

# Discovering Hidden Algebraic Structures via Transformers with Rank-Aware Beam GRPO

Anonymous authors  
Paper under double-blind review

## Abstract

Recent efforts have extended the capabilities of transformers in logical reasoning and symbolic computations. In this work, we investigate their capacity for non-linear latent pattern discovery in the context of functional decomposition, focusing on the challenging algebraic task of multivariate polynomial decomposition. This problem, with widespread applications in science and engineering, is proved to be NP-hard, and demands both precision and insight. Our contributions are threefold: First, we develop a synthetic data generation pipeline providing fine-grained control over problem complexity. Second, we train transformer models via supervised learning and evaluate them across four key dimensions involving scaling behavior and generalizability. Third, we propose Beam Grouped Relative Policy Optimization (BGRPO), a rank-aware reinforcement learning method suitable for hard algebraic problems. Fine-tuning with rank-aware BGRPO lifts beam-search accuracy by 34–37 percentage points over the SFT initialization and by 1.7–3.7 points over vanilla GRPO, with non-overlapping  $\pm 1\sigma$  bands at every scale. After RL, even greedy decoding surpasses the SFT model’s best beam-search score by 23.8–25.8 percentage points at every scale. Additionally, our model demonstrates competitive performance in polynomial simplification, outperforming Mathematica in various cases.

## 1 Introduction

Transformers, initially developed for natural language processing (Vaswani et al., 2017), have shown remarkable versatility across diverse domains such as vision (Dosovitskiy et al., 2020) and protein folding (Jumper et al., 2021). More recently, their applications in formal reasoning, symbolic mathematics and algorithmic tasks start to gain traction. Several works have demonstrated transformer-based architectures’ ability to tackle highly structured problems, including theorem proving (Polu & Sutskever, 2020; Trinh et al., 2024), integration (Lample & Charton, 2020), matrix multiplication (Fawzi et al., 2022) and equation solving (Drori et al., 2022).

In this work, we investigate the transformer’s capacity for non-linear latent pattern discovery in the context of functional decomposition, i.e. decomposing a complex function as the composition of simpler sub-functions. In contrast to step-by-step logical deduction, or pattern recognition in data analysis, functional decomposition poses significant new challenges to the transformer, because the forms of the sub-functions that we try to discover can be totally hidden or obscured in the final compact form of the original function. Furthermore, it requires extreme precision without any margin of error. Unlike more forgiving classification tasks, the decomposition problem admits only a sparse set of correct solutions: even minor deviations in signs or coefficients can render outputs completely invalid.

Beyond its theoretical interest, functional decomposition has ubiquitous applications in software engineering (Tempero et al., 2024), systems biology (Mori et al., 2023), mechanical design (She et al., 2024), systems engineering (Hernandez et al., 2024) and digital logic design (Adamski et al., 2005; Lin et al., 2008), where capturing hidden substructures within high-dimensional functions leads to more tractable and efficient models. However, identifying a function’s latent compositional structure requires models to look past surface-level correlations, attending instead to deep algebraic symmetries and invariants.

A particularly rich case of functional decomposition arises in multivariate polynomial functions. The polynomial decomposition problem over a ring  $k$  seeks to decompose a given polynomial  $f \in k[x_1, \dots, x_n]$  into polynomials  $g \in k[y_1, \dots, y_m]$  and  $h_1, \dots, h_m \in k[x_1, \dots, x_n]$  such that

$$f(x_1, \dots, x_n) = g(h_1(x_1, \dots, x_n), \dots, h_m(x_1, \dots, x_n)). \quad (1)$$

It has wide-ranging applications from cryptography (Patarin & Goubin, 1997) to dynamical modeling (Dang & Testylier, 2012), signal processing (Demirtas et al., 2012) and robotics (Elias & Wen, 2025; Manocha & Canny, 1992).

The multivariate polynomial decomposition problem has been proved to be NP-hard by Dickerson (Dickerson, 1987; 1993), although efficient algorithms for various special cases are discussed in (Gathen et al., 2003; Von Zur Gathen, 1990a;b; Faugère & Perret, 2009a;b; Zhao et al., 2012). To illustrate the difficulty of the problem for the models, let us consider the following expression

$$\begin{aligned} f = & 2a_1^3b_1^3 + 25a_1^2b_1^2 + 6a_1^2a_2b_2b_1^2 + 6a_1^2a_3b_3b_1^2 + 6a_1a_2^2b_2^2b_1 + 6a_1a_3^2b_3^2b_1 \\ & + 96a_1b_1 + 50a_1a_2b_2b_1 + 50a_1a_3b_3b_1 + 12a_1a_2a_3b_2b_3b_1 + 2a_2^3b_2^3 + 2a_3^3b_3^3 \\ & + 25a_2^2b_2^2 + 25a_3^2b_3^2 + 6a_2a_3^2b_2b_3^2 + 96a_2b_2 + 6a_2^2a_3b_2^2b_3 + 96a_3b_3 \\ & + 50a_2a_3b_2b_3 + 12 \end{aligned}$$

It has a hidden  $O(3)$ -symmetry, which can be revealed by decomposing  $f = g \circ h$ , with  $g(y) = y^2 + 2(4 + y)^3$  and  $h = a_1b_1 + a_2b_2 + a_3b_3$ . This is a highly nontrivial task to identify the inner function  $h$  directly from the expanded form of  $f$ , as its structure becomes completely obscured after polynomial substitution, expansion and simplification. Even in this relatively constrained case where  $g$  is univariate, discovering the decomposition requires recognizing non-linear latent patterns across dozens of terms. When  $g$  becomes multivariate, the complexity increases substantially, making the problem even more challenging.

To tackle the polynomial decomposition problem, we develop a systematic approach with four key components. First, we create a backward synthetic data generation pipeline that allows fine-grained control over polynomial complexity involving range of coefficients, degree, and number of variables. Second, we train lightweight transformer models on these synthetic datasets using supervised learning and analyze how performance scales across four axes (performance complexity scaling, architecture scaling, distribution adaptation, search strategy analysis). Third, we discover that both multi-sampling and greedy search methods struggle with the sparse solution space of the polynomial decomposition problem, and we implement a beam search strategy to effectively extract the models’ capabilities. Finally, we develop a rank-aware variant of the Grouped Relative Policy Optimization (GRPO) reinforcement learning algorithm, which encodes rank information directly in the reward function.

Our study makes the following contributions to neural approaches for polynomial decomposition. First, our backward data generation pipeline enables targeted training across varying levels of decomposition difficulty. Second, our evaluation across four axes ( $\mathcal{D}_1$ – $\mathcal{D}_4$ ) establishes the first baselines for transformers on polynomial decomposition. Third, our rank-aware Beam Grouped Relative Policy Optimization (BGRPO) lifts beam-search accuracy by 34–37 percentage points over the SFT initialization and by 1.7–3.7 points over vanilla GRPO (non-overlapping  $\pm 1\sigma$  bands at every scale); after BGRPO, greedy decoding alone exceeds the SFT model’s best beam-search score by 23.8–25.8 points at every scale. Additionally, our model demonstrates competitive performance in polynomial simplification, outperforming Mathematica in various cases. This shows that neural models can complement and extend classical symbolic computation.

## 2 Method

### 2.1 Backward Synthetic Data Generation

We generate synthetic data for supervised learning using a backward approach, starting from the decomposed form. First, we generate the inner functions ( $h_1, \dots, h_m$  in Eq. (1)) and the outer function ( $g$  in Eq. (1)) with random monomial terms of bounded degree and random coefficients within a given range. Then, we obtain the composed function ( $f$  in Eq. (1)) via substitution, expansion, and term collection. See Appendix A for

the detailed algorithm. For each generated instance, we create a training pair consisting of the expanded polynomial  $f$  as input and its decomposed components  $\{g, h_1, \dots, h_{v_{\text{outer}}}\}$  as the target output. The model is trained to minimize the standard negative log-likelihood loss function.

