# Understanding the Staged Dynamics of Transformers in Learning Latent Structure

Rohan Saha [1]    Farzane Aminmansour [1]    Alona Fyshe [1,2]

## Abstract

While transformers can discover latent structure from context, the dynamics of how they acquire different components of the latent structure remain poorly understood. In this work, we use the Alchemy benchmark, to investigate the dynamics of latent structure learning. We train a small decoder-only transformer on three task variants: 1) inferring missing rules from partial contextual information, 2) composing simple rules to solve multi-step sequences, and 3) decomposing complex multi-step examples to infer intermediate steps. By factorizing each task into interpretable events, we show that the model acquires capabilities in discrete *stages*, first learning the coarse grained rules, before learning the complete latent structure. We also identify a crucial asymmetry, where the model can compose fundamental rules robustly, but struggles to decompose complex examples to discover the fundamental rules. These findings offer new insights into understanding *how* a transformer model learns latent structures, providing a granular view of how these capabilities evolve during training.

## 1. Introduction

In the context of language modeling, it is clear that transformers can learn latent structure. For example, models learn syntactic hierarchies without supervision, with attention heads encoding structures such as dependency trees (Hewitt & Manning, 2019; Manning et al., 2020; Ravishankar et al., 2021), and encode causal graphs in their attention mechanisms (Nichani et al., 2024). However, because of the complexity of latent structures in language, it is difficult to carefully study the learning progression of a task with latent structure.

---
[1]Department of Computing Science, University of Alberta, Edmonton, Canada [2]Department of Psychology, University of Alberta, Edmonton, Canada. Correspondence to: Rohan Saha <rohansaha60@gmail.com>.

Here, we explore a setting with a clear and controlled latent structure: the Alchemy benchmark (Wang et al., 2021). The latent structures in the Alchemy benchmark allow us to study the learning progression in multiple settings by carefully controlling what the model observes during training. In each of these settings, we observe staged training dynamics: performance improves in plateaus, followed by sudden jumps, suggesting that models acquire capabilities in stages rather than continuously. Although previous studies have reported staged learning in a variety of settings (Wei et al., 2022; Singh et al., 2024; Chen et al., 2025; Song et al., 2025; Gopalani & Hu, 2025), the complexity of these settings makes it difficult to attribute the stages to actual learning. Our analysis allows us to attribute the stages to specific features of the latent structure.

In the Alchemy setting, players start with a stone whose properties can be altered by applying various potions, changing the stone's features accordingly. The set of stones and the transitions between them, induced by potions, form a "chemistry" represented as a directed graph with eight vertices arranged in a cubic structure (Figure 1, middle). Each vertex corresponds to a stone state, and each edge represents a potion that causes a transition between stone states. Through controlled experiments and careful analyses, we demonstrate that rapid improvements in various prediction tasks align with the successful learning of different aspects of the chemistry. We investigate latent structure learning dynamics by withholding observations of specific latent features, increasing task complexity, and varying compositional demands.

We train a small decoder-only transformer model ($\sim$ 2M parameters) (Vaswani et al., 2017) on three formulations of the Alchemy task. 1) We strategically withhold information about specific potion effects and test the model's ability to infer the effects of these missing potions (Figure 1, left). 2) We investigate compositionality by providing the model with all fundamental, single potion effects and test its ability to compose these single potions to solve novel, multi-step potion sequences (Figure 1, top right). 3) We test the model's decomposition ability by providing it with multi-step potion sequences and requiring it to infer the hidden single-step potion effects (Figure 1, bottom right).
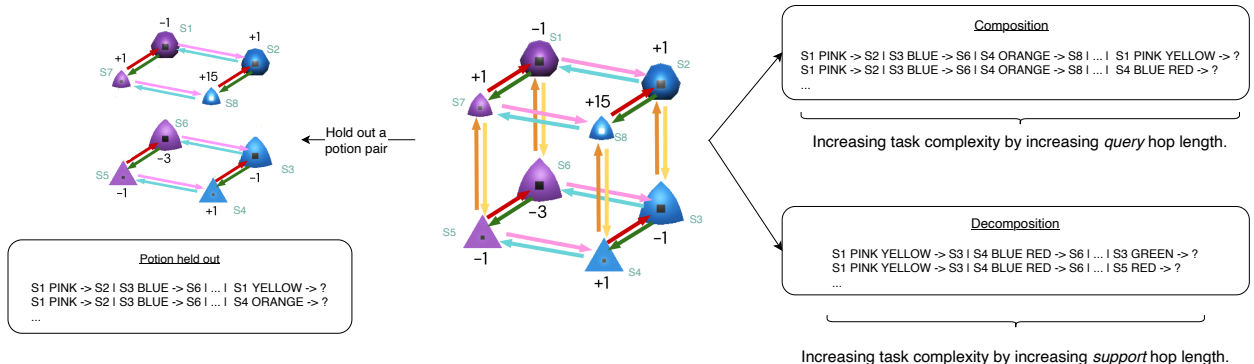
*Figure 1.* Overview of Alchemy chemistry structure and experimental tasks. Middle: Example chemistry with vertices representing stone states connected with bidirectional edges (potions). A chemistry consists of eight stones, and application of potions changes the stone features. Left: Experiment to investigate the staged dynamics of latent structure learning: given a chemistry, all samples for a randomly selected potion pair are withheld (e.g., YELLOW/ORANGE). The model needs to correctly predict the effect of the withheld potions when applied to a query stone. Right: Tasks examining complexity: For composition (top right), each sample consists of the concatenation of all 1-hop (single-step) support transitions for the chemistry, followed by a multi-hop query. For decomposition (bottom right), the model receives the concatenation of all multi-hop transitions for a given hop length in the support, followed by a 1-hop query. Note that in all our experiments, we provide the full enumeration of the stone features in the support and query. Figure adapted from the original Alchemy benchmark (Wang et al., 2021).

Our central contribution is to analyze the staged dynamics during learning of latent structure, by factorizing the model's overall accuracy into interpretable components, corresponding to different latent properties. This approach allows us to demonstrate that the plateaus and jumps in the validation accuracy are not random, but correspond to the model acquiring specific interpretable sub-skills. Concretely, we show:

1. In the potion effect discovery task, the model learns to infer the effect of a withheld potion in distinct stages: it first identifies the general set of possible outcomes, then gradually masters the correct effect.

2. When composing the effects of single-step potions to solve multi-step potion sequences, the model shows no noticeable performance degradation as the potion sequence length increases. Furthermore, our analysis shows that the learning stages for all sub-skills occur at approximately the same time during training, irrespective of the potion sequence length.

3. When decomposing multi-step potion sequences to solve a single-step potion sequence, the model performance degrades with increasing potion sequence length. Our analysis demonstrates that while the model learns the general set of outcomes quickly irrespective of the potion sequence length, the degradation of the model performance is driven by a delayed learning of fine-grained components of the task.

We release our code to encourage further analyses exploring staged dynamics of transformer models.

## 2. Methods

In this section, we introduce the Alchemy dataset, its latent structures, and properties. We then introduce notation to motivate the experiments. Finally, we describe the transformer model architecture and training details.

### 2.1. Dataset

We use the Alchemy dataset (Wang et al., 2021), which is a compositional meta-reinforcement learning benchmark designed to test a model's ability to discover and exploit latent structures. In Alchemy, the player's task is to apply potions to a stone, changing its state. The interactions between the stones and potions is defined by a chemistry, and chemistries differ between rounds. All chemistries share a latent cubic graph structure that describes the interactions of potions and stones (Figure 1, middle).

To solve the Alchemy task, the model cannot memorize the transitions between stones from the cubic graph structure, instead, it must learn to discover and exploit the invariant latent structures that are consistent across all chemistries. While the original alchemy dataset contains chemistries with missing edges (incomplete graphs), in our experiments, we only consider complete graphs[1], where each vertex has

---

[1]Since incomplete chemistries are subgraphs of complete ones, we exclude them to ensure strict disjointness between train and validation sets, and thus preventing any data leakage between them.

three outgoing edges.

The chemistries have the following invariances:

- Each chemistry graph has 8 possible stone states positioned on the vertices of the cubic structure. Each stone is defined by four features: color, size, roundness, and reward. The first three perceptual features ('color', 'size', 'roundness') have three possible values. The 'color' feature can take values pink, violet, blue, 'size' can take values small, medium, large, and 'roundness' can take values pointy, medium_round, and round. The 'reward' feature can take one of four values. This creates 108 unique possible stones across all chemistries. A chemistry can only contain 8 out of the 108 stones.

- Potions come in fixed complementary pairs (RED/-GREEN, YELLOW/ORANGE, PINK/BLUE), and these pairings are consistent across all chemistries. Applying a potion changes the stone's features, effectively moving it along an edge connecting two vertices. For any given chemistry, each potion color corresponds to a single latent transition rule that can be applied to a vertex. In addition, it is crucial to note that a potion's transition rule corresponding to one axis of the chemistry, may not align with a single perceptual feature, and thus a single potion can cause multiple perceptual features to change at once. For example, in the chemistry shown in Figure 1 (middle), the PINK potion leaves the stone's size and roundness unchanged while shifting its color by two levels (pink $\rightarrow$ violet $\rightarrow$ blue). By contrast, YELLOW potions in this chemistry adjust the roundness, and size by a single level but keep the color unchanged. A potion can also have different effects on the stone's reward value. For example, applying the PINK potion to the stone $S1$, increases its reward feature from -1 to +1, but applying PINK to $S6$ changes its reward feature from -3 to -1. Although the specific effect of a potion may differ across chemistries, the complementary paired potion always induces an inverse transition (e.g., if PINK moves a stone from one color state to another, its complementary BLUE potion always returns it to the previous color state).

