

# Transformers Can Achieve Length Generalization But Not Robustly

Yongchao Zhou<sup>1,2</sup>, Uri Alon<sup>1</sup>, Xinyun Chen<sup>1</sup>, Xuezhi Wang<sup>1</sup>, Rishabh Agarwal<sup>1</sup> and Denny Zhou<sup>1</sup>

<sup>1</sup>Google DeepMind, <sup>2</sup>University of Toronto

Length generalization, defined as the ability to extrapolate from shorter training sequences to longer test ones, is a significant challenge for language models. This issue persists even with large-scale Transformers handling relatively straightforward tasks. In this paper, we test the Transformer’s ability of length generalization using the task of addition of two integers. We show that the success of length generalization is intricately linked to the data format and the type of position encoding. Using the right combination of data format and position encodings, we show for the first time that standard Transformers can extrapolate to a sequence length that is 2.5× the input length. Nevertheless, unlike in-distribution generalization, length generalization remains fragile, significantly influenced by factors like random weight initialization and training data order, leading to large variances across different random seeds.

## 1. Introduction

Transformer-based models have revolutionized natural language understanding and generation across diverse applications (Gemini et al., 2023; OpenAI, 2023). Despite their impressive abilities in mathematical reasoning (Lewkowycz et al., 2022), code synthesis (Li et al., 2022), and theorem proving (Wu et al., 2022), Transformers often struggle with length generalization, an ability that requires the model to generalize to longer sequences than seen during training (Abbe et al., 2023; Anil et al., 2022; Zhou et al., 2023). This limitation raises an essential question: do Transformers genuinely grasp the correct underlying algorithms for a given task, or are they merely resorting to superficial memorization or shortcuts that fail to scale to more complex problems (Liu et al., 2023b)?

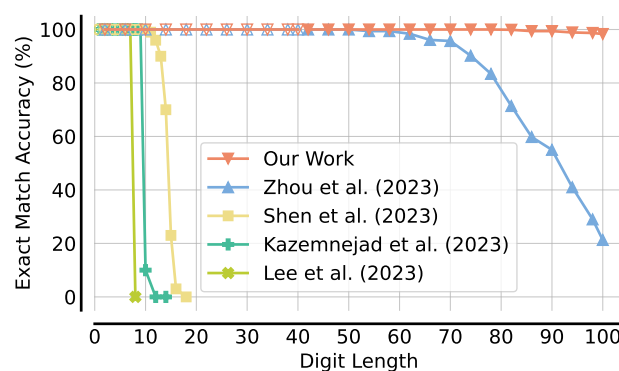


Figure 1 | Using an appropriate position encoding and data formatting, we demonstrate that Transformers can generalize to 100-digit decimal addition tasks with more than 98% of accuracy when trained up to 40-digit addition, resulting in a length extension ratio of 2.5×, which is much more than the ratio of Lee et al. (2023) (1.0×), Kazemnejad et al. (2023) (1.125×), Shen et al. (2023) (1.1×), and Zhou et al. (2023) (1.5×). Unfilled markers (▽) denote in-distribution test results, filled markers (▼) denote out-of-distribution results. In Zhou et al. (2023) and Our Work, each curve is the best out of 10 trials. For the other three methods, we report the value from their corresponding paper.

Recent work has scrutinized Transformers’ shortcomings in length generalization across formal language learning (Deletang et al., 2023) and algorithmic reasoning tasks (Anil et al., 2022; Dziri et al., 2023; Veličković et al., 2022; Zhang et al., 2022). These investigations consistently indicate a notable deficiency in length generalization capabilities. This recurring issue raises a crucial question: Is there an inherent limitation in Transformers’ design preventing effective length generalization?

In this paper, we systematically examine the Transformer’s capability of length generalization, specifically focusing on the  $N$ -digit decimal addition problem. We view the addition problem as a form of synthetic language learning, which despite its relative simplicity compared to natural language, provides valuable insights into the Transformer’s ability to internalize fundamental algorithms. Notwithstanding its simplicity, recent work has demonstrated that Transformers exhibit limited length generalization in this task (Kazemnejad et al., 2023; Lee et al., 2023; Shen et al., 2023).