Our synthetic data generation process provides fine-grained control over problem complexity through eight parameters:  $C_{\text{inner}}$  (coefficient range for inner polynomials),  $d_{\text{inner}}$  (maximum degree of inner polynomials),  $v_{\text{inner}}$  (number of variables in inner polynomials),  $t_{\text{inner}}$  (maximum number of terms in inner polynomials), and similarly  $C_{\text{outer}}$ ,  $d_{\text{outer}}$ ,  $v_{\text{outer}}$ , and  $t_{\text{outer}}$  for the outer polynomial.

## 2.2 Beam Search

Beam search is a breadth-first search algorithm that approximates optimal decoding by keeping track of the  $w$  most probable sequences at each step (Freitag & Al-Onaizan, 2017). For each of the  $w$  current sequences, the algorithm considers the top- $w$  token extensions per sequence. These  $w^2$  candidate continuations are then ranked by the sum of log probabilities of all tokens in the sequence, and only the top- $w$  sequences with the highest cumulative log probability are retained for the next step. In this paper, we refer to  $w$  as the beam width, and to the position (1st, 2nd, etc.) of an output in the final beam as its rank.

Our analysis across all model outputs identified a specific error pattern in polynomial decomposition: the model achieves approximately 90% accuracy for predicting non-sign tokens (operators, numbers, variables), but exhibits near-random performance for deciding between positive and negative signs. This creates a unique inference challenge where exploration needs to be constrained for high-confidence structural elements while simultaneously expanded for uncertain sign choices.

Beam search is particularly well-suited for this situation as it maintains the high-confidence structural backbone while systematically exploring variations in the uncertain components. Our experiments demonstrate that beam search significantly outperforms greedy decoding and random sampling for polynomial decomposition tasks. See Appendix C for a detailed error analysis and an explanation of beam search effectiveness for this task.

## 2.3 BGRPO: A Rank-Aware Reinforcement Learning Method

Supervised training alone leaves the model near-random on a small subset of tokens (the  $\pm$  signs in our case; see Section 2.2), which beam search at inference partially compensates for. To raise accuracy at the policy level rather than only at decoding time, we introduce Beam Grouped Relative Policy Optimization (BGRPO), a reinforcement learning method that extends GRPO by generating its group of outputs through beam search and by rewarding correct answers in proportion to their rank in the beam.

Reinforcement learning enables models to explore solution spaces more effectively than supervised learning alone, enhancing the model’s capabilities by addressing specific weaknesses through a reward mechanism. This approach encourages correct answers while discouraging incorrect ones based on an advantage function—the difference between a solution’s reward and a baseline reward. Group Relative Policy Optimization (GRPO) (Shao et al., 2024) estimates this baseline for each question by sampling a group of outputs, and has shown promising results for reinforcement learning in language generation tasks due to its sample efficiency and stability (DeepSeek-AI, 2025). We chose GRPO over traditional Proximal Policy Optimization (PPO) (Schulman et al., 2017) because it eliminates the need for a separate value network or reward model, reducing training complexity while improving stability, and its group-wise baseline calculation naturally fits tasks with a clear binary reward structure like polynomial decomposition.

Our proposed Beam Grouped Relative Policy Optimization (BGRPO) extends this approach by using beam search rather than independent sampling for generating the group of outputs. While this significantly alters the distribution of outputs, making their average reward less suitable as a traditional baseline, it still provides valid training signals by reinforcing correct answers and penalizing incorrect ones. BGRPO is particularly effective for our task because beam search generates outputs with identical structure that differ only in the confusing elements (signs), creating a focused learning signal.

Additionally, BGRPO incorporates rank information directly into the reward function by applying an exponential decay factor based on the position in the beam. This incentivizes correct answers to appear at earlier positions in the beam search, effectively pushing correct solutions toward the top of the beam ranking.

**Training Objective** For a prompt  $x$ , let  $\mathcal{B}(x) = \{y_1, \dots, y_w\}$  be the set of beam search outputs with beam width  $w$  generated by the old policy  $\pi_{\theta_{\text{old}}}$ , where each  $y_i$  has length  $|y_i|$  tokens. Each output sequence  $y_i$  receives a sequence-level reward  $r_i$ , where  $r_i = 0$  for incorrect polynomial decomposition and  $r_i = 1$  for correct decomposition. In BGRPO, we incorporate rank information by scaling the reward for correct decompositions by  $e^{-\text{rank}/b}$ . The decay base  $b$  is a hyperparameter independent of the beam width  $w$ ; we ablate it in Section 4.5. We optimize the policy model  $\pi_\theta$  for  $\mu$  inner gradient steps by maximizing the following token-level surrogate, following the GRPO formulation of Shao et al. (2024):

$$\mathcal{J}_{\text{BGRPO}}(\theta) = \frac{1}{\sum_{i=1}^w |y_i|} \sum_{i=1}^w \sum_{t=1}^{|y_i|} (\min(\rho_{i,t} A_i, \text{clip}(\rho_{i,t}, 1 - \varepsilon, 1 + \varepsilon) A_i) - \beta \mathbb{D}_{\text{KL}}(\pi_\theta \| \pi_{\text{ref}})_{i,t}), \quad (2)$$

where  $\rho_{i,t} = \pi_\theta(y_{i,t} | x, y_{i,<t}) / \pi_{\theta_{\text{old}}}(y_{i,t} | x, y_{i,<t})$  is the per-token importance ratio,  $\varepsilon$  is the PPO clipping parameter, and  $\beta$  controls the KL regularization. The per-token KL is approximated as

$$\mathbb{D}_{\text{KL}}(\pi_\theta \| \pi_{\text{ref}})_{i,t} = \frac{\pi_{\text{ref}}(y_{i,t} | x, y_{i,<t})}{\pi_\theta(y_{i,t} | x, y_{i,<t})} - \log \frac{\pi_{\text{ref}}(y_{i,t} | x, y_{i,<t})}{\pi_\theta(y_{i,t} | x, y_{i,<t})} - 1. \quad (3)$$

Here,  $\pi_{\text{ref}}$  is the reference policy, fixed to the SFT initialization before BGRPO training. The advantage  $A_i$  is shared across all tokens of beam  $y_i$  and is computed without standard-deviation normalization as  $A_i = r_i - \bar{r}$ , where  $\bar{r} = \frac{1}{w} \sum_{j=1}^w r_j$ , following Liu et al. (2025).

## 3 Experimental Setup

### 3.1 Evaluation Axes

To systematically analyze our models’ capabilities for the polynomial decomposition problem, we consider four key evaluation dimensions.

**Problem Complexity Scaling ( $\mathcal{D}_1$ ).** We analyze how the model performance varies with respect to changes in the complexity parameters for synthetic data generation. We vary the number of variables  $v_{\text{inner}}, v_{\text{outer}}$ , and the maximum degrees  $d_{\text{inner}}, d_{\text{outer}}$  for both the inner and outer polynomials.

**Architecture Scaling ( $\mathcal{D}_2$ ).** We investigate how model performance scales with key architectural hyperparameters of the transformer. In particular, we measure  $\mathcal{P}(M(d, l, a))$ , the performance of models with embedding dimension  $d$ , number of layers  $l$ , and number of attention heads  $a$ . Our goal is to characterize how these hyperparameters influence model capabilities.

**Distribution Adaptation ( $\mathcal{D}_3$ ).** A practical challenge in applying transformers to symbolic computation is their sensitivity to the numerical ranges present in the training data. For example, models trained on specific coefficient ranges tend to struggle with polynomials outside these ranges. On the other hand, we found that models can rapidly adapt to new coefficient distributions with minimal additional training, suggesting that they manage to learn generalizable pattern recognition rather than merely memorizing specific numerical relationships.

To quantify the model’s ability to transfer its polynomial decomposition skills to numerically distinct but structurally identical problems, we prepare the model  $M_{C_1 \rightarrow C_2}^n$ . This model is initially trained on 1M polynomial decomposition examples with  $C_{\text{outer}} = C_1$  and then fine-tuned with  $n$  examples with  $C_{\text{outer}} = C_2$  where  $C_1 \cap C_2 = \emptyset$ . We measure the performance of model  $M_{C_1 \rightarrow C_2}^n$  on a test set of polynomial decomposition problems with  $C_{\text{outer}} = C_2$ :

$$\mathcal{G}(n) = \mathcal{P}(M_{C_1 \rightarrow C_2}^n, \text{test set with } C_{\text{outer}} = C_2) \quad (4)$$

**Search Strategy Analysis ( $\mathcal{D}_4$ ).** We investigate how beam search enhances model performance on polynomial decomposition tasks, analyzing its effectiveness across different model architectures and levels of problem complexity.