## 2.2. Notation and Task Formulation:

We frame the Alchemy task as a supervised learning problem. We provide the model with a set of **support** (contextual) examples from a specific chemistry, using which the model must respond to a **query** example. We formalize the notions of **support** and **query** below.

Given a complete chemistry graph, the **support** set is a triplet consisting of K example transitions: $S = \{(x_{s_k}, z_{s_k}, y_{s_k})\}_{k=1}^{K}$. Each support example represents a transformation in the chemistry where $x_{s_k}$ is the *start stone*

*state*, $z_{s_k}$ is the *sequence of potions applied*, and $y_{s_k}$ is the resulting *end stone state*. The model then predicts the output stone ($y$) for a **query** example ($q$): $q = (x_q, z_q)$. where $x_q$ is the start stone state and $z_q$ is the query potion sequence.

To systematically vary task complexity, we introduce 'hop length' (hl), which is the length of the potion sequence $z$. For instance, a 2-hop example (hl $= 2$) corresponds to applying two potions in succession to move from the start stone state to the end stone state on the chemistry (Figure 1, right, for an example sequence). We can vary the hop length for both the support ($hl_{support}$), and the query ($hl_{query}$). We train the model to learn $p(y \mid x_q, z_q, S)$. In the other words, the model needs to correctly predict the output stone $y$, when $z_q$ is applied to $x_q$, using the entire support set $S$ as the context to infer the latent chemistry. In our experiments, we ensure that $S$ and $q$ are disjoint, and we shuffle the examples in $S$ to prevent any ordering biases.

We use this task formulation to design three systematic experiments. 1) **Latent-structure discovery under partial support**: we vary the content of $S$ with $hl_{support} = 1$ but withhold all examples for a randomly selected potion pair. The query $q$ with $hl_{query} = 1$ then tests the model on the effect of the withheld potions (Figure 1, left). (2) **Composition**: we test the model's ability to compose 1-hop support transitions ($hl_{support} = 1$), to solve multi-hop queries ($hl_{query} \in \{2, 3, 4, 5\}$) (Figure 1, top right). (3) **Decomposition**: we test the model's ability to decompose complex multi-hop support examples ($hl_{support} \in \{2, 3, 4, 5\}$) to solve a 1-hop query ($hl_{query} = 1$) (Figure 1, bottom right).

## 2.3. Training Details

We use a transformer model (Vaswani et al., 2017) with $\sim$ 2M parameters. We use 4 hidden layers, embedding dimension of 256, feedforward dimension of 512, a dropout rate of 0.1, and use causal attention.

As in recent work training transformers to solve tasks from contextual examples (Garg et al., 2022; Chan et al., 2022; Singh et al., 2023; Lake & Baroni, 2023), we provide the model with the concatenation of all support examples $S$ and the query $(x_q, z_q)$, to predict the correct output stone state $y$ (1 out of 108 possible stone states as defined in Section 2.1). We use standard cross-entropy loss on the output of the last token representation to train the model.

We use a 90 / 10 split to create the train and validation sets, ensuring that the validation chemistries are unseen during training. We train all models for 1000 epochs with a batch size of 32. We report the average validation accuracy over three seeds. For all experiments, we use a batch size of 32, AdamW optimizer (Loshchilov & Hutter, 2019), and 4 NVIDIA L40S GPUs with FP32 precision training. We tune multiple hyperparameters, including the initial learning rate,
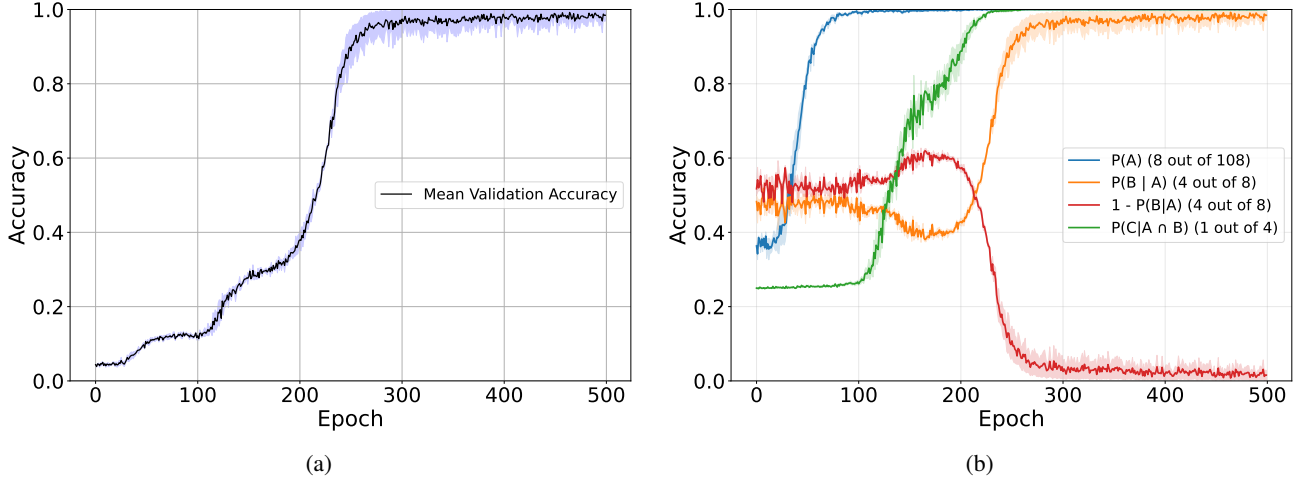
(a)

(b)

*Figure 2.* (a) Validation accuracy for latent structure learning (withheld potion pair with $hl_{support} = hl_{query} = 1$). Perfect performance demonstrates the model's ability to learn and generalize latent structures, but via distinct plateaus and jumps, indicating the acquisition of the latent structure through various stages. (b) Validation accuracy after factorizing the task into different events / components. Blue: $\mathbb{P}[A]$ (in-support). Orange: $\mathbb{P}[B|A]$ (correct half given in-support). Red: complementary incorrect half accuracy. Green: $\mathbb{P}[C|A \cap B]$ (exact match given correct half). The product of the events denoting the blue, green, and orange curves recreates the overall accuracy in Figure 2a. We only show 500 epochs because the model reached convergence. Error bars denote the standard error of the mean.

minimum learning rate, weight decay, and the learning rate scheduler. We provide the hyperparameter tuning details in Appendix A.

## 3. Staged Dynamics of Latent Structure Learning

In this experiment, we explore latent structure discovery under partial support, wherein the model must infer the effect of a withheld potion pair using the support transitions of the remaining potion pairs.

Intuitively, withholding transitions for a potion pair creates two *disconnected* halves of the chemistry graph (Figure 1, left). To correctly predict the output of a query, the model must correctly align and orient the halves relative to one another. This contrast with previous work where models might rely on surface-level patterns or counting statistics (Mittal et al., 2024). Instead, our experimental setup requires the model to learn the underlying transition rules, preventing it from exploiting superficial statistical shortcuts. We show the results in Figure 2a.

### 3.1. Analytical Study of Learning Stages

We observe multiple stages in the learning curve shown in Figure 2a. To systematically study these stages, we factorize the task to define multiple interpretable events.

**Sets and Events:** We define $T_{chem}$ as the set of eight stones in the support $S$ for any given chemistry $chem$. Because we withhold a potion pair from $T_{chem}$, we denote the target

bearing half with $H_{chem}^{(other)}$ and the non-target bearing half with $H_{chem}^{(same)}$. For these sets, we define the following events for a given chemistry $chem$ and query $q$:

$$A := \{\hat{y} \in T_{chem}\} \quad \text{(in-support)},$$
$$B := \{\hat{y} \in H_{chem}^{(other)}\} \quad \text{(correct half)},$$
$$C := \{\hat{y} = y\} \quad \text{(exact match)}.$$

These events correspond to meaningful structural properties of the task and *sub-skills* that the model must achieve to solve the task. For instance, $A$ is a meaningful event because the correct answer $(y)$ is always in the support set $(S)$. Similarly, $B$ is a meaningful event because the application of the withheld potion pair always moves stones between halves, so identifying the correct half is a necessary intermediate step to learn the task[2].

By the chain rule, we obtain the total accuracy as the product of the conditional probabilities:

$$\boxed{\mathbb{P}[C] = \mathbb{P}[A] \, \mathbb{P}[B \,|\, A] \, \mathbb{P}[C \,|\, A \cap B]}$$

Note that according to the definition of $T_{chem}$ and $H_{chem}^{(other)}$, when $\hat{y} = y$, $C \subset B \subset A$ always holds. We plot these components in Figure 2b. Note that the multiplication of all the events reproduces the overall accuracy in Figure 2a. For completeness, we also show $\mathbb{P}[A \cap B^c \,|\, A] = 1 - \mathbb{P}[B \,|\, A]$ (red) - predicting the same (incorrect) half $H_{chem}^{(same)}$.

