero" in the legend corresponds to a predictor that always predicts 0.

not terminating the final segment in the haystack at all and simply letting it continue on as the test sequence. We see that the black dots perform worse than the red points for merely continuing the last system in the haystack but better than the yellow points showing what happens the first time we see this system. This clearly shows the ability to do recall, but there is worse prediction performance for the second observation in the test segment than the first (see Fig. 12). The conceptual difference represented by the red dots is that performing well here involves no recall at all, merely being able to continue predicting within one sequence.

Figure 46b shows what this looks like for the identity systems using the black dots. The recall of the constant is a bit spotty in the first position but we get more than 90% of the way there with a squared-error of only 0.01 instead of the 1 we would have if we just guessed the all-zero vector. But after that initial observation in the test segment, the quality of understanding the constant-nature of these segments is quite good with all squared-losses around $10^{-4}$ or below.

We can also see how well the model can use recalled information as a function of the position we are trying to predict within the test segment. This parallels the natural language questions in [35] that spawned widespread "needle in haystack" testing for LLMs generally [27]. Figs. 46c and 46d show the effect of the needle position on the prediction performance of the orthogonal and identity models respectively. Here, we can see the first observation prediction quality in black showing higher quality when the needle is in the earlier positions within the context (larger position values in the figure). This is somewhat different from the "U-shaped" curve found for language-based LLMs although both have better performance when asking about the first thing seen. Meanwhile, the observation prediction quality for the second position is actually slightly better if the needle is closest to the test segment.

In Fig. 47, we look at how the effect of the needle position develops throughout the training of the orthogonal model. Again, the median squared error of its predictions on different indices in the test segment are shown for 19 needle positions for many training checkpoints. For indices 7 and 8 in the test segment, the trend of better predictions for closer positions to the test segment is the first trend to develop in training and is sustained throughout. This same trend is also seen in indices 2 and 3 after $4.10 \times 10^6$ training examples have been processed, but by $3.17 \times 10^7$ training examples the model's predictions for the needle positions furthest away from the test segment start to improve.