### 3.2 Synthetic Dataset Setup

For the axis  $\mathcal{D}_1$  of the problem complexity scaling, we first examine degree scaling by training a model on 2M polynomial decomposition examples with different inner and outer degrees as described in Table 1. We then evaluate this model on separate test datasets with the same configuration parameters, each corresponding to one of nine different  $(d_{\text{inner}}, d_{\text{outer}})$  pairs to assess performance across varying problem complexities.

For the second part of the  $\mathcal{D}_1$  axis, we train a model for each combination of  $v_{\text{inner}}$  and  $v_{\text{outer}}$  varying from 2 to 4 while fixing the other parameter at 3. For each combination, we use 1M examples to train the model.

For the axis  $\mathcal{D}_2$  of architecture scaling, we train multiple models with varying architectural configurations, all using the same dataset of 2M examples with polynomial parameters as described in Table 1.

For the axis  $\mathcal{D}_3$  of distribution adaptation, we train initial models on 1M examples with  $C_{\text{outer}} = C_1 = [-5, 5]$  and then adapt them to examples with  $C_{\text{outer}} = C_2 = [-10, -6] \cup [6, 10]$ . Other parameters are the same across both datasets as described in Table 1.

For the second part of  $\mathcal{D}_1$  (Variable Scaling) and  $\mathcal{D}_2$ , we set  $t_{\text{inner}} = t_{\text{outer}} = 3$  to prevent expressions from becoming too long. We describe our tokenization in Appendix B.

Table 1: Synthetic Dataset Configuration Across Evaluation Axes

Evaluation Axis	Inner Coeff.	Outer Coeff.	Inner Degrees	Outer Degrees	Inner Vars	Outer Vars
$\mathcal{D}_1$ (Degree Scaling)	$[-20, 20]$	$[-20, 20]$	$\{2, 3, 4\}$	$\{2, 3, 4\}$	1	1
$\mathcal{D}_1$ (Variable Scaling)	$[-5, 5]$	$[-5, 5]$	3	3	$\{2, 3, 4\}$	$\{2, 3, 4\}$
$\mathcal{D}_2$ (Architecture)	$[-5, 5]$	$[-5, 5]$	3	3	3	3
$\mathcal{D}_3$ (Adaptation)	$[-20, 20]$	$C_1 = [-5, 5]$	$\{1, 2\}$	$\{1, 2, 3, 4\}$	1	1
	$[-20, 20]$	$C_2 = [-10, -6] \cup [6, 10]$	$\{1, 2\}$	$\{1, 2, 3, 4\}$	1	1

### 3.3 Architecture Configuration

We employ a decoder-only transformer architecture following standard design principles (Vaswani et al., 2017). Table 2 summarizes our task-specific configurations across all experimental axes. For lightweight and effective training, we developed our own model and training pipeline based on `minGPT` (Karpathy, 2020).

Table 2: Transformer Model Configuration Across Experiments

Experiment	Context Window	Embedding Dim.	Layers	Heads
$\mathcal{D}_1$ (Degree Scaling)	256	512	6	8
$\mathcal{D}_1$ (Variable Scaling)	850	512	6	8
$\mathcal{D}_2$ (Architecture)	850	$\{256, 512, 768\}$	$\{4, 6\}$	8
$\mathcal{D}_2$ (Attention Heads)	850	512	6	$\{4, 8, 16\}$
$\mathcal{D}_3$ (Distribution)	256	512	4	8

*Common settings: GELU activation, learned positional embeddings, multi-head attention with causal masking, MLP hidden dimension =  $4 \times$  embedding dimension.*

### 3.4 Supervised Learning Details

We train our models using the Adam optimizer with an initial learning rate of  $6 \times 10^{-4}$ , incorporating a 10% warmup period followed by cosine decay. Each configuration initially trains on 1M instances, with additional 1M training examples added incrementally until performance saturation. We use a batch size of 200 throughout training. We train models with enough epochs until it saturates with the given dataset.

### 3.5 BGRPO Implementation

For the BGRPO reinforcement learning phase we generate candidate solutions using beam search with a width of 32 and temperature of 1.0. The rank-aware reward uses decay base  $b = 20$  as the default; the advantage is computed without standard-deviation normalization as  $A_i = r_i - \bar{r}$ , matching the specification in Section 2 and following Liu et al. (2025). We set the PPO clipping parameter  $\epsilon$  to 0.2 and the KL divergence coefficient  $\beta$  to 0.01, with a learning rate of  $1 \times 10^{-5}$ . Each outer training step samples 8 distinct polynomial decomposition problems from a held-in pool of 200 non-repeating problems, disjoint from the evaluation set, and applies 5 inner gradient steps. We train for 420 outer steps and report results averaged across 3 random seeds (148, 1, 2). Evaluation uses greedy decoding on 200 held-out problems and beam search of width 30 on the same 200 problems, drawn from the multi-variable test set used in  $\mathcal{D}_2$ .

## 4 Experimental Results

### 4.1 Problem Complexity Scaling ( $\mathcal{D}_1$ )

In the first part of  $\mathcal{D}_1$ , we examine how model performance varies with the degrees of inner and outer polynomials. The result is shown in Figure 1. We use greedy search for the inference. Regardless of the degrees of the polynomials, our model achieves a remarkable single-output accuracy. Notably, when using beam search with a width of 10, the model’s accuracy reaches 100% for these configurations.

Our analysis reveals a pattern: performance remains invariant to increases in the outer polynomial’s degree, while decreasing when the inner polynomial’s degree increases. This demonstrates that the transformer’s decomposition capability is primarily limited by the complexity of the inner polynomial rather than that of the outer polynomial.

In the second part of  $\mathcal{D}_1$ , we investigate how the performance scales with  $v_{\text{inner}}$  and  $v_{\text{outer}}$ , the number of variables in the inner and outer polynomials. Figures 2 and 3 present these results.

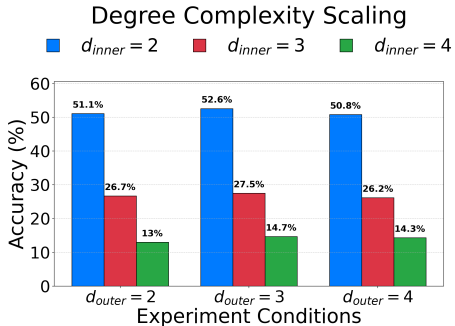


Figure 1: Performance across different  $d_{\text{inner}}, d_{\text{outer}}$

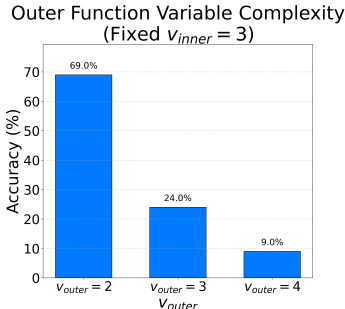


Figure 2: Performance across different  $v_{\text{outer}}$

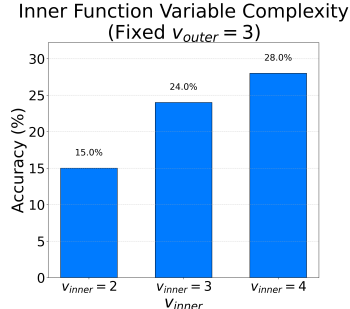


Figure 3: Performance across different  $v_{\text{inner}}$

Given the challenging nature of multivariate polynomial decomposition, we evaluate the model’s performance using beam search with a width of 30, considering a prediction correct if at least one of the 30 candidate outputs is correct decomposition.