---

[2]Our goal is not to outline every possible stage, but rather study how the stages are learned. While interesting, outlining every possible stage is intractable.

We now describe the stages in Figure 2b.

1. *Stage I: In-support (epochs ≈20–100).* $\mathbb{P}[A]$ (blue) quickly rises achieving perfect accuracy, indicating that the model quickly learns to restrict its predictions to the valid support set. However, it cannot yet distinguish between the halves or the exact query targets, indicated by the respective chance accuracies for $\mathbb{P}[B \mid A]$ and $\mathbb{P}[C \mid A \cap B]$.

2. *Stage II: Exact match given the correct half (epochs ≈100–200).* $\mathbb{P}[C \mid A \cap B]$ (red) is learned quickly indicating that given the correct half ($\mathbb{P}[B \mid A]$), the model can identify the correct target for the query $q$. Yet, the model does not learn $\mathbb{P}[B \mid A]$, which is at chance (∼ 50% accuracy).

3. *Stage III: Correct half given in-support (epochs ≈200–300).* The model completes learning $\mathbb{P}[B \mid A]$ resulting in the complete learning of the task.

We highlight a counterintuitive observation, where the exact match accuracy ($\mathbb{P}[C \mid A \cap B]$) rises before the correct half accuracy ($\mathbb{P}[B \mid A]$). This effect is caused by an adjacency bias, where the model exploits correlations in reward features to identify the set of stones adjacent to the query start stone $x_q$. Since $hl_{support} = 1$, there are three stones adjacent to $x_q$. For instance, in Figure 1 (left), if $x_q = S8$ is the start query stone with a +15 reward feature value, the adjacent stones always have a reward feature of +1 (in both same and other halves). Similarly, for $x_q = S3$, the adjacent stones always have a reward feature value of either +1 or -3. This correlation allows the model to identify adjacent neighbors without learning the full latent graph structure.

The transient reduction in the correct half performance ($H_{chem}^{(other)}$ - orange curve), epochs ≈ 120-200, confirms this effect. Crucially, two out of three stones adjacent to any $x_q$ reside in the same half (the incorrect half), whereas the true target is located in the other half (the correct half). Therefore, as the model learns to predict the adjacent stones based on the presence of the 1-hop transitions in the support set $S$, it becomes biased towards the same half, causing the observed dip in the correct half accuracy. In stage III, the model correctly learns the correct half sub-goal, mastering the full task. We provide a detailed empirical analysis of this adjacency effect in Appendix B.

## 4. Robustness of Composition

We now discuss the results for the composition task, where the support set $S$ has all possible 1-hop transitions, and the model needs to solve queries with multi-hop potion sequences. Here $\texttt{hl}_{support} = 1$, and $\texttt{hl}_{query} \in \{2, 3, 4, 5\}$. We show the results in Figure 3.
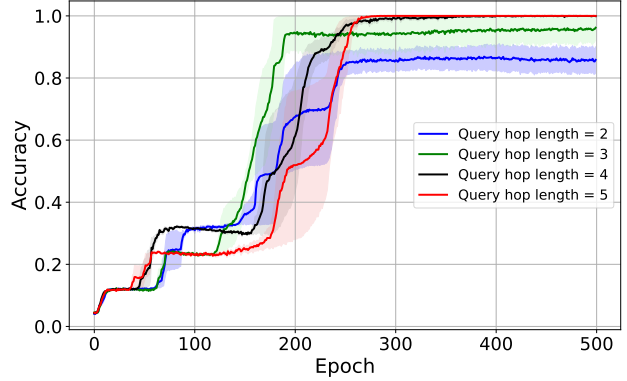


*Figure 3.* Model performance for composing 1-hop transitions to solve multi-hop queries. There is no noticeable difference in model performance with increasing in the value of $hl_{query} \in \{2, 3, 4, 5\}$. The x-axis denotes epochs, and the y-axis denotes the validation accuracy. Performance is averaged over three seeds. The error bars denote the standard error of the mean. We only show the first 500 epochs as all hops reached near-perfect accuracy.

We observe that while there is no noticeable difference in the overall learning dynamics and performance of the model for each $hl_{query} \in \{2, 3, 4, 5\}$, there is still staged learning. Furthermore, depending on the value of $hl_{query}$, the stages and plateaus occur at different accuracy values. For example for $hl_{query} \in \{2, 4\}$, the stage at ≈ 80-150 epochs occurs at ≈ 37.5% accuracy, while for $hl_{query} \in \{3, 5\}$, the same stage occurs at ≈ 25% accuracy. Next, we study these stages.

### 4.1. Analytical study of composition learning stages

To study the stages, we define three meaningful events that naturally follow from the composition task construction.

Events / metrics, for a given chemistry $chem$ and a query $q$, we define:

$$
\begin{aligned}
A &:= \{\hat{y} \in T_{chem}\} &&\text{(in-support)}, \\
R &:= \{\hat{y} \in R_k(x_q)\} &&\text{(reachable stones)}, \\
C &:= \{\hat{y} = y\} &&\text{(exact match)}.
\end{aligned}
$$

Event $R_k(x_q)$ defines the set of stones that can be "reached" by composing k-hop transitions for a given query stone $x_q$. For instance, in Figure 1 (middle), for $x_q = S1$, and $k = hl_{query} = 2$, the target $y$ can only be one of $R_2(x_q) = \{S3, S8, S5\}$.

Because by construction $y \in R_k(x_q) \subset T_{chem}$, we have $C \subset R \subset A$, and thus

$$\boxed{\mathbb{P}[C] = \mathbb{P}[A] \, \mathbb{P}[R \mid A] \, \mathbb{P}[C \mid A \cap R]}$$
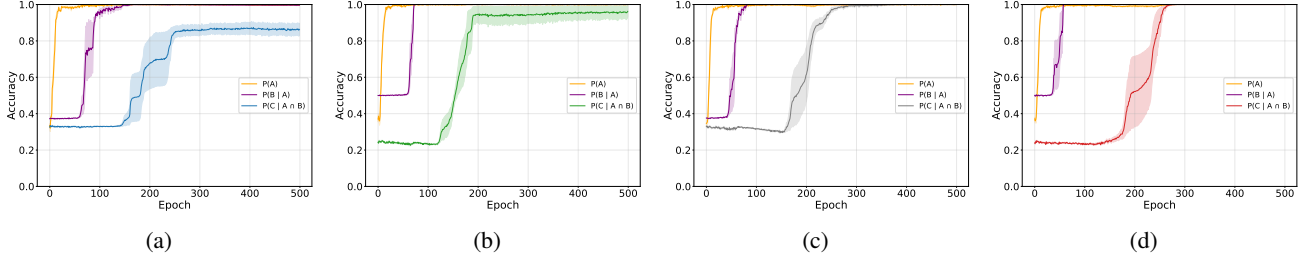
*Figure 4.* Staged learning dynamics for the composition task (first 500 epochs) with $hl_{query} \in \{2, 3, 4, 5\}$. (a) Staged dynamics for 2-hop composition. (b) Staged dynamics for 3-hop composition. (c) Staged dynamics for 4-hop composition. (d) Staged dynamics for 5-hop composition. In each subfigure, "orange" plots $\mathbb{P}[A]$ (in-support), "purple" plots $\mathbb{P}[R \mid A]$ (within reachable stones given in-support), "colored" plots (for each $k = hl_{query}$) $\mathbb{P}[C \mid A \cap R]$ (exact match given the reachable stones); their product recreates the overall accuracy in Figure 3. The x-axis denotes epochs and the y-axis denotes validation accuracy. Error bars show the standard error of the mean.

We now describe the stages in Figure 4.

1. *Stage I: In-support (epochs < ≈ 30 ).* $\mathbb{P}[A]$ rises rapidly to perfect accuracy, correctly restricting the predictions to the eight support stones. This stage mirrors "in-support" behavior in Section 3.1.

2. *Stage II: within reachable stones given in-support (epochs ≈ 30-90).* $\mathbb{P}[R|A]$ increases from near-chance levels to perfect accuracy as the model learns which subset of the eight support stones is reachable from the query start stone $x_q$ via $k$ hops. We note the different chance levels for each hop length. For $k \in \{2, 4\}$, the reachable set contains three stones, setting chance to $3/8 = 37.5\%$. For $k \in \{3, 5\}$, the reachable set contains four stones, and thus chance is $4/8 = 50\%$. These chance levels are reflected in the purple curves, in each subfigure of Figure 4. This stage is unique to composition, reflecting the need to identify the reachable stones before disambiguating among them.

3. *Stage III: Exact match given reachable stones.* $\mathbb{P}[C \mid A \cap R]$ reaches near perfect accuracy as the model learns to correctly disambiguate the target among the reachable set for a given $x_q$ and $z_q$.

Our analysis of composition and the respective stages indicates that the query hop-length has very little effect on the speed of learning. This suggests that given all 1-hop support transitions, the model can systematically compose them to solve multi-hop queries, at least for $hl_{query} \in \{2, 3, 4, 5\}$.

The sequential nature of the stages supports a *coarse-to-fine* learning hypothesis, where the model first learns a set of coarse solutions (or partial solutions) before learning the fine-grained solutions (understanding the exact target). This finding aligns with the recent work studying emergent behavior in neural transformer models (e.g,. Gopalani et al. (2024); Nanda et al. (2023)).