Previous attempts to improve Transformer’s length generalization ability primarily focus on two areas: refining position encodings (Press et al., 2022; Shen et al., 2023) and optimizing data formats (Lee et al., 2023; Zhou et al., 2023). Therefore, we perform an extensive empirical evaluation of combinations of widely used position encoding and various data formats, resulting in a recipe for successful length generalization. Our final recipe consists of: FIRE position encodings (Li et al., 2023), with randomized positions (Ruoss et al., 2023), in reversed format, with index hints (Zhou et al., 2023).

As shown in Figure 1, when trained on only 40 digits, our model successfully extrapolates to sequences of up to 100 digits, exceeding the input length by  $2.5\times$ . To the best of our knowledge, this is the strongest known generalization result for text-based Transformers on addition. Nevertheless, we observe that the robustness of this length generalization is fragile, significantly swayed by variables such as random initialization and the training data order.

Our key contributions are summarized as follows:

- (i) We demonstrate that the success in length generalization is markedly influenced by position encoding and data format. Through careful selection of these factors, we achieved extrapolation to lengths that are  $2.5\times$  longer than those seen during training.
- (ii) Our exploration of established data formatting and augmentation techniques indicates that their effectiveness in length generalization is primarily contingent on the choice of position encoding.
- (iii) Despite remarkable generalization to lengths  $2.5\times$  longer than training, we found this generalization to be fragile and heavily relying on factors like random weight initialization and training data order.

## 2. Position Encoding and Data Formats

Recently proposed improvements in architectural design, notably in position encoding (Kazemnejad et al., 2023; Ruoss et al., 2023; Shen et al., 2023) and attention mechanisms (Duan and Shi, 2023; Dubois et al., 2019), aim to address the challenge of length generalization in arithmetic computations with Transformers. However, the effectiveness of such modifications is often constrained, either due to their overly ad-hoc nature or their poor performance on longer sequences. Although scaling the size of models and datasets has been recognized as a generally effective strategy to improve performance, prior research (Anil et al., 2022; Brown et al., 2020) suggests that relying solely on scale might not be sufficient for handling test sequences that are longer than training. Concurrently, with the rising focus on data-centric AI (Motamedi et al., 2021), recent work has investigated refining the data format to enhance the learning efficacy of existing Transformer models. In this section, we review some of the most common position encodings (Section 2.1) and relevant data formats (Section 2.2)

## 2.1. Position Encoding for Length Generalization

The inability of transformers to extrapolate to longer sequences has been primarily attributed to position encoding (PE; [Shaw et al., 2018](#)). In this section, we review existing positional encoding approaches with an emphasis on their length generalization abilities.

**Absolute Positional Encoding (APE).** APE enhances Transformer models with positional information by attaching a positional vector  $\mathbf{p}_i$  to each position  $i$ . This is achieved through a predefined sinusoidal function ([Vaswani et al., 2017](#)) or a learnable approach ([Devlin et al., 2018](#)). Then, the vector  $\mathbf{p}_i$  is combined with the token embedding  $\mathbf{e}_i$  before entering the transformer’s first layer. Although straightforward, APE often struggles with generalizing to longer sequences, as observed in both NLP ([Press et al., 2022](#)) and algorithmic tasks ([Kazemnejad et al., 2023](#)).

**Additive Relative Positional Encoding (RPE).** [Shaw et al. \(2018\)](#) pioneered the additive RPEs, diverging from standard input-level integration by modifying keys and, optionally, values in each attention layer. This concept was advanced by T5, which employed scalar biases to directly affect pre-softmax attention logits, a method noted for its simplicity yet criticized for limited efficiency and positional differentiation in long sequences ([Press et al., 2022](#); [Raffel et al., 2020](#)). Later approaches such as Alibi ([Press et al., 2022](#)), Kerple ([Chi et al., 2022](#)) and FIRE ([Li et al., 2023](#)) build on the idea of learned additive bias, proposing different functions to model the scalar bias as a function of the key- and query-indices. Most pre-softmax attention logits of additive RPEs can be generally written as ([Li et al., 2023](#)):

$$\mathbf{A}_{\text{RPE}}(\mathbf{X}) = \mathbf{X}\mathbf{W}_Q(\mathbf{X}\mathbf{W}_K)^\top + \mathbf{B}, \quad (1)$$

where  $\mathbf{X}$ ,  $\mathbf{W}_Q$ ,  $\mathbf{W}_K$  denote the input and weight matrices for queries and keys. The bias matrix  $\mathbf{B} \in \mathbb{R}^{n \times n}$  is induced by the **position encoding function**  $b : \mathbb{N}^{*2} \rightarrow \mathbb{R}$ , with its  $(i, j)$ -th entry defined as  $b(i, j)$ . Instances of  $b(i, j)$  include:

- T5 ([Raffel et al., 2020](#)):  $b(i, j) = r_{\min}\{i - j, K\}$ , where  $K$  is a hyperparameter and  $r_i$  are learned scalars.
- Alibi ([Press et al., 2022](#)):  $b(i, j) = -r|i - j|$ , where  $r > 0$  is a hyperparameter.
- KerpleLog ([Chi et al., 2022](#)):  $b(i, j) = -r_1 \log(1 + r_2|i - j|)$ , where  $r_1, r_2 > 0$  are learnable scalars.
- FIRE ([Li et al., 2023](#)):  $b(i, j) = f_\theta\left(\frac{\psi(i-j)}{\psi(\max\{L, i\})}\right)$ , where  $f_\theta : \mathbb{R} \rightarrow \mathbb{R}$  is a learnable MLP parameterized by  $\theta$ ,  $\psi : \mathbb{N} \rightarrow \mathbb{R}_+$  is  $\psi(x) = \log(cx + 1)$  and  $c > 0, L > 0$  are learnable scalars.

Additional background on additive RPEs is provided in [Appendix A.1](#)

**Rotary Positional Encoding (RoPE).** RoPE ([Su et al., 2024](#)) encodes position information in attention logits through rotational encoding of query and key vectors based on their relative positions. Despite being simple and effective, RoPE exhibits limited length generalization ([Kazemnejad et al., 2023](#); [Press et al., 2022](#)). While extensions like Position Interpolation [Chen et al. \(2023\)](#); [Peng et al. \(2023\)](#); [Su \(2023\)](#) enhance RoPE’s context length, they do not necessarily improve length generalization on algorithmic tasks where learning the underlying algorithm is crucial.

**No Positional Encoding (NoPE).** While encoder-only Transformers (e.g., BERT ([Devlin et al., 2018](#))) are permutation equivariant without positional encodings, decoder-only counterparts with causal attention, as shown by [Haviv et al. \(2022\)](#), acquire positional understanding autonomously, even without explicit PE. Interestingly, recent findings by [Kazemnejad et al. \(2023\)](#) further reveal that a model without PE outperforms those with specialized PEs on simple algorithmic tasks.

**Randomized Position Encoding.** [Ruoss et al. \(2023\)](#) introduced Randomized PE to enhance existing PEs by randomly sampling encodings from a range exceeding test-time lengths while pre-

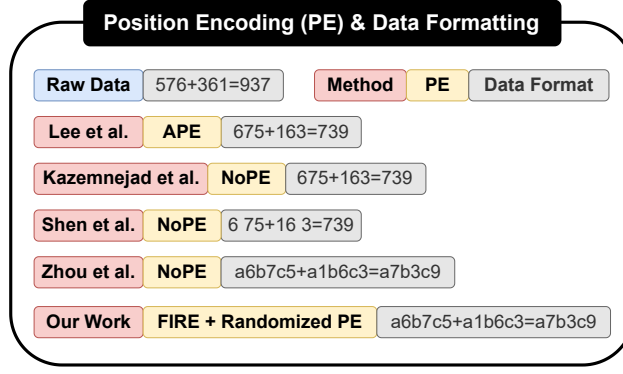


Figure 2 | Comparative overview of PEs and data formats: While most related studies focus on APE or NoPE, our approach integrates FIRE (Li et al., 2023) and Randomized PE (Ruoss et al., 2023). All studies utilize a reversed format. Shen et al. (2023) enhance this with random space augmentation, and both Zhou et al. (2023) and Our Work incorporate index hints.

serving the order. Transformers trained this way adapt to larger positional encodings, effectively eliminating OOD position encodings during testing.

## 2.2. Data Formats

Data format plays a pivotal role in enhancing Transformers’ length generalization capabilities, primarily by transforming the data into a format that could be more easily learned. We give an overview of the existing techniques below.