Our results reveal two trends: performance decreases dramatically as  $v_{\text{outer}}$  increases, yet counter-intuitively improves as  $v_{\text{inner}}$  increases. This observation aligns with the following heuristic understanding: higher  $v_{\text{outer}}$  creates an information bottleneck, requiring the model to simultaneously resolve multiple interdependent inner functions. In contrast, higher  $v_{\text{inner}}$  provides more dimensions of input variation with additional structural indicators that can guide the decomposition process.

## 4.2 Architecture Scaling ( $\mathcal{D}_2$ )

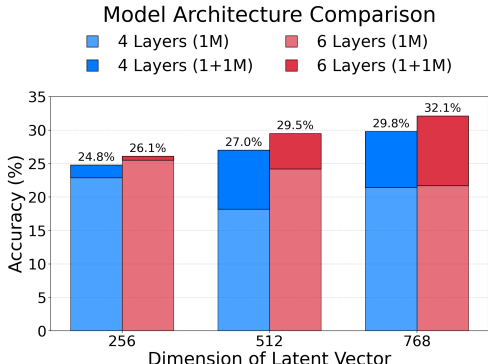


Figure 4: Accuracies on different number of layer and dimension.

Figure 4 reveals the scaling behavior (Kaplan et al., 2020) of transformer architectures on polynomial decomposition. As model capacity increases through higher embedding dimensions and additional layers, performance consistently improves. Notably, our results demonstrate the presence of a data-dependent scaling threshold. With limited training data (1M examples), larger models initially underperform their simpler counterparts, particularly evident in the 6-layer configurations with higher embedding dimensions. However, this pattern reverses completely with additional training data, confirming that larger models possess superior capacity for mathematical pattern recognition when provided with sufficient examples to leverage their parametric advantage.

In  $\mathcal{D}_2$ , we also examine model performance with different numbers of attention heads. Our experiments reveal that increasing the number of attention heads while maintaining constant total embedding dimension leads to progressively deteriorating performance on polynomial decomposition tasks. Models with 4 heads achieved 32.0% accuracy, while those with 8 and 16 heads reached only 28.0% and 25.0% accuracy, respectively. This suggests that for our specific task of mathematical pattern recognition, fewer, more expressive attention heads with larger per-head dimensions provide better performance than numerous specialized heads with smaller dimensions.

## 4.3 Distribution Adaptation ( $\mathcal{D}_3$ )

We evaluate  $\mathcal{G}(n)$  as defined in Eq. 4, which measures how quickly models adapt to new coefficient distributions as a function of adaptation sample size  $n$ . For this experiment, we train a model with 4 layers and 512 embedding dimension on the dataset described in Section 3.2. The initial training used 1M examples with outer polynomial coefficient range  $C_1$ , followed by fine-tuning on  $n$  examples with coefficient range  $C_2$  for a single epoch. We report the variance in accuracy based on three independent trials.

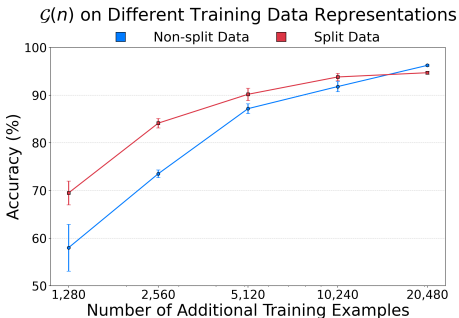


Figure 5: Performance recovery when adapting to a new coefficient distribution

randomly select terms from the expanded form and split their coefficients. For example:

Models trained exclusively on the first dataset achieve only 5.67% accuracy on the new distribution, despite reaching nearly 100% accuracy on the original distribution. Figure 5 illustrates how performance recovers during adaptation. Notably, despite using only  $\approx 2\%$  of the original training data size, the model rapidly recovers its accuracy from single digits to over 90%. This rapid adaptation indicates successful transfer learning, suggesting that the model develops a general mathematical understanding of polynomial substructures rather than memorizing specific numerical relationships.

We further investigate whether alternative data representations could enhance this adaptation capability. We propose "split" representation of polynomials, where we

$$\begin{aligned}
f_{\text{non-split}}(a) &= -63 + 23a - 71a^2 - 11a^3 - 14a^4 - 12a^5 - 2a^6 \\
f_{\text{split}}(a) &= -63 + 23a - 4a^2 - 67a^2 - 8a^3 - 3a^3 - 7a^4 - 7a^4 - 12a^5 - a^6 - a^6
\end{aligned}
\tag{5}$$

In Figure 5, the red line demonstrates  $\mathcal{G}(n)$  of the model trained on data with both normal and split representation. Models trained on this mixed data including split representation demonstrate significantly faster adaptation, requiring only 70% of the additional training examples to reach equivalent performance on the new distribution.

This enhanced generalization likely stems from the model being forced to recognize mathematically equivalent but differently represented polynomials, compelling it to develop a deeper understanding of polynomial structure rather than memorizing specific patterns.

#### 4.4 Search Strategy Analysis ( $\mathcal{D}_4$ )

We evaluate how search strategies impact model performance on polynomial decomposition, comparing greedy decoding against beam search across model scales and problem complexities. Figure 6 and 7 illustrate the accuracy achieved across different beam widths for polynomials with varying numbers of variables.

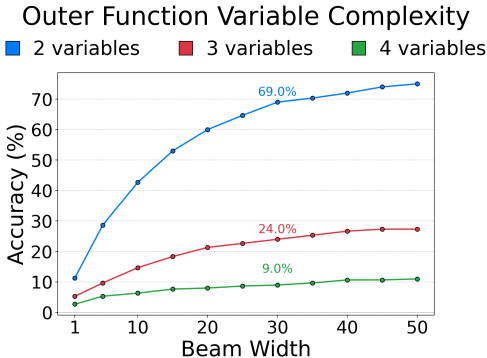


Figure 6: Beam width scaling with varying  $v_{\text{outer}}$  ( $v_{\text{inner}} = 3$ )

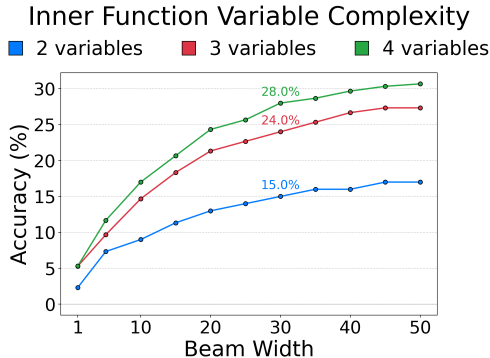


Figure 7: Beam width scaling with varying  $v_{\text{inner}}$  ( $v_{\text{outer}} = 3$ )

Our results reveal an unusually large impact of beam search for polynomial decomposition compared to typical NLP tasks. For two-variable polynomials, accuracy improves from 11% with greedy search to 69% with a beam width of 30, a  $6.3\times$  gap. This is much larger than in standard neural machine translation, where beam search typically yields BLEU improvements of only 2–4 points (Huang et al., 2018; Ranzato et al., 2016) and most systems show diminishing returns with beam widths beyond 5–10 (Freitag & Al-Onaizan, 2017). We return to this gap in Section 4.5, where rank-aware BGRPO closes it: greedy decoding alone after BGRPO already exceeds the SFT model’s best beam-search score at every scale.

#### 4.5 BGRPO Results

Rank-aware BGRPO raises beam@30 accuracy over the SFT initialization by 34–37 accuracy points at every scale. On the three 6-layer SFT models from  $\mathcal{D}_2$  (embedding dimensions 256, 512, 768), beam@30 climbs from 16.5%, 20.0%, and 20.0% to  $53.5 \pm 0.4\%$ ,  $54.5 \pm 0.4\%$ , and  $53.7 \pm 0.2\%$  respectively. All numbers in this section are mean  $\pm$  standard deviation across 3 training seeds, evaluated on 200 held-out multi-variable test problems (beam width 30; see Section 3.5).