## 5. Decomposition under Increasing Task Complexity

We now establish the difference in model performance when systematically varying $hl_{support} \in \{2, 3, 4, 5\}$ to solve 1-hop queries ($hl_{query} = 1$). Increasing $hl_{support}$ tests the model's ability to infer the underlying chemistry from multi-step potion sequences to solve a 1-hop query. For a given $hl_{support}$ value, we provide all transitions to the model.

The results for decomposition appear in Figure 5 and show that increasing $hl_{support}$ delays model convergence, with $hl_{support} = 2$ reaching convergence sooner than $hl_{support} = 5$. This indicates that increasing the support hop length delays the model's ability to learn and generalize latent structures from contextual examples. This finding echos recent works on understanding language model performance when solving tasks with increasing complexity (Dziri et al., 2023; Shojaee et al., 2025), and extends it to a task with a well-defined latent structure.

### 5.1. Analytical study of decomposition learning stages

Figure 5 shows multiple stages for all $hl_{support} \in \{2, 3, 4, 5\}$, which we investigate next.

We first define the sets. Let $T_{chem}$ denote the set of eight stones in the support $S$ for any given chemistry $chem$. For a given potion $z_q$, the target $y$ is in a specific half of chemistry $H_{chem}^{(1)}$, determined by the axis of $z_q$. For example, in Figure 1 (middle), applying a YELLOW 1-hop potion on the stone $x_q = S1$, can only result in a stone from $\{S6, S5, S4, S3\}$ (because YELLOW can only go to these four stones).

For these sets, we define the following events (for query $q$):

$$A := \{\hat{y} \in T_{chem}\} \qquad \text{(in-support)},$$
$$B := \{\hat{y} \in H_{chem}^{(other)}\} \qquad \text{(correct half)},$$
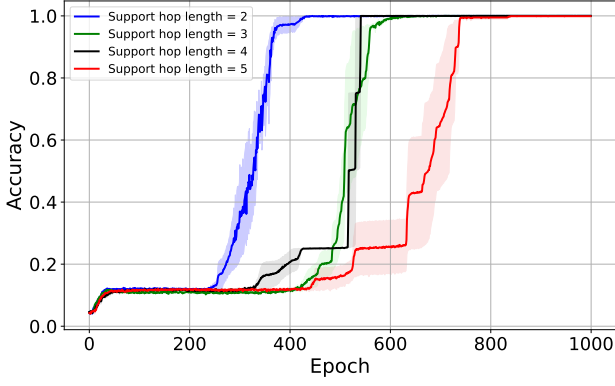$$C := \{\hat{y} = y\} \qquad \text{(exact match)}.$$

*Figure 5.* Decomposition results for various $hl_{support}$ $\in$ $\{2, 3, 4, 5\}$. We see delayed convergence with increasing task complexity ($hl_{support}$), where the 2-hop converges the earliest and the 5-hop converges the latest. The x-axis denotes epochs, and the y-axis denotes the validation accuracy. Performance is averaged over three seeds, and the error bars show the standard error of the mean.

the multiplication of which gives the exact match accuracy.

$$\mathbb{P}[C] = \mathbb{P}[A]\,\mathbb{P}[B \mid A]\,\mathbb{P}[C \mid A \cap B]$$

We now describe the stages in Figure 6.

1. *Stage I: In-support (epochs $<\approx$ 20).* Across $k \in \{2, 3, 4, 5\}$, the model quickly learns to restrict predictions to the support set, as shown by the rapid rise of $\mathbb{P}_k[A]$ to near perfect accuracy.

2. *Stage II: Correct half given in-support.* For all $hl_{support} \in \{2, 3, 4, 5\}$, $\mathbb{P}[B \mid A]$ remains near chance ($\approx 50\%$) for multiple epochs, indicating that the model cannot learn the correct half immediately after learning $\mathbb{P}[A]$. This corresponds to the long stagnation segment at 12.5% in Figure 5. Interestingly, increasing $hl_{support}$ delays learning the correct half as shown by the delayed rise of the purple curve.

3. *Stage III: Exact match given the correct half.* For all $hl_{support} \in \{2, 3, 4, 5\}$, we observe that after learning the correct half, the model does not show a significant delay in learning the final exact match accuracy, i.e., the model can correctly distinguish the target from the set of the four stones. This learning is reflected by the rapid rise of the final event $\mathbb{P}[C \mid A \cap B]$ from chance accuracy of 25% to perfect accuracy.

In summary, increasing $hl_{support}$ primarily delays the learning of the correct half accuracy, which ultimately delays model convergence. Nonetheless, similar to composition task (Section 4.1). The stages are learned in a coarse-to-fine

sequence where the model first learns to restrict the prediction to the set stones in the support ($\mathbb{P}[A]$), then learns to predict the correct half where the target belongs to ($\mathbb{P}[B \mid A]$), and finally learns to correctly disambiguate within the correct half ($\mathbb{P}[C \mid A \cap B]$).

We undertook a deeper analysis to understand why learning the correct half stage ($\mathbb{P}[B \mid A]$) is delayed with increasing $hl_{support}$. We examined the presence of intermediate stages, by defining relevant metrics for events, such as initially narrowing down the search space in the support set $S$ to a subset of candidate stones before learning the correct half. However, the simultaneous rise of the correct half accuracy and the relevant metrics indicated the absence of intermediate stages. We provide a detailed analysis of this investigation in Appendix D.

**Sensitivity of Decomposition Learning Dynamics**: While our results showed that increasing task complexity delays the model's ability to learn the latent structure, we further observed that this ability is highly sensitive to the hyperparameter configuration. We identify failure modes, where choosing sub-optimal hyperparameters can not only alter the staged learning dynamics, but also prevent the model from converging within a compute budget. Consequently, sub-optimal hyperparameters can cause the model to stagnate for an extended period without learning some components of the latent structure. We provide a case study of hyperparameter sensitivity in Appendix C.

# 6. Related Work

## 6.1. Latent Structure Learning

Transformers can learn to encode graph like latent structures during training, going beyond memorization. For example, a two-layer transformer can encode the latent causal graph in the first attention layer (Nichani et al., 2024), and attention weights can be effectively seen as equivalent to edges of a graph structure (Henderson et al., 2023). Transformer based language models such as BERT (Devlin et al., 2019) can also encode linguistic structure in the attention heads (Manning et al., 2020; Ravishankar et al., 2021). However, analyzing the progression of this learning in natural language setting is difficult due to the complexities of language. We address this gap by using a controlled benchmark to track the learning of specific aspects of latent structure.

## 6.2. Staged Dynamics

Previous works on understanding the learning dynamics of transformers (Barak et al., 2022; Chen et al., 2024; Gopalani et al., 2024; Chen et al., 2025; Gopalani & Hu, 2025; Song et al., 2025; Zhang et al., 2025), have studied the connection between different phases in the loss curves and learning different components of the task, often attributed to the
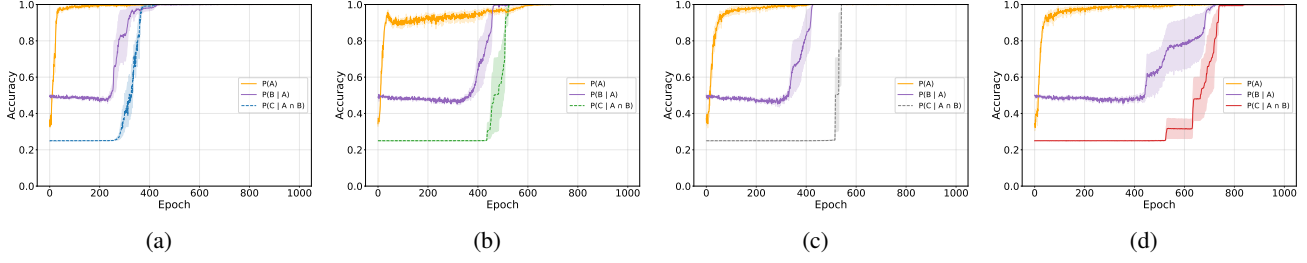
*Figure 6.* Staged learning dynamics for the decomposition experiments with $hl_{support} \in \{2, 3, 4, 5\}$ (a) Staged learning dynamics for 2-hop decomposition, (b) Staged learning dynamics 3-hop decomposition. (c) Staged learning dynamics for 4-hop decomposition. (d) Staged learning dynamics for 5-hop decomposition. In each subfigure, "orange" plots $\mathbb{P}_k[A]$ (in-support), "purple" plots $\mathbb{P}[B \mid A]$ (correct half given in-support), and the "colored" plots $\mathbb{P}_k[C \mid A \cap B]$ (exact match given correct half); their product recreates the overall accuracy in Figure 5. We show the model performance for 1000 epochs. Error bars denote the standard error of the mean across three seeds.

emergence of "induction heads" (Reddy, 2023; Singh et al., 2024; Olsson et al., 2022). A closely related work by Dai et al. (2025) investigates graph reasoning through circuit analysis (Ameisen et al., 2025), but does not study the training dynamics of a task with a well-defined latent structure. In contrast, we analyze the training dynamics of such a task, demonstrating that plateaus and stages in accuracy correspond to the acquisition of interpretable sub-skills.