**Reversed Format.** Computing addition in an algorithmic way (as taught in elementary school) requires starting with the least significant digit (LSD) and proceeds to the most significant digit (MSD). This sequence contrasts with the standard printed format ( $A_3A_2A_1 + B_3B_2B_1 = C_3C_2C_1$ , where  $A_1$  and  $B_1$  are the LSDs, which is not ideally suited for autoregressive models due to their outputting the MSD first. However, the reversed format ( $A_1A_2A_3 + B_1B_2B_3 = C_1C_2C_3$ ) aligns better with these the natural order of computing the digits. It simplifies the learning task to a function that depends only on the two corresponding operand digits and the carry from the previous step (Lee et al., 2023; Shen et al., 2023; Zhou et al., 2023).

**Index Hints.** Zhou et al. (2023) introduced “index hints” in both the query and response of arithmetic tasks. For example,  $42 + 39 = 81$  is represented as  $a4b2 + a3b9 = a8b1$  during training and inference, enabling transformers to execute indexing via induction heads (Olsson et al., 2022).

**Random Space Augmentation.** Shen et al. (2023) explored the impact of random spacing between digits in addition, aiming to disrupt the model’s reliance on absolute positional information. Their results show successful generalization from 10-digit to 11-digit addition, but falters with longer sequences.

Figure 2 lists the position encodings and data formats used in some of the most related work to ours.

## 3. A Recipe for Length Generalization in Decimal Addition

The task of decimal addition is composed of two critical subtasks: (a) the identification of the right operands to add; and (b) the summation of these operands with the preceding carry. While the

summation step ((b)) is relatively easier because it has a finite set of possible inputs, the primary generalization challenge lies in the operand identification ((a)), where precise positional access is crucial.

Our best model, which leads to the results in Figure 1, uses the following combination:

1. **FIRE position encodings** (Li et al., 2023): We believe that FIRE position encodings are helpful for length generalization because they are more expressive than other PEs, as shown by Li et al. (2023).
2. **Randomized position encodings** (Ruoss et al., 2023): We believe that randomized position encodings are crucial to avoid overfitting on the position indices and index differences that were seen during training.
3. **Reversed format**: The reversed format makes it easier for the model to *decompose* the long computation to local, “markovian”, steps that depend only on the single previous step.
4. **Index hints** (Zhou et al., 2023): We believe that index hints are useful because they ease the task of *operand identification* (discussed in (b)), of matching the right operands to add at a certain step.

We ablate each of these decisions and some other alternative choices in Section 4.

## 4. Experiments

### 4.1. Setup

**Data.** As shown in Figure 2, we adopt the reversed format with index hints as our default data format. During training, we randomly sample consecutive index hints from a pre-defined ordered set of hints with 102 symbols, thereby enhancing the learning of hint sequences and their order. We generated a dataset comprising 30M examples on input lengths 1-40 for training and 1,000 examples per input length for testing.

**Model.** Our base model, following Zhou et al. (2023), is a 25M parameter Transformer featuring 6 blocks, a 512 hidden size, and a feedforward layer with a hidden dimension of 2048. We also adopt RMSNorm, integrating both PreNorm and PostNorm layers, following the Primer architecture (So et al., 2021). We use the AdamW optimizer (Loshchilov and Hutter, 2017) to train the model with a weight decay value of 0.1 and no dropout, for 50,000 steps. The learning rate schedule incorporates an initial 500-step linear warm-up, followed by a cosine decay, starting at  $3e-4$ . The hyperparameters are chosen based on Appendix C.10.

**Randomized PE and Random Space Augmentation.** As will be demonstrated in Figures 7 and 8, the success of these techniques is markedly PE-dependent. Hence, we tailor the default hyperparameter choice to best suit each PE. Further, instead of using random spaces, we use another special token to prevent automatic merging by the tokenizer.

Due to the high variance (which we discuss in the next section), we repeat each experiment five times unless mentioned otherwise. More implementation details are provided in Appendix B.