Figure 8 traces the full training trajectory for each scale, comparing four conditions at matched compute (420 outer training steps, 8 problems per step): the SFT initialization (dashed gray), vanilla GRPO, BGRPO with a binary reward, and BGRPO with the rank-aware reward at decay base  $b = 20$ . Lines are mean  $\pm 1\sigma$  across 3 training seeds. Two observations stand out. First, vanilla GRPO already closes most of the gap from the SFT initialization to the rank-aware ceiling, reaching  $49.8 \pm 1.0\%$ ,  $52.3 \pm 0.2\%$ , and  $52.0 \pm 0.4\%$

at the three scales. The rank-aware reward adds another 1.7–3.7 accuracy points on top of GRPO, with non-overlapping  $\pm 1\sigma$  bands at the converged plateau (outer steps 300–420) at every scale. Second, BGRPO with a flat binary reward is the weakest of the three RL variants at  $d = 512$  and  $d = 768$ , falling 13–15 points below vanilla GRPO. Beam-search rollout offers no advantage over independent sampling when the policy update cannot distinguish among the correct beams it returns: a flat reward assigns the same advantage to every correct beam, discarding the rank signal that the beam rollout collected.

### BGRPO training trajectory at beam@30

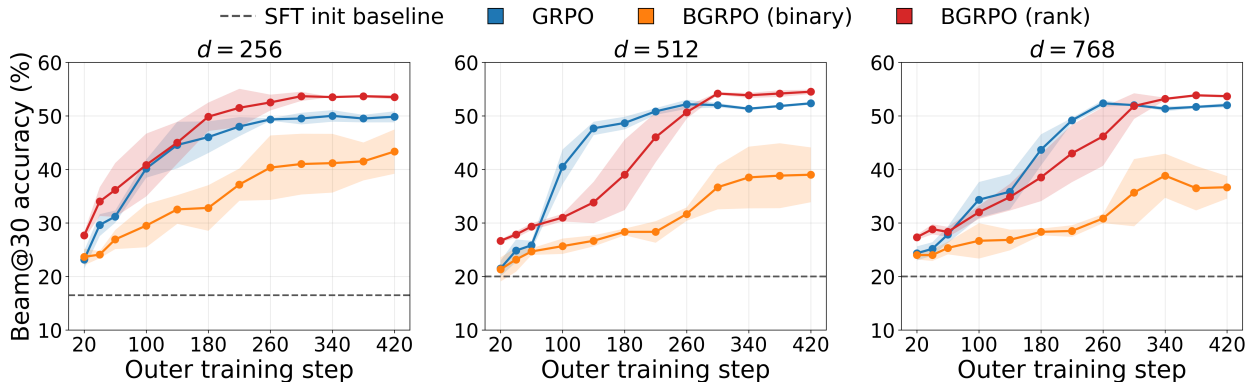


Figure 8: Beam@30 accuracy over outer training step for the three SFT initializations. Each panel compares the SFT initialization (dashed gray) against vanilla GRPO, BGRPO with a binary reward, and BGRPO with the rank-aware reward at decay base  $b = 20$ . Shaded bands are  $\pm 1\sigma$  across 3 training seeds, evaluated on 200 held-out multi-variable problems.

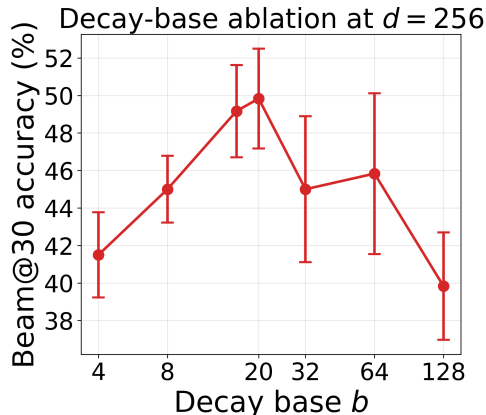


Figure 9: Decay-base ablation on the  $d = 256$  initialization at outer step 180. Mean  $\pm 1\sigma$  across 3 training seeds. The peak sits at  $b = 20$ , near the beam width 32.

**Decay-base ablation.** The rank-aware reward scales correct outputs by  $e^{-\text{rank}/b}$  (Section 2); the decay base  $b$  is the only free hyperparameter introduced by the rank signal, and BGRPO is highly sensitive to it. Figure 9 sweeps  $b \in \{4, 8, 16, 20, 32, 64, 128\}$  on the  $d = 256$  initialization at outer step 180 (the common training budget across all sweep runs). The peak-to-trough range spans roughly 10 accuracy points: the curve is unimodal, peaking at  $b = 20$  ( $49.8 \pm 2.6\%$ ) and dropping to  $41.5 \pm 2.3\%$  at  $b = 4$  and  $39.8 \pm 2.9\%$  at  $b = 128$ . The peak sits near the beam width 32, which is the scale at which the exponential can meaningfully distinguish the top of the beam from the bottom. A sharper decay ( $b = 4$ ) concentrates all credit on rank 1 and starves the rest of the beam, while a flatter decay ( $b = 128$ ) approaches the binary case where every correct beam looks identical to the policy update. The remainder of this section uses  $b = 20$ .

**Greedy decoding after RL.** After rank-aware BGRPO, greedy decoding alone surpasses the SFT model’s best beam-search score at every scale (Table 3). At  $d = 256$ , the rank-aware model decodes greedily at  $42.3 \pm 2.5\%$ , exceeding the SFT model’s

beam@30 score of 16.5% by 25.8 points. The corresponding greedy scores at  $d = 512$  and  $d = 768$  are  $45.7 \pm 1.2\%$  and  $43.8 \pm 2.8\%$ , against SFT-beam@30 of 20.0% and 20.0% respectively. Notably, the asymmetry between beam search and greedy decoding flips between the two RL variants: vanilla GRPO is 2.0–3.5 points higher than rank-aware BGRPO at greedy decoding (e.g.,  $48.7 \pm 0.6\%$  vs  $45.7 \pm 1.2\%$  at  $d = 512$ ), while rank-aware BGRPO is 1.7–3.7 points higher at beam@30. This pattern is consistent with the design of the rank decay, which redistributes probability mass into the upper region of the beam rather than strictly to rank 1; the small cost at the very top is recovered and exceeded once more than one beam position is examined.

Table 3: Greedy decoding after rank-aware BGRPO compared with SFT-beam@30. Greedy values are mean  $\pm 1\sigma$  across 3 training seeds, evaluated on 200 held-out multi-variable problems.

	SFT (beam@30)	rank-aware BGRPO (greedy)	gain
$d = 256$	16.5%	$42.3 \pm 2.5\%$	+25.8
$d = 512$	20.0%	$45.7 \pm 1.2\%$	+25.7
$d = 768$	20.0%	$43.8 \pm 2.8\%$	+23.8

**Architecture-scaling compression.** The SFT models exhibit a 3.5-point spread in beam@30 accuracy across the three embedding dimensions (16.5% / 20.0% / 20.0% at  $d = 256$  /  $d = 512$  /  $d = 768$ ). After rank-aware BGRPO the spread collapses to 1.0 point (53.5% / 54.5% / 53.7%), a factor of roughly  $3.5\times$  compression. Greedy decoding shows the same pattern: an SFT greedy spread of 5.0 points (10.0% / 10.0% / 15.0%) becomes 3.4 points (42.3% / 45.7% / 43.8%) after rank-aware BGRPO. The smallest model gains the most under RL:  $d = 256$  gains 37 accuracy points at beam@30, against 34.5 and 33.7 points for  $d = 512$  and  $d = 768$ . RL post-training partially substitutes for parameter count on this task; whether the architecture-scaling gap reopens at larger model sizes is left to future work.

#### 4.6 Simplification Comparison with Mathematica

While polynomial simplification and polynomial decomposition represent two distinct mathematical objectives, simplification frequently arises as a consequence of decomposition, since decomposed forms generally exhibit reduced algebraic complexity compared to the original expression. In this subsection, we briefly explore the capabilities of our models for this related problem, and benchmark against the most powerful symbolic computation engine Mathematica. Despite our lightweight parameter budgets and the absence of any explicit simplification objective in our training, the models were able to reduce the leaf count (Wolfram Research, Inc., 1996) of complex expressions, with performance on par with — and in two of five complexity regimes surpassing — Mathematica’s state-of-the-art FullSimplify function (see Table 4, competitive performances are bolded).