The "phase change" behavior is also related to grokking (Power et al., 2022), where a system suddenly achieves generalization after overfitting. In contrast to grokking, we do not explore the case of overfitting on the training data.

### 6.3. Compositionality and Decompositionality

Compositionality is the test of a model's ability to combine components in novel ways, and remains a fundamental challenge in machine learning. Existing works also highlight the shortcomings of transformers in compositional tasks with increasing complexity (Hupkes et al., 2020; Dziri et al., 2023; Petty et al., 2025), and in tasks requiring compositional generalization (Ontañón et al., 2022; Kobayashi et al., 2024).

Other studies have investigated the ability of transformers to solve compositionality when they are trained from scratch (Lake & Baroni, 2023; Thomm et al., 2024), but they focus on the final learned model, and on tasks with an unclear latent structure or graph like structure, making it difficult to understand the learning behavior in such settings. To the best of our knowledge, our work is the first to show how a model learns composition in a task with a well-defined latent structure.

Decomposition is the ability of language models to break down complex tasks into simpler components. Existing work has studied the decomposition ability of large language models through the lens of prompting (Khot et al., 2023; Pasewark et al., 2024; Wei et al., 2023; Prasad et al., 2024) to improve model performance. However, the dynamics

of how decomposition capabilities emerge during training remains unclear, which we address in the current work.

## 7. Conclusion

We studied how a small decoder-style transformer model learns a task with a well-defined latent structure, using the Alchemy benchmark. Using a simple event-based factorization, we showed that the model goes through different learning *stages*, corresponding to different aspects of the latent structure. Specifically, the sequence of stages indicates that the model learns the latent structure by first mastering the subset of possible targets, before refining its predictions to the exact target. In the latent structure discovery under partial support task, we identified the model's ability to leverage the adjacency bias in the latent structure, which resulted in learning the correct target before the correct half. Through the design of composition and decomposition tasks, we tested the effect of increasing task complexity and found a fundamental asymmetry in performance. In composition, the model exhibited invariance to task complexity, robustly solving simple rules without noticeable performance degradation. However, in the decomposition task, the model suffered from a complexity bottleneck, where increasing task complexity delayed model convergence. Our principled analysis pinpointed this failure to the delay in learning one of the stages, causing stagnation for multiple epochs.

Future work should investigate whether pre-trained Large Language Models (LLMs) exhibit similar staged dynamics during pre-training or when fine-tuning on structured tasks. Investigating alternative training strategies to mitigate the complexity is also a promising research direction. Furthermore, designing more datasets and benchmarks allowing for a careful study of latent structure learning in transformer based models remains an open problem. Finally, extending this analysis to sophisticated architectures, such as graph neural networks, remains a promising direction for understanding the architectural limits of latent structure learning.

## Impact Statement

We contributed towards a better understanding of how a transformer learns a structured task. We showed that a decoder-style transformer goes through various stages of acquiring latent structures, which we believe will improve our understanding of how their capabilities emerge during training.

## References

Ameisen, E., Lindsey, J., Pearce, A., Gurnee, W., Turner, N. L., Chen, B., Citro, C., Abrahams, D., Carter, S., Hosmer, B., Marcus, J., Sklar, M., Templeton, A., Bricken, T., McDougall, C., Cunningham, H., Henighan, T., Jermyn, A., Jones, A., Persic, A., Qi, Z., Ben Thompson, T., Zimmerman, S., Rivoire, K., Conerly, T., Olah, C., and Batson, J. Circuit tracing: Revealing computational graphs in language models. *Transformer Circuits Thread*, 2025. URL https://transformer-circuits.pub/2025/attribution-graphs/methods.html.

Barak, B., Edelman, B., Goel, S., Kakade, S., Malach, E., and Zhang, C. Hidden Progress in Deep Learning: SGD Learns Parities Near the Computational Limit. *Advances in Neural Information Processing Systems*, 35:21750–21764, December 2022.

Biewald, L. Experiment tracking with weights and biases, 2020. URL https://www.wandb.com/. Software available from wandb.com.

Chan, S. C. Y., Santoro, A., Lampinen, A. K., Wang, J. X., Singh, A., Richemond, P. H., McClelland, J., and Hill, F. Data Distributional Properties Drive Emergent In-Context Learning in Transformers, November 2022.

Chen, A., Shwartz-Ziv, R., Cho, K., Leavitt, M. L., and Saphra, N. Sudden Drops in the Loss: Syntax Acquisition, Phase Transitions, and Simplicity Bias in MLMs, March 2025.

Chen, S., Sheen, H., Wang, T., and Yang, Z. Training Dynamics of Multi-Head Softmax Attention for In-Context Learning: Emergence, Convergence, and Optimality, June 2024.

Dai, X., Lo, C.-H., Guo, K., Zeng, S., Luo, D., and Tang, J. Uncovering Graph Reasoning in Decoder-only Transformers with Circuit Tracing, September 2025.

Devlin, J., Chang, M.-W., Lee, K., and Toutanova, K. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding, May 2019.

Dziri, N., Lu, X., Sclar, M., Li, X. L., Jiang, L., Lin, B. Y., West, P., Bhagavatula, C., Bras, R. L., Hwang, J. D., Sanyal, S., Welleck, S., Ren, X., Ettinger, A., Harchaoui, Z., and Choi, Y. Faith and Fate: Limits of Transformers on Compositionality, 2023.

Garg, S., Tsipras, D., Liang, P. S., and Valiant, G. What can transformers learn in-context? A case study of simple function classes. In Koyejo, S., Mohamed, S., Agarwal, A., Belgrave, D., Cho, K., and Oh, A. (eds.), *Advances in Neural Information Processing Systems*, volume 35, pp. 30583–30598. Curran Associates, Inc., 2022.

Gopalani, P. and Hu, W. What Happens During the Loss Plateau? Understanding Abrupt Learning in Transformers, October 2025.

Gopalani, P., Lubana, E. S., and Hu, W. Abrupt learning in transformers: A case study on matrix completion. *Advances in Neural Information Processing Systems*, 37:55053–55085, 2024.

Henderson, J., Mohammadshahi, A., Coman, A., and Miculicich, L. Transformers as Graph-to-Graph Models. In Elazar, Y., Ettinger, A., Kassner, N., Ruder, S., and A. Smith, N. (eds.), *Proceedings of the Big Picture Workshop*, pp. 93–107, Singapore, December 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.bigpicture-1.8.

Hewitt, J. and Manning, C. D. A Structural Probe for Finding Syntax in Word Representations. In Burstein, J., Doran, C., and Solorio, T. (eds.), *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pp. 4129–4138, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1419.

Hupkes, D., Dankers, V., Mul, M., and Bruni, E. Compositionality Decomposed: How do Neural Networks Generalise? *Journal of Artificial Intelligence Research*, 67:757–795, April 2020. ISSN 1076-9757. doi: 10.1613/jair.1.11674.

Khot, T., Trivedi, H., Finlayson, M., Fu, Y., Richardson, K., Clark, P., and Sabharwal, A. Decomposed Prompting: A Modular Approach for Solving Complex Tasks, April 2023.

Kobayashi, S., Schug, S., Akram, Y., Redhardt, F., von Oswald, J., Pascanu, R., Lajoie, G., and Sacramento, J. When can transformers compositionally generalize in-context?, July 2024.

Lake, B. M. and Baroni, M. Human-like systematic generalization through a meta-learning neural network. *Nature*, 623(7985):115–121, November 2023. ISSN 0028-0836, 1476-4687. doi: 10.1038/s41586-023-06668-3.

Loshchilov, I. and Hutter, F. Decoupled weight decay regularization. In *International Conference on Learning Representations*, 2019. URL https://openreview.net/forum?id=Bkg6RiCqY7.

Manning, C. D., Clark, K., Hewitt, J., Khandelwal, U., and Levy, O. Emergent linguistic structure in artificial neural networks trained by self-supervision. *Proceedings of the National Academy of Sciences*, 117(48):30046–30054, December 2020. ISSN 0027-8424, 1091-6490. doi: 10.1073/pnas.1907367117.

Mittal, S., Elmoznino, E., Gagnon, L., Bhardwaj, S., Sridhar, D., and Lajoie, G. Does learning the right latent variables necessarily improve in-context learning?, May 2024.

Nanda, N., Chan, L., Lieberum, T., Smith, J., and Steinhardt, J. Progress measures for grokking via mechanistic interpretability, October 2023.

Nichani, E., Damian, A., and Lee, J. D. How Transformers Learn Causal Structure with Gradient Descent, 2024.

Olsson, C., Elhage, N., Nanda, N., Joseph, N., DasSarma, N., Henighan, T., Mann, B., Askell, A., Bai, Y., Chen, A., Conerly, T., Drain, D., Ganguli, D., Hatfield-Dodds, Z., Hernandez, D., Johnston, S., Jones, A., Kernion, J., Lovitt, L., Ndousse, K., Amodei, D., Brown, T., Clark, J., Kaplan, J., McCandlish, S., and Olah, C. In-context Learning and Induction Heads, September 2022.

Ontañón, S., Ainslie, J., Cvicek, V., and Fisher, Z. Making Transformers Solve Compositional Tasks, March 2022.