### 4.2. Results

**FIRE enables significantly better length generalization.** Figure 3 compares the length generalization capabilities of four positional encodings in the best of 10 trials (See Appendix C.1 for all trials). Trained exclusively on sequences of lengths 1-40, the best trial of FIRE exhibit near-perfect

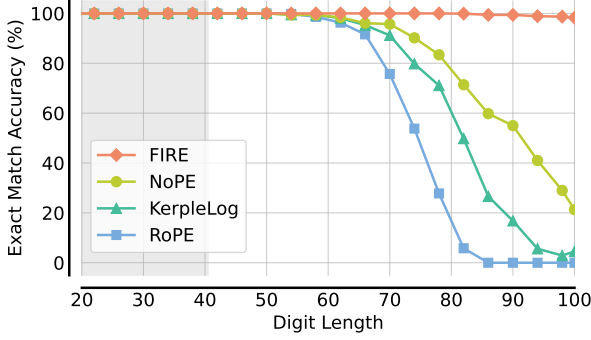


Figure 3 | EM accuracy (best of 10 trials), trained exclusively on sequences of lengths 1 to 40, the best trials involving FIRE exhibit near-perfect generalization on 100-digit addition.

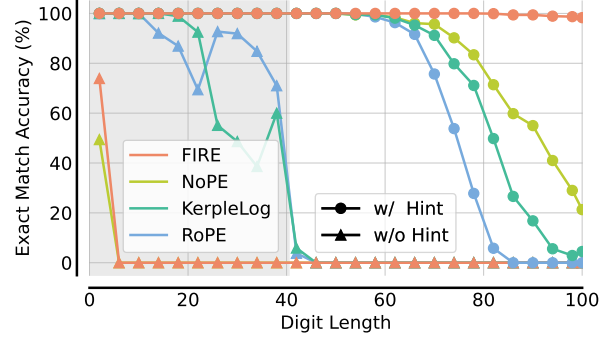


Figure 4 | EM accuracy of models trained with and without index hints (best of 10 trials): Without index hints, all PE methods fail in generalization, both within and beyond trained lengths.

generalization to sequences up to the length of 100. In contrast, other PEs show a visible degradation in generalization accuracy beyond the sequence length of 60. This finding counters the findings of [Kazemnejad et al. \(2023\)](#) that no positional encoding (NoPE) surpasses complex PE techniques for length generalization. Our findings suggest that a well-designed PE, such as FIRE, is essential for optimal length generalization.

**Index hints are crucial.** We compare models trained with and without index hints. As shown in Figure 4, index hints significantly enhance length generalization across various PEs, corroborating the findings of [Zhou et al. \(2023\)](#). Notably, without index hints, NoPE and FIRE demonstrate poor in-distribution generalization for 40-digit additions, a marked deviation from their reasonable performance when trained on 10-digits, as shown in Figure C.8(a). Figure D.1 shows that this phenomenon occurs across all random seeds. Conversely, RoPE and KerpleLog exhibit moderate in-distribution generalization but falter in out-of-distribution scenarios. Appendices D.1 and D.2 shows the training loss and test accuracy of these runs.

Analyzing errors in 11-digit additions from models trained on 10-digits revealed a common misalignment issue: the Transformer often adds operands adjacent to the correct ones. An attempt to rectify this by reformatting addition ( $A_1B_1, A_2B_2, A_3B_3 = C_1C_2C_3$ , with 1 as the least significant bit) failed to improve length generalization, merely shifting the error to adjacent output positions. This highlights the Transformer’s inherent limitations in precise position identification.

**Standard format vs reversed format.** As shown in Figure 5, standard formatting shows limited length generalization in all PEs compared to the reversed format. FIRE excels in length generalization even with the standard format, even matching RoPE in reverse format. However, FIRE’s performance (with standard format) declines beyond 60-digit additions, likely due to increased carry propagation challenges exceeding the model’s capacity.

Looking at the training loss and training next-token accuracy in both formats also shows interesting differences. As shown in Figures 6 and C.3, the standard format training leads to gradual improvement, whereas reverse format yields a sharp performance transition. This transition, which is a reminiscent of “grokking” phenomenon [Power et al. \(2022\)](#), shows in this case the “Eureka moment” in which the Transformer learns the right addition algorithm.