Table 4: Average leaf count comparison (Beam width = 30)

Problem Complexity		Leaf Count (mean)		
$v_O$	$v_S$	Transformer	Mathematica	$\Delta$
2	3	<b>27.28</b>	30.03	<b>-2.75</b>
3	3	22.85	<b>22.12</b>	<b>0.73</b>
4	3	22.52	20.00	2.52
3	2	17.27	<b>17.10</b>	<b>0.17</b>
3	4	<b>26.04</b>	27.56	<b>-1.52</b>

These findings highlight that transformers’ inherent ability to uncover latent patterns rivals that of the most advanced symbolic computation methods.

## 5 Conclusion

Our investigation into transformers for polynomial decomposition shows how neural networks can infer hidden algebraic structures.

We find that model performance depends asymmetrically on polynomial complexity parameters ( $\mathcal{D}_1$ ): inner polynomial degree plays a dominant role, while outer polynomial complexity has limited impact. Counter-intuitively, increasing the number of inner variables improves accuracy by imposing structural constraints, whereas more outer variables create information bottlenecks.

From an architectural viewpoint ( $\mathcal{D}_2$ ), we confirm that performance scales with model size. We observe that fewer but more expressive attention heads are especially effective for this task. In terms of distribution

adaptation ( $\mathcal{D}_3$ ), models transfer rapidly to new coefficient distributions, requiring as little as 2% of the original training data, indicating that they internalize generalizable principles rather than rely on memorization. Moreover, we can enhance this generalization capability through strategic dataset design.

Beam search analysis ( $\mathcal{D}_4$ ) yields up to  $6.3\times$  improvement over greedy decoding for the SFT models, due to the sparse, precise nature of mathematical solutions. After fine-tuning with rank-aware BGRPO, beam-search accuracy increases by 34–37 percentage points over the SFT initialization and by 1.7–3.7 points over vanilla GRPO with non-overlapping  $\pm 1\sigma$  bands at every scale; greedy decoding alone after BGRPO already exceeds the SFT model’s best beam-search score by 23.8–25.8 points at every scale. Lastly, our model demonstrates competitive performance in polynomial simplification compared with symbolic computation tools in Mathematica.

Our work provides, for the first time, a systematic analysis of transformer capabilities for polynomial decomposition through controlled experiments across four dimensions. Our methodologies can serve as a road map for exploring neural models in other domains that require non-local latent pattern discovery, such as functional decomposition problems ranging from systems engineering and mechanical design to digital logic design. While we developed BGRPO specifically for the polynomial decomposition problem, similar techniques may prove useful in other domains with sparse solution spaces where models can identify correct structures but struggle with specific details.

**Limitations** Our generalizability investigation was constrained to univariate polynomials with relatively narrow coefficient ranges and limited maximum degrees. Computational constraints restricted our architecture scaling experiments to relatively small models (maximum 6 layers, 768-dimensional embeddings); however, the consistent performance improvements without accuracy saturation suggest that further scaling would yield additional gains. Finally, our three rank-aware models converge to within 1 point of each other at the end of training (53.5% / 54.5% / 53.7% at  $d = 256$  /  $d = 512$  /  $d = 768$ ), suggesting that the 200-problem multi-variable test set may impose a ceiling near 55% that limits our ability to distinguish the larger models from the smaller one; a more challenging evaluation set may show wider separation.

### Broader Impact Statement

This work investigates transformers applied to the polynomial decomposition problem, which is foundational research not tied to immediate societal applications. The synthetic data used in our experiments carries no privacy concerns, and the models are trained for a well-defined mathematical task with no direct path to harmful applications.

### Author Contributions

### Acknowledgments

### References

- Marian Andrzej Adamski, Andrei Karatkevich, Marek Wegrzyn, Mariusz Rawski, Tadeusz Łuba, Zbigniew Jachna, and Paweł Tomaszewicz. The influence of functional decomposition on modern digital design process. *Design of Embedded Control Systems*, pp. 193–204, 2005.
- Thao Dang and Romain Testylier. Reachability analysis for polynomial dynamical systems using the bernstein expansion. *Reliab. Comput.*, 17(2):128–152, 2012.
- DeepSeek-AI. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning, 2025. URL <https://arxiv.org/abs/2501.12948>.
- Sefa Demirtas, Guolong Su, and Alan V Oppenheim. Sensitivity of polynomial composition and decomposition for signal processing applications. In *2012 Conference Record of the Forty Sixth Asilomar Conference on Signals, Systems and Computers (ASILOMAR)*, pp. 391–395. IEEE, 2012.
- Matthew T Dickerson. Polynomial decomposition algorithms for multivariate polynomials. Technical report, Cornell University, 1987.

- Matthew T Dickerson. General polynomial decomposition and the s-1-decomposition are np-hard. *International Journal of Foundations of Computer Science*, 4(02):147–156, 1993.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.
- Iddo Drori, Sarah Zhang, Reece Shuttleworth, Leonard Tang, Albert Lu, Elizabeth Ke, Kevin Liu, Linda Chen, Sunny Tran, Newman Cheng, et al. A neural network solves, explains, and generates university math problems by program synthesis and few-shot learning at human level. *Proceedings of the National Academy of Sciences*, 119(32):e2123433119, 2022.
- Alexander J Elias and John T Wen. Ik-geo: Unified robot inverse kinematics using subproblem decomposition. *Mechanism and Machine Theory*, 209:105971, 2025.
- Jean-Charles Faugère and Ludovic Perret. An efficient algorithm for decomposing multivariate polynomials and its applications to cryptography. *Journal of Symbolic Computation*, 44(12):1676–1689, 2009a.
- Jean-Charles Faugère and Ludovic Perret. High order derivatives and decomposition of multivariate polynomials. In *Proceedings of the 2009 international symposium on Symbolic and algebraic computation*, pp. 207–214, 2009b.
- Allhussein Fawzi, Khurram Kozhasov, Micah Goldblum, Jens Behrmann, Chaoning Zhang, Fabian Fuchs, Po-Sen Huang, Lala Li, and Pushmeet Kohli. Discovering faster matrix multiplication algorithms with reinforcement learning. In *International Conference on Machine Learning (ICML)*, 2022.
- Markus Freitag and Yaser Al-Onaizan. Beam search strategies for neural machine translation. In *Proceedings of the First Workshop on Neural Machine Translation*. Association for Computational Linguistics, 2017. doi: 10.18653/v1/w17-3207. URL <http://dx.doi.org/10.18653/v1/W17-3207>.
- Joachim von zur Gathen, Jaime Gutierrez, and Rosario Rubio. Multivariate polynomial decomposition. *Applicable Algebra in Engineering, Communication and Computing*, 14(1):11–31, 2003.
- Isabella Hernandez, Bryan C Watson, Marc J Weissburg, and Bert Bras. Using functional decomposition to bridge the design gap between desired emergent multi-agent-system resilience and individual agent design. *Systems Engineering*, 27(5):911–930, 2024.
- Liang Huang, Kai Zhao, and Mingbo Ma. When to finish? optimal beam search for neural text generation (modulo beam size), 2018. URL <https://arxiv.org/abs/1809.00069>.
- John Jumper, Richard Evans, Alexander Pritzel, Tim Green, Michael Figurnov, Olaf Ronneberger, Kathryn Tunyasuvunakool, Russ Bates, Augustin Žídek, Anna Potapenko, et al. Highly accurate protein structure prediction with alphafold. *nature*, 596(7873):583–589, 2021.
- Jared Kaplan, Sam McCandlish, Tom Henighan, and Tom B. Brown. Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361*, 2020.
- Andrej Karpathy. mingpt. <https://github.com/karpathy/minGPT>, 2020.
- Guillaume Lample and François Charton. Deep learning for symbolic mathematics. *arXiv preprint arXiv:2006.02974*, 2020.
- Hsuan-Po Lin, Jie-Hong R Jiang, and Ruei-Rung Lee. To sat or not to sat: Ashenurst decomposition in a large scale. In *2008 IEEE/ACM International Conference on Computer-Aided Design*, pp. 32–37. IEEE, 2008.
- Zichen Liu, Changyu Chen, Wenjun Li, Penghui Qi, Tianyu Pang, Chao Du, Wee Sun Lee, and Min Lin. Understanding r1-zero-like training: A critical perspective, 2025. URL <https://arxiv.org/abs/2503.20783>.