Pasewark, E., Montgomery, K., Duan, K., Song, D., and Wang, C. Re-Tuning: Overcoming the Compositionality Limits of Large Language Models with Recursive Tuning. In Ku, L.-W., Martins, A., and Srikumar, V. (eds.), *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 10422–10437, Bangkok, Thailand, August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.acl-long.561.

Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., Desmaison, A., Köpf, A., Yang, E., DeVito, Z., Raison, M., Tejani, A., Chilamkurthy, S., Steiner, B., Fang, L., Bai, J., and Chintala, S. Pytorch: An imperative style, high-performance deep learning library, 2019. URL https://arxiv.org/abs/1912.01703.

Petty, J., Hu, M. Y., Wang, W., Ravfogel, S., Merrill, W., and Linzen, T. RELIC: Evaluating Compositional Instruction Following via Language Recognition, 2025.

Power, A., Burda, Y., Edwards, H., Babuschkin, I., and Misra, V. Grokking: Generalization beyond overfitting on small algorithmic datasets, 2022. URL https://arxiv.org/abs/2201.02177.

Prasad, A., Koller, A., Hartmann, M., Clark, P., Sabharwal, A., Bansal, M., and Khot, T. ADaPT: As-Needed Decomposition and Planning with Language Models, April 2024.

Ravishankar, V., Kulmizev, A., Abdou, M., Søgaard, A., and Nivre, J. Attention Can Reflect Syntactic Structure (If You Let It). In Merlo, P., Tiedemann, J., and Tsarfaty, R. (eds.), *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pp. 3031–3045, Online, April 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.eacl-main.264.

Reddy, G. The mechanistic basis of data dependence and abrupt learning in an in-context classification task, December 2023.

Shojaee, P., Mirzadeh, I., Alizadeh, K., Horton, M., Bengio, S., and Farajtabar, M. The Illusion of Thinking: Understanding the Strengths and Limitations of Reasoning Models via the Lens of Problem Complexity, June 2025.

Singh, A. K., Chan, S. C. Y., Moskovitz, T., Grant, E., Saxe, A. M., and Hill, F. The Transient Nature of Emergent In-Context Learning in Transformers, December 2023.

Singh, A. K., Moskovitz, T., Hill, F., Chan, S. C. Y., and Saxe, A. M. What needs to go right for an induction head? A mechanistic study of in-context learning circuits and their formation, April 2024.

Song, J., Xu, Z., and Zhong, Y. Out-of-distribution generalization via composition: A lens through induction heads in Transformers. *Proceedings of the National Academy of Sciences*, 122(6):e2417182122, February 2025. doi: 10.1073/pnas.2417182122.

Thomm, J., Camposampiero, G., Terzic, A., Hersche, M., Schölkopf, B., and Rahimi, A. Limits of Transformer Language Models on Learning to Compose Algorithms, November 2024.

Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ukasz Kaiser, Ł., and Polosukhin, I. Attention is All you Need. In *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017.

Wang, J., King, M., Porcel, N., Kurth-Nelson, Z., Zhu, T., Deck, C., Choy, P., Cassin, M., Reynolds, M., Song, F., Buttimore, G., Reichert, D., Rabinowitz, N., Matthey, L., Hassabis, D., Lerchner, A., and Botvinick, M. Alchemy:

A structured task distribution for meta-reinforcement learning. *arXiv preprint arXiv:2102.02926*, 2021. URL https://arxiv.org/abs/2102.02926.

Wei, J., Tay, Y., Bommasani, R., Raffel, C., Zoph, B., Borgeaud, S., Yogatama, D., Bosma, M., Zhou, D., Metzler, D., Chi, E. H., Hashimoto, T., Vinyals, O., Liang, P., Dean, J., and Fedus, W. Emergent Abilities of Large Language Models, October 2022.

Wei, J., Wang, X., Schuurmans, D., Bosma, M., Ichter, B., Xia, F., Chi, E., Le, Q., and Zhou, D. Chain-of-thought prompting elicits reasoning in large language models, 2023. URL https://arxiv.org/abs/2201.11903.

Zhang, Y., Singh, A. K., Latham, P. E., and Saxe, A. Training Dynamics of In-Context Learning in Linear Attention, May 2025.

## A. Training details

We implement all model using PyTorch (Paszke et al., 2019) version 2.4. For hyperparameter tuning, we conduct a large sweep using Weights & Biases (Biewald, 2020) over multiple hyperparameters with the following values for three random seeds. We show the hyperparameter configuration in Table 1.

| Hyperparameter | Values |
| --- | --- |
| learning rate | [1e-3, 4e-4, 5e-4, 1e-4, 9e-5, 7e-5, 1e-5] |
| scheduler | cosine, cosine w/ restarts, multi step, no scheduler |
| weight decay | [0.1, 0.01, 0.001] |
| minimum learning rate (cosine and cosine w/ restarts) | 7e-5, 8e-5, 9e-5, 0.000085, 0.000095, 1e-5 |
| multi step gamma | [0.2, 0.4, 0.5, 0.6, 0.7] |

*Table 1.* Set of hyperparameters sweeped over for the experiments.

## B. Reward-Based Analysis of Staged Latent-Structure Discovery with Partial Support

To investigate the effect of adjacency bias (Section 3) on the learning curve of $\mathbb{P}[C \mid A \cap B]$, we plot the same events for each individual $x_{q_r}$ with a distinct reward feature value, $r \in \{+15, +1, -1, -3\}$. One underlying reward-related structure in Alchemy is that there always exists exactly one stone state with reward $+15$ and $-3$ per chemistry, and three stone states with reward $+1$ and $-1$ per chemistry, which in total adds up to $1 + 1 + 3 + 3 = 8$ total vertices of a cube.

From the definition of the events $A$, $B$, and $C$, we can infer that $C \subset B \subset A$, and thus $\mathbb{P}[C \mid A \cap B]$ can be equivalently written as $\mathbb{P}[C \mid B]$. Figure 7a demonstrates the averages $\mathbb{P}_{+15}[C \mid B]$ (brown) and $\mathbb{P}_{-3}[C \mid B]$ (pear) over all chemistries, as well as $\mathbb{P}_{+1}[C \mid B]$ (pink) and $\mathbb{P}_{-1}[C \mid B]$ (blue), which are averaged over all chemistries and an additional average over all three stone states $x_{q_r}$ with the reward feature $r \in \{+1, -1\}$, per chemistry. We can see that $\mathbb{P}_{r \in \{+15, -3\}}[C \mid B]$ is learned faster than $\mathbb{P}_{r \in \{+1, -1\}}[C \mid B]$. Similarly, Figure 7b demonstrates the correct target prediction for the in-support predicted stones, i.e. $\mathbb{P}_r[C \mid A]$, for all query start stone states $x_{q_r}$ with a distinct reward value, $r \in \{+15, +1, -1, -3\}$. The graphs show faster rise of the learning curve $\mathbb{P}_{r \in \{+15, -3\}}[C \mid A]$ comparing to the $\mathbb{P}_{r \in \{-1, +1\}}[C \mid A]$. Through the remainder of this section, we explain these observations further by exploring the adjacency reward bias that exists in the underlying Alchemy latent structure.

### B.1. Investigating the Staged Learning of Latent Structure Based on Reward-Adjacency

According to the Alchemy underlying latent structure, the reward feature can only increase or decrease by one step per potion transition. For example, $+15$ and $-3$ only transits to stone states $+1$ and $-1$, respectively, while stones with rewards $+1$ and $-1$ can move towards stones with a higher step or lower step rewards of $\{+15, -1\}$ and $\{+1, -3\}$, respectively.