**Random space augmentation and randomized position encoding.** Figure 7 reveals divergent impacts of random space augmentation on four PEs. The augmentation’s efficacy is notably contingent upon the chosen PE. While Random Spaces marginally enhances RoPE and KerpleLog’s performance,



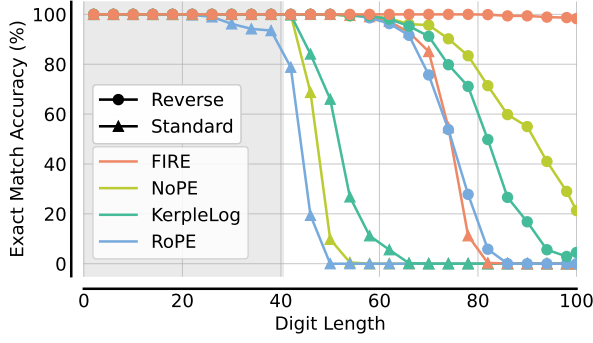


Figure 5 | EM accuracy of the standard vs. the reversed format: Consistently with prior studies, the reversed format excels over the standard format across all PEs.

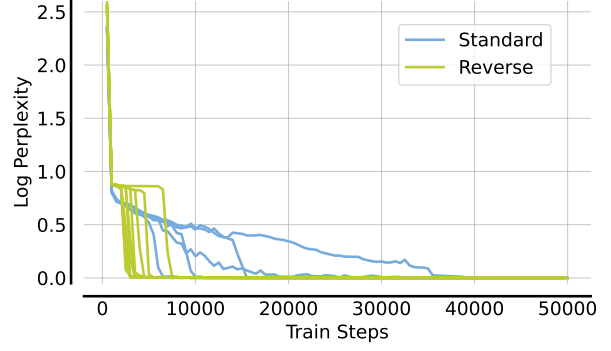


Figure 6 | The reversed format shows distinct grokking during training, unlike the gradual enhancement in the standard format. This phenomenon is observed across all PEs (Figure C.3)

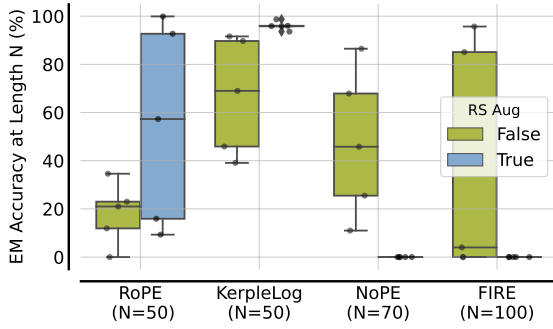


Figure 7 | Effects of Random Space Augmentation (RS Aug): Random space augmentation is beneficial for RoPE and KerpleLog; adverse for NoPE and FIRE.

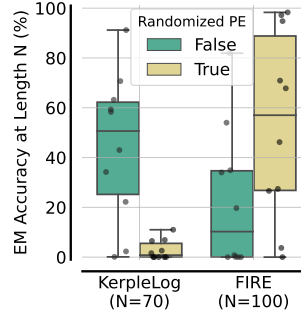


Figure 8 | Effects of Randomized PE: Randomized PE enhances FIRE but degrades KerpleLog

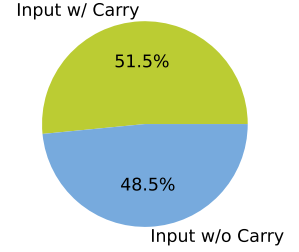


Figure 9 | Error Distribution: Errors appear almost equally with and without carry.

it markedly deteriorates NoPE and FIRE. A similar PE-specific pattern is evident in Randomized PE, as Figure 8 demonstrates. Randomized PE significantly degrades KerpleLog’s effectiveness, yet it substantially boosts FIRE. See Appendices D.4 and D.5 for training loss and EM accuracy for all trials in each setting.

**Length generalization is not robust to neither weight initialization nor training data order.** Figure 10 illustrates the varying performance of 10 FIRE trials using identical training data order but distinct weight initializations. Notably, while all trials achieve similar close-to-zero training losses after 10K training steps (Figure C.2) and exhibit perfect in-distribution generalization, their out-of-distribution (OOD) length generalization shows significant variance. Moreover, the length generalization performance fluctuates significantly across training steps (Appendix C.3). This observation contrasts with earlier studies suggesting in-distribution loss as a reliable OOD generalization predictor (Nagarajan et al., 2020).