- Dinesh Manocha and John F Canny. Real time inverse kinematics for general 6r manipulators. In *ICRA*, pp. 383–389, 1992.
- Matteo Mori, Chuankai Cheng, Brian R Taylor, Hiroyuki Okano, and Terence Hwa. Functional decomposition of metabolism allows a system-level quantification of fluxes and protein allocation towards specific metabolic functions. *Nature Communications*, 14(1):4161, 2023.
- Catherine Olsson, Nelson Elhage, Neel Nanda, Nicholas Joseph, Nova DasSarma, Tom Henighan, Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, Dawn Drain, Deep Ganguli, Zac Hatfield-Dodds, Danny Hernandez, Scott Johnston, Andy Jones, Jackson Kernion, Liane Lovitt, Kamal Ndousse, Dario Amodei, Tom Brown, Jack Clark, Jared Kaplan, Sam McCandlish, and Chris Olah. In-context learning and induction heads. *Transformer Circuits Thread*, 2022. <https://transformer-circuits.pub/2022/in-context-learning-and-induction-heads/index.html>.
- Jacques Patarin and Louis Goubin. Asymmetric cryptography with s-boxes is it easier than expected to design efficient asymmetric cryptosystems? In *International Conference on Information and Communications Security*, pp. 369–380. Springer, 1997.
- Stanislas Polu and Ilya Sutskever. Generative language modeling for automated theorem proving. *arXiv preprint arXiv:2009.03393*, 2020.
- Marc’Aurelio Ranzato, Sumit Chopra, Michael Auli, and Wojciech Zaremba. Sequence level training with recurrent neural networks, 2016. URL <https://arxiv.org/abs/1511.06732>.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms, 2017. URL <https://arxiv.org/abs/1707.06347>.
- Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, Y. K. Li, Y. Wu, and Daya Guo. Deepseekmath: Pushing the limits of mathematical reasoning in open language models, 2024. URL <https://arxiv.org/abs/2402.03300>.
- Jinjuan She, Elise Belanger, and Caroline Bartels. Evaluating the effectiveness of functional decomposition in early-stage design: development and application of problem space exploration metrics. *Research in Engineering Design*, 35(3):311–327, 2024.
- Ewan Tempero, Paul Denny, James Finnie-Ansley, Andrew Luxton-Reilly, Diana Kirk, Juho Leinonen, Asma Shakil, Robert Sheehan, James Tizard, Yu-Cheng Tu, et al. On the comprehensibility of functional decomposition: An empirical study. In *Proceedings of the 32nd IEEE/ACM International Conference on Program Comprehension*, pp. 214–224, 2024.
- Trieu H Trinh, Yuhuai Wu, Quoc V Le, He He, and Thang Luong. Solving olympiad geometry without human demonstrations. *Nature*, 625(7995):476–482, 2024.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- Joachim Von Zur Gathen. Functional decomposition of polynomials: the tame case. *Journal of Symbolic Computation*, 9(3):281–299, 1990a.
- Joachim Von Zur Gathen. Functional decomposition of polynomials: the wild case. *Journal of Symbolic Computation*, 10(5):437–452, 1990b.
- Kevin Wang, Alexandre Variengien, Arthur Conmy, Buck Shlegeris, and Jacob Steinhardt. Interpretability in the wild: a circuit for indirect object identification in gpt-2 small, 2022. URL <https://arxiv.org/abs/2211.00593>.
- Wolfram Research, Inc. `ComplexityFunction` – Wolfram Language Documentation. <https://reference.wolfram.com/language/ref/ComplexityFunction.html>, 1996. Accessed 12 May 2025.



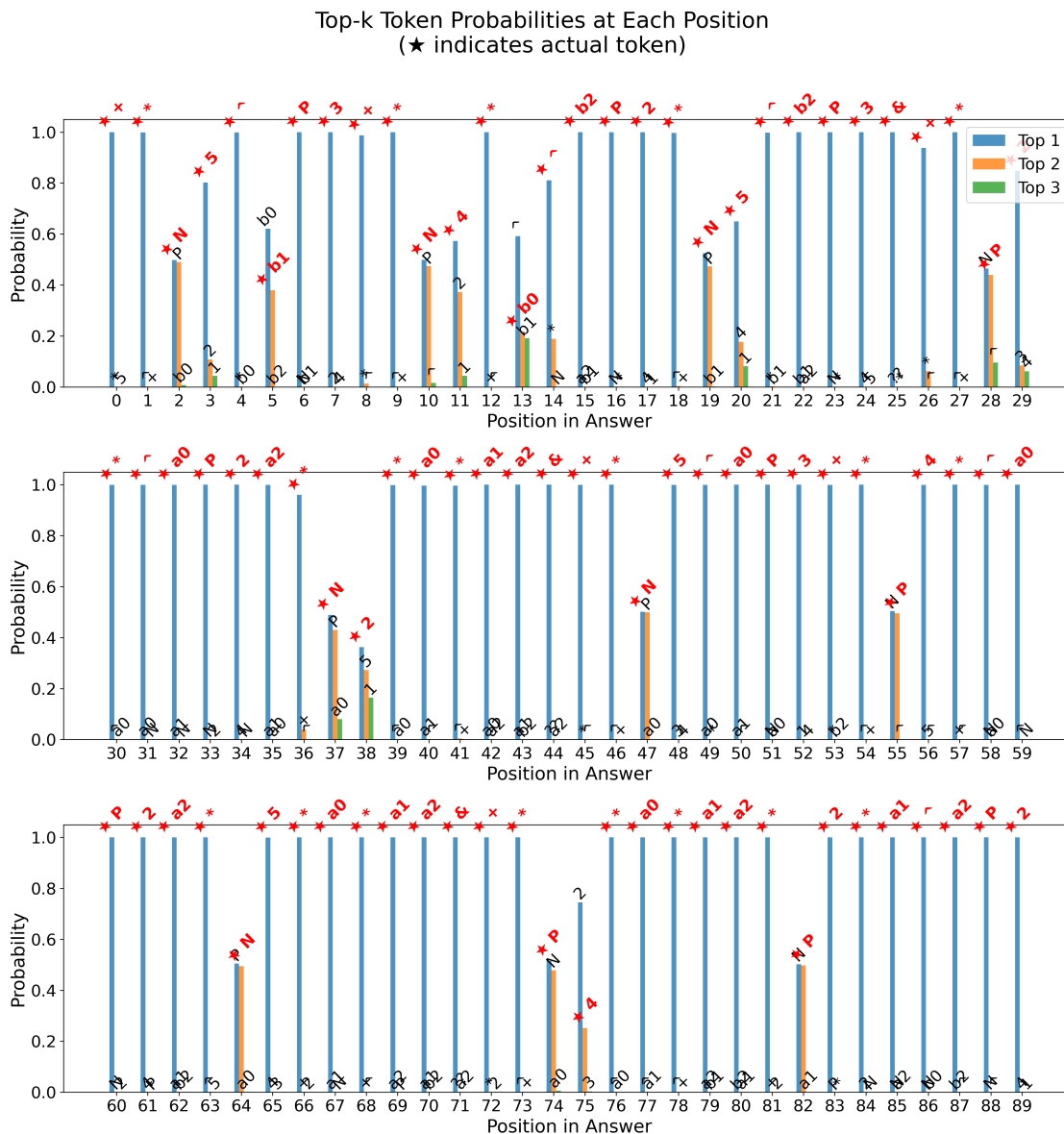


Figure 10: Top-3 probability for each token position in the answer sequence where