Let the set of possible stones that can be reached with a 1-hop transition from $x_{q_{+15}}$ be called $\mathcal{T}_{+15}$, consisting of the three stone states with a reward value of $r = +1$. Similarly for $x_{q_{-3}}$, $\mathcal{T}_{-3}$ always consists of the three stone states with $r = -1$. Each stone state with reward $r = +1$ is always mapped to either the three stone states with reward $r = -1$ or the stone state with reward $r = +15$, resulting in a $|\mathcal{T}_{+1}| = 4$. Similarly, each of the three stones with reward $-1$ can only map to either the three stones with reward $+1$ or the stone with reward $-3$, resulting in $|\mathcal{T}_{-1}| = 4$. Suppose the model learns this relatively easy underlying latent structure as a distinguishing signal for predicting the correct half properly, given the start stone state in a query, $x_{q_r}$. In that case, we should see an almost uniform distribution over $\hat{y} \in \mathcal{T}_r$ at the beginning of learning. Mathematically, we plot the graphs for $\mathbb{P}_r[\hat{y} = y \mid y \in \mathcal{T}_r]$ for each $x_{q_r}$ for all values of $r$. This is shown in Figure 8 in light blue, where we can see that the initial values after one epoch are still almost $\frac{1}{\mathcal{T}_{r \in \{-3, +15\}}} \approx \frac{1}{3} \approx 0.33$ and $\frac{1}{\mathcal{T}_{r \in \{-1, +1\}}} \approx \frac{1}{4} = 0.25$.
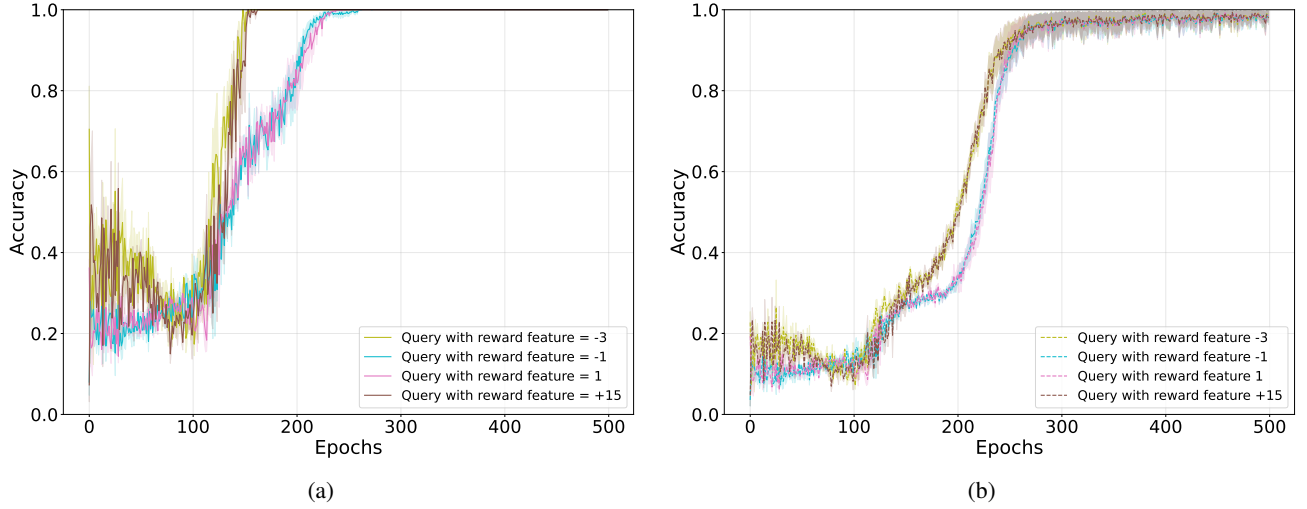
(a)                                    (b)

*Figure 7.* Prediction accuracy grouped by the reward value of the query. (a) Exact match accuracy given correct half. (b) Exact match accuracy given in-support. We observe that the model learns the exact match accuracy for queries with -3, and +15 reward features faster than those with -1, and +1 reward features.

We also show how the model learns the correct set of $\mathcal{T}_r$ (within reward adjacent metric in Figure 8) as the reward-based query prediction candidate set, given the 8 stone states seen in the support of each sample, $\mathbb{P}_r[\hat{y}_{q_r} \in \mathcal{T}_r \mid \hat{y}_{q_r} \in S]$. This is shown in dark blue with an initial accuracy of $\frac{|\mathcal{T}_{r \in \{-1,+1\}}|}{|S|} = 4/8 = 0.5$, and $\frac{|\mathcal{T}_{r \in \{-3,+15\}}|}{|S|} = 3/8 = 0.375$ for rewards $\{-3, +15\}$. Note that for stones with reward $\{-3, +15\}$, the model is initially heavily biased towards predicting $\mathcal{T}_r$ ($\forall r \in \{-3, +15\}$) because of the presence of the respective transitions in the support $S$.

So far, we have defined the subset $\mathcal{T}_r$ to denote the stone states with plausible 1-hop neighboring reward features. Considering the cubic latent structure of chemistry, we already know that each vertex can maximally connect to three other vertices. Let the geometrically neighboring vertices subset for each stone state $x_{q_r}^i$ to be shown with $\mathcal{N}_r$, which is always size 3 due to the cubic chemistry structure. In Figure 8, we additionally plot the learning curve of geometrically neighboring (within true adjacent) vertices based on each reward feature value, which is depicted in orange. It naturally becomes clear that for the reward features $r \in \{-3, +15\}$, $\mathcal{N}_r$ and $\mathcal{T}_r$ perfectly overlap:

$$\mathbb{P}_r[\hat{y}_{q_r} \in \mathcal{N}_r \mid \hat{y}_{q_r} \in S] = \mathbb{P}_r[\hat{y}_{q_r} \in \mathcal{T}_r \mid \hat{y}_{q_r} \in S] = 3/8 = 0.375 \ , \ \forall r \in \{-3, +15\}.$$

However, the geometrical adjacency subset does not overlap with the reward-based adjacency subset of a query start stone state with rewards $-1$ and $+1$; rather, $\mathcal{N}_r \subset \mathcal{T}_r$ for $r \in \{-1, +1\}$:

$$\mathbb{P}_r[\hat{y}_{q_r} \in \mathcal{N}_r \mid \hat{y}_{q_r} \in S] = \frac{|\mathcal{N}_{r \in \{-1,+1\}}|}{|S|} = 3/8 = 0.375$$

$$\neq$$

$$\mathbb{P}_r[\hat{y}_{q_r} \in \mathcal{T}_r \mid \hat{y}_{q_r} \in S] = \frac{|\mathcal{T}_{r \in \{-1,+1\}}|}{|S|} = 4/8 = 0.5$$

We also studied the selection rate of the in-support non-target reward-based adjacent stone states (within adjacent in support), $\mathcal{R}_r$, which consists of two geometrically adjacent vertices in the same half as the $x_{q_r}$'s half.

To clarify this definition further, for $r \in -3, +15$, there are two reward-based adjacent nodes in the same half of the query's start state. For $r \in -1, +1$, the in-support reward-based adjacent states consists of three stones, which two of them are geometrically adjacent in the same half as $x_{q_r}$. In other words, these are the stones for which there are transitions from $x_{q_r}$ in the support $S$.

Therefore, we have $\mid \mathcal{R}_{r \in \{-1,+1\}} \mid = 2$, and $\mid \mathcal{R}_{r \in \{+15,-3\}} \mid = 2$.

13

The stone states in $\mathcal{R}_r$ should not be predicted by the learner when it captures the underlying latent structure fully. However, since they have relevant reward features to the true target, we can expect a temporary rise in their probability of prediction. The red curve shown in Figure 8 demonstrates the probability of model predicting $\hat{y}_{q_r} \in \mathcal{R}_r$ given the in-support stone states $S$:

$$\mathbb{P}_r[\hat{y}_{q_r} \in \mathcal{R}_r \mid \hat{y}_{q_r} \in S] = \frac{|\mathcal{R}_{r \in \{-1,+1\}}|}{|S|} = 2/8 = 0.25$$

and

$$\mathbb{P}_r[\hat{y}_{q_r} \in \mathcal{R}_r \mid \hat{y}_{q_r} \in S] = \frac{|\mathcal{R}_{r \in \{-3,+15\}}|}{|S|} = 2/8 = 0.25$$

The red learning curve begins with a rise that demonstrates the increasing reliance of the learner on reward features to make better predictions.

This explains the dip in the learning curve of correct half prediction in the latent structure discovery under partial support (Section 3) component analysis shown in Figure 2b (orange). The dip in correct half prediction happens due to the model bias on predicting the target according to the reward adjacency structure for a while right before epoch 200 (i.e. the rise in the red curve here in Figure 2b). According to the reward-based underlying latent structure, the model first learns to predict the stone-states in each adjacent reward-based subset randomly. This increases the model's performance for a while until the model predicting of the stones in the same half becomes damaging to the overall performance. So the model unlearns to predict any reward-based adjacent stones randomly, but consider other perceptual features as well. These are the epochs when we see the rise in orange curve in Figure 2b and the decline in the red curve in Figure 8.
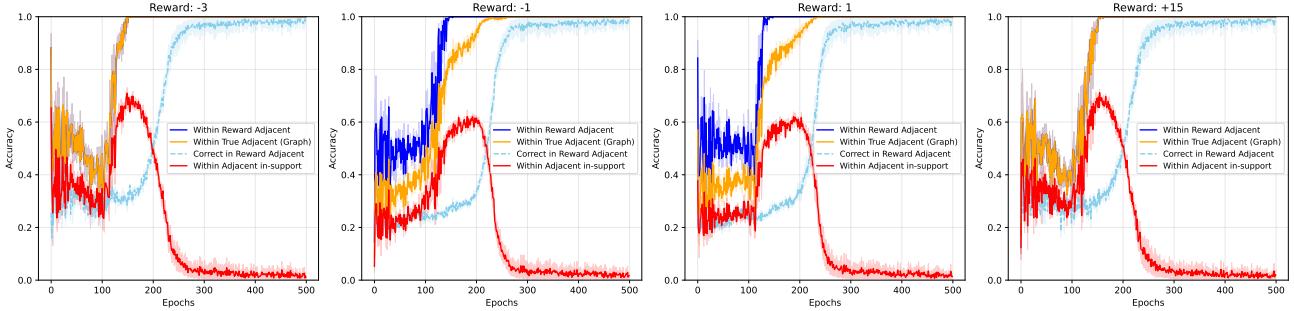


*Figure 8.* Adjacency analysis for each individual reward feature value: i) Dark blue shows the learning curve of the adjacent candidate vertices based on only the reward-related latent structure of the 1-hop transitions (correct set of $\mathcal{T}_r$); ii) orange demonstrates the geometrical adjacency ($\mathcal{N}_r$) subset of neighboring vertices, which overlaps the reward adjacency subset, when the start state has reward values of $\{-3,+15\}$; iii) light blue belongs to the learning of the correct target prediction given the reward-based adjacent subset (correct given $\mathcal{T}_r$); iv) red depicts the selection probability of in-support geometrical 1-hop vertices $\mathcal{R}_r$ (only two vertices in latent structure discovery under partial support), which are in the wrong half and should not get selected (goes to zero) as the model learns the underlying latent structure.

## C. Case Study on Decomposition Hyperparameter Sensitivity

Throughout our analysis, we observed that model performance is highly sensitive to hyperparameter configuration. In this section, we examine how sub-optimal hyperparameter choices can lead to three distinct failure modes: 1) differing optimal configurations for different task complexities, and 2) failure to learn a specific stage within a compute budget, and 3) learning a stage before mastering a previous stage. We use the decomposition experiment (Section 5) to discuss these points.

**Sensitivity of model convergence across hyperparameters** We first observe that optimal hyperparameters are not transferrable across different hop lengths. We use the weight decay value to demonstrate our point (while controlling for other hyperparameters such as learning rate and scheduler). As shown in Figure 9a, the overall accuracy for the 2-hop

decomposition task achieves optimal convergence with a weight decay value of 0.001, and using a weight decay value of 0.01 results in later convergence. On the other hand, for the 3-hop decomposition task shown in Figure 9b, the optimal weight decay value is 0.1, and using a 0.001 weight decay value does not lead to convergence within the compute budget of 1000 epochs, making the model stay in a plateau. The plateau behavior suggests that the model may learn some stages but not others. We examine the effect of hyperparameters in staged learning next.

**Suboptimal hyperparameters can prevent the observation of stages:** We demonstrate the effect of hyperparameters on learning the different components using the 3-hop decomposition task in Figure 10.

For the optimal weight decay value = 0.1, the model learns all the components of the task (Figure 10a). For another weight decay value (= 0.01), the model learns the in-support event, and the correct half event, but there's a delay in learning the exact match give the correct half event (Figure 10b). Finally, for the suboptimal weight decay value of 0.001, the model fails to learn some stages, shown by the stagnation at chance levels plateau for the entire 1000 epoch budget (Figure 10c). This highlights a critical experimental insight: a failure to converge may not indicate a lack of model capacity, but rather a choice of poor hyperparameter values that delayed the learning of a latent structure component.

All in all, these findings indicate that while the model may have the capacity to learn the latent structure, sub-optimal hyperparameters can prevent it from learning all the task components.
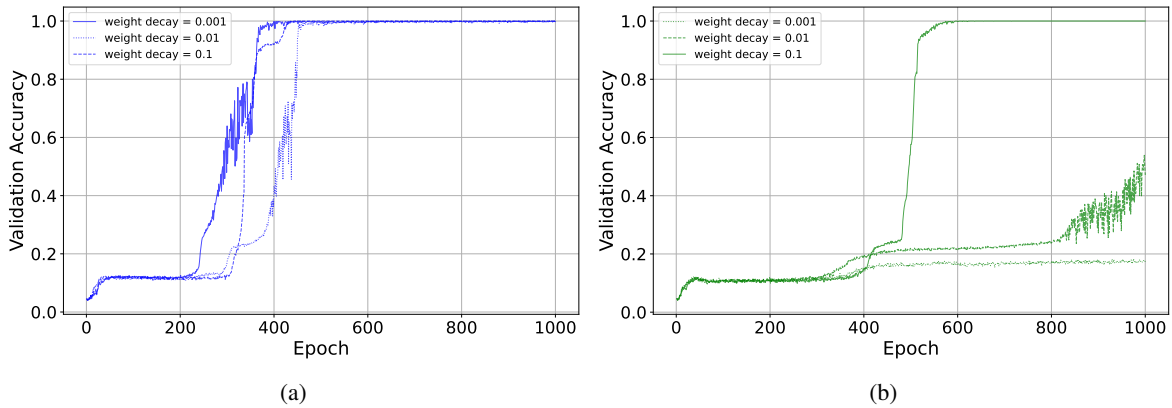


(a)                                          (b)

*Figure 9.* Examples of hyperparameter sensitivity for the 2-hop and 3-hop decomposition experiments. (a) 2-hop decomposition with different weight decay values delays convergence. (b) 3-hop decomposition experiment, where suboptimal weight decay values might cause the model to get stuck in a plateau for extended periods during training, delaying convergence.
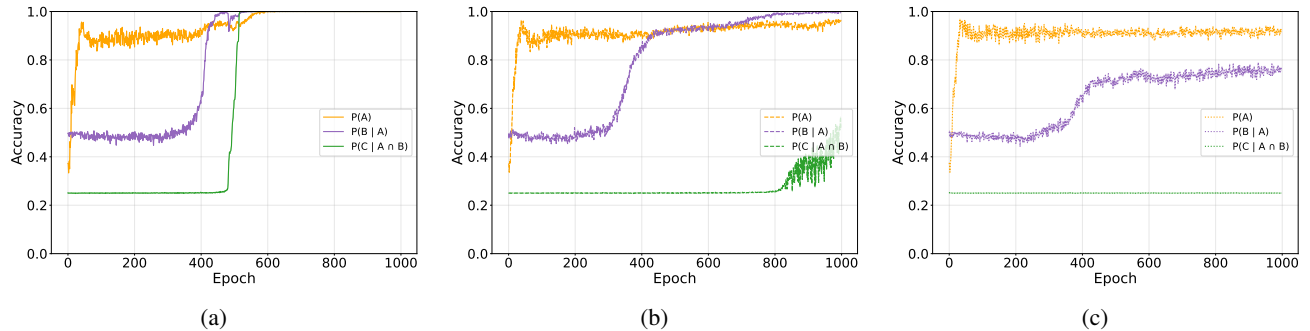


(a)                                    (b)                                    (c)

*Figure 10.* Example of hyperparameter sensitivity on the 3-hop decomposition learning stages. (a) Stages with weight decay = 0.1. (b) Stages with weight decay = 0.01. (c) Stages with weight decay = 0.001. In each subfigure, 'orange' denotes in-support prediction $\mathbb{P}[A]$, 'purple' denotes correct half prediction $\mathbb{P}[B \mid A]$, and 'green' denotes exact match given correct half $\mathbb{P}[C \mid A \cap B]$. Suboptimal hyperparameters can cause the model to delay learning of the stages, resulting in long stagnation segments.

## D. Investigating the delay of learning the correct half in decomposition

In the decomposition experiment, we observed that increasing $hl_{support}$ delayed the learning of the correct half event ($\mathbb{P}[B \mid A]$). To further investigate this effect, we examined the presence of other intermediate stages.

We define $\mathcal{N}(x_q)$ as the set of immediate 1-hop neighbors of the query $x_q$. We then define the following sub-events:

$$EN := \{\hat{y} \in H_c(y) \cup \mathcal{N}(x_q)\} \text{ (extended neighborhood)}$$
$$NR := \{\hat{y} \in H_c(y)\} \qquad \text{(neighborhood refinement)}$$

Event $EN$ shows whether the model has narrowed down the search space to the set of stones for which there are transitions to in the support set, and the set of stones to which one can go to using the potion $z_q$, irrespective of the $x_q$. Event $NR$ tracks whether the model can successfully filter out the neighbors for which the transition exists in the support. We plot the extended neighborhood accuracy $\mathbb{P}[EN \mid A]$ (cyan), and the neighborhood refinement accuracy $\mathbb{P}[NR \mid EN]$ (pink) in Figure 11.

We do not observe any noticeable lag between these events, because the cyan and pink curves rise in immediate succession, indicating a sharp phase transition towards predicting the correct half. This indicates that the multi-hop nature of the support prevents the model from successfully extracting the structure of the chemistry, shown by the chance performance of 12.5% in Figure 5 for an extended period of time. This implies that as soon as the model learns that the target is in the extended neighborhood, it can almost immediately refine the predictions to the correct half.
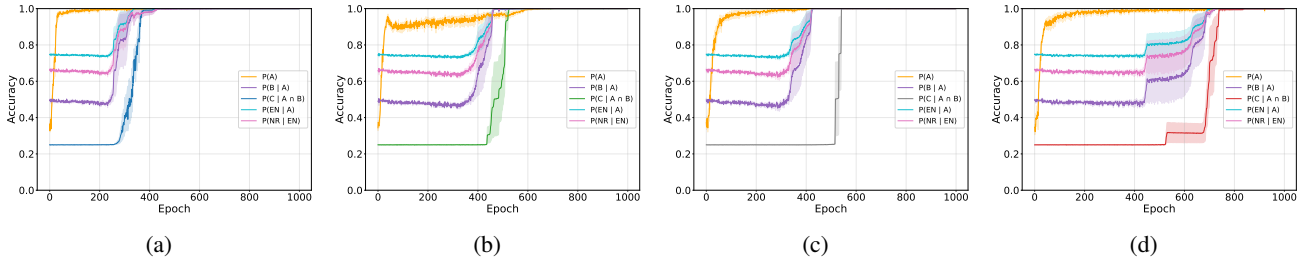


(a)  (b)  (c)  (d)

*Figure 11.* Decomposition staged learning dynamics with the extended neighborhood $EN$ curves and neighborhood refinement $NR$ events. We observed that $EN$ and $NR$ rise simultaneously with the correct half prediction accuracy $\mathbb{P}[B \mid A]$, suggesting the absence of any intermediate stages.