**Answer:** + \* N 5 ^ b1 P 3 + \* N 4 \* b0 ^ b2 P 2 \* N 5 ^ b2 P 3 & + \* P 2 \* ^ a0 P 2 a2 \* N 2 \* a0 \* a1 a2 & + \* N 5 ^ a0 P 3 + \* P 4 \* ^ a0 P 2 a2 \* N 5 \* a0 \* a1 a2 & + \* P 4 \* a0 \* a1 a2 \* P 2 \* a1 ^ a2 P 2

**Question:** + \* P 6 2 5 ^ a0 P 9 + \* N 1 5 0 0 \* ^ a0 P 8 a2 + \* P 1 8 7 5 \* ^ a0 P 7 \* a1 a2 + \* P 1 2 0 0 \* ^ a0 P 7 ^ a2 P 2 + \* N 3 0 0 0 \* ^ a0 P 6 \* a1 ^ a2 P 2 + \* P 1 8 7 5 \* ^ a0 P 5 \* ^ a1 P 2 ^ a2 P 2 + \* N 3 2 0 \* ^ a0 P 6 ^ a2 P 3 + \* P 1 2 0 0 \* ^ a0 P 5 \* a1 ^ a2 P 3 + \* N 1 6 2 8 \* ^ a0 P 4 \* ^ a1 P 2 ^ a2 P 3 + \* P 4 3 3 \* ^ a0 P 3 ^ a1 P 3 ^ a2 P 3 + \* N 1 2 8 \* ^ a0 P 3 \* ^ a1 P 2 ^ a2 P 4 + \* N 3 5 2 \* ^ a0 P 2 \* ^ a1 P 3 ^ a2 P 4 + \* N 3 2 \* ^ a0 P 2 \* ^ a1 P 2 ^ a2 P 5 + \* N 2 0 8 \* a0 \* ^ a1 P 3 ^ a2 P 5 \* N 4 0 \* ^ a1 P 3 ^ a2 P 6 ?

Table 5: Token Type Analysis Across Different Model Architectures

Token Type	Metric	4 Layers			6 Layers		
		256 dim	512 dim	768 dim	256 dim	512 dim	768 dim
Sign	Probability	$0.489 \pm 0.001$	$0.489 \pm 0.001$	$0.493 \pm 0.001$	$0.491 \pm 0.001$	$0.490 \pm 0.001$	$0.490 \pm 0.001$
	Accuracy	$0.519 \pm 0.006$	$0.531 \pm 0.006$	$0.530 \pm 0.006$	$0.522 \pm 0.006$	$0.523 \pm 0.006$	$0.521 \pm 0.006$
Operator	Probability	$0.920 \pm 0.002$	$0.915 \pm 0.002$	$0.919 \pm 0.002$	$0.927 \pm 0.002$	$0.925 \pm 0.002$	$0.925 \pm 0.002$
	Accuracy	$0.937 \pm 0.002$	$0.934 \pm 0.002$	$0.935 \pm 0.002$	$0.943 \pm 0.002$	$0.941 \pm 0.002$	$0.942 \pm 0.002$
Number	Probability	$0.880 \pm 0.002$	$0.870 \pm 0.002$	$0.878 \pm 0.002$	$0.890 \pm 0.002$	$0.885 \pm 0.002$	$0.884 \pm 0.002$
	Accuracy	$0.901 \pm 0.002$	$0.893 \pm 0.003$	$0.897 \pm 0.002$	$0.911 \pm 0.002$	$0.905 \pm 0.002$	$0.903 \pm 0.002$

Note: Values shown as mean  $\pm$  standard error of the mean. The sign token probabilities are near-random, while operators and numbers show high confidence and accuracy.

token over a 0.55 probability token. Since our polynomial expressions typically contain fewer than 10 sign decisions, beam search with a width of approximately 30 can efficiently cover most viable sign permutations while maintaining the correct monomial structure identified with high confidence.

## D Attention Score Analysis: Monomial Heads

Attention mechanism analysis has provided valuable insights into transformer model behaviors, with studies identifying specialized attention heads that serve specific functions. For example, Olsson et al. (2022) identified "Induction Heads" that play a central role in in-context learning, while Wang et al. (2022) provided a detailed analysis of indirect object identification in GPT-2 Small.

In our analysis of attention patterns in polynomial decomposition models, we identified specialized attention heads that recognize the structure of polynomials, particularly focusing on monomial identification. We call these "Monomial Heads," and they appear consistently across all model sizes in our architecture scaling experiments ( $\mathcal{D}_2$ ).

Monomial Heads manifest in two distinct patterns in our models. First, in layer 0, several attention heads consistently attend to tokens 1-5 positions behind the current position, as shown in the leftmost plot of Figure 11. Second, in layer 1, we observe specialized behavior where certain heads focus attention on specific tokens within each monomial of the input polynomial (middle plot), while others specifically attend to delimiter tokens in the decomposition output (rightmost plot).

We hypothesize that this represents a two-stage process: in layer 0, the model identifies tokens that serve as monomial indicators by examining local context (1-5 tokens behind). In layer 1, tokens within each monomial attend to these indicator tokens to establish their monomial membership. While this pattern is most clear in the encoding of the input polynomial, the decomposition output shows evidence of boundary recognition, particularly at the transitions between inner functions marked by delimiter tokens.

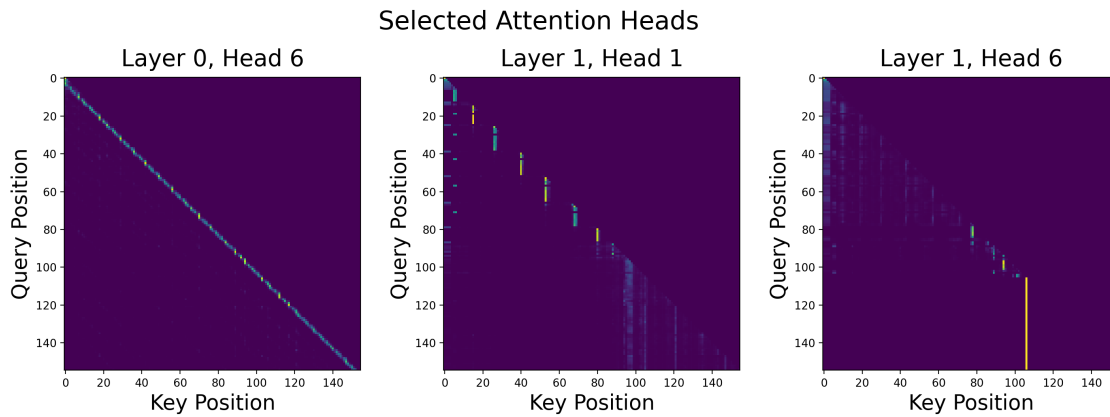


Figure 11: Attention score visualization of selected attention heads from our 6-layer transformer model with embedding dimension 768. The visualization shows attention patterns for a tokenized polynomial sequence and its decomposition.

**Input polynomial:**  $+ * P^2 5 6^{\wedge} a_0 P^9 + * N^1 9 2 *^{\wedge} a_0 P^8 a_1 + * P^4 8 *^{\wedge} a_0 P^7^{\wedge} a_1 P^2 + * N^4 *^{\wedge} a_0 P^6^{\wedge} a_1 P^3 + * N^6 4 *^{\wedge} a_0 P^3^{\wedge} a_1 P^6 + * P^1 6 *^{\wedge} a_0 P^2^{\wedge} a_1 P^7 * P^6 4^{\wedge} a_1 P^9 ?$

**Model's decomposition output:**  $+ * N^4^{\wedge} b_0 P^3 + * b_0^{\wedge} b_2 P^2 * N^1^{\wedge} b_2 P^3 \& + * N^4^{\wedge} a_0 P^3 *^{\wedge} a_0 P^2 a_1 \& + * N^3^{\wedge} a_1 P^3 + * N^2 * a_1^{\wedge} a_2 P^2 * N^4^{\wedge} a_2 P^3 \& * N^4^{\wedge} a_1 P^3$  The visualization reveals how different attention heads focus on specific structural elements when decomposing polynomials.