We further examine 15 unique combinations, resulting from 3 weight initialization seeds and 5 data input orders. As shown in Figure 11, there is significant variance across training data orders even when the weight initialization is constant. Intriguingly, certain weight initializations demonstrate remarkable resilience to changes in data input order. This observation is reminiscent of the Lottery Ticket Hypothesis (Frankle and Carbin, 2018), which posits the existence of a sparse, equally effective sub-network within a larger neural network. Our findings suggest the presence of “fortunate” weight

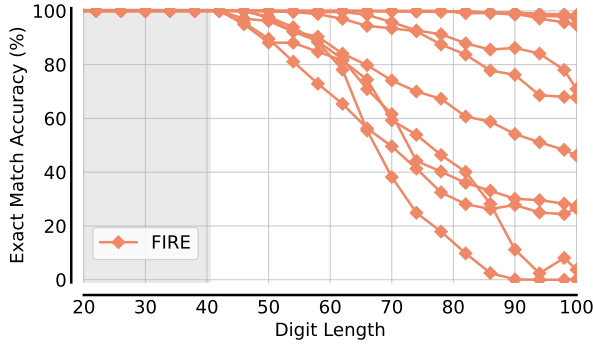


Figure 10 | Exact match across 10 trials using FIRE. While transformers can achieve near-perfect accuracy in 100-digit addition, the variance across different random seeds is high.

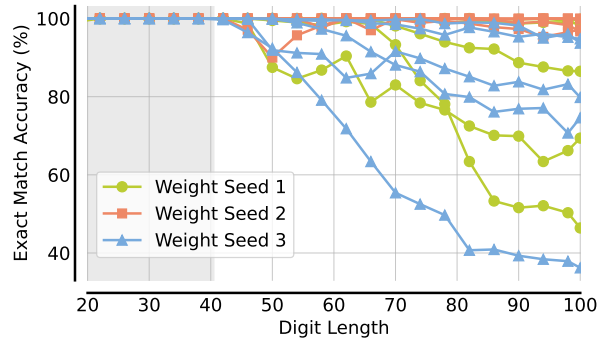


Figure 11 | Effects of weight initialization and data input order: 15 models trained on a combination of three weight initialization seeds and five data input order seeds.

configurations that exhibit robust length generalization, akin to a “lucky weight ticket.”

While Anil et al. (2022) also noticed similar in-distribution accuracy but marked differences in OOD behavior on parity tasks, their OOD performance was quite poor across all runs. Moreover, contrary to the findings of Anil et al. (2022) on the impact of hyperparameter variation, our experiments reveal considerable performance fluctuations even with different random seeds. This inconsistency appears unrelated to position encoding (refer to Figure C.1 for different PEs), and is more likely due to variations in random weight initialization and data order.

## 5. Analysis

**Error analysis.** In examining Transformers’ error characteristics, we classified erroneous predictions into two categories: those with and without carry. Figure 9 shows no significant difference between these categories, thus carry propagation does not majorly impede length generalization.

Additionally, we analyzed the error distribution in 100-digit addition using FIRE, illustrated in Figure C.10. As shown, Figure C.10 indicates an overall uniform error distribution across all indices, despite some individual model checkpoints showing errors at specific positions. Excluding two near-zero accuracy runs, over 90% of errors in incorrect examples are single-digit mistakes, following an exponential distribution. Additional results are shown in Figures C.11 and C.12.

Despite the imperfect calculation, the FIRE model does not show any systematic error. Random errors may stem from phenomena such as attention glitches Liu et al. (2023a). Conversely, other PEs systematically fail to identify the start or end of addition, leading to premature termination.

**Performance evolution during training.** Figure 12 shows that while transformers achieve near-perfect in-distribution accuracy early in training, they explore different extrapolation strategies. This ability is remarkable considering the inherent unpredictability and architecture-dependent nature of OOD accuracy. Notably, transformers with FIRE exhibit a generally steady increase in OOD accuracy during training, suggesting that FIRE’s inductive bias may be helpful in finding solutions that generalize to different lengths. In contrast, other PE methods display more volatile OOD performance. Interestingly, some methods exhibit a “grokking-like” phenomenon, where there is a sudden surge in the OOD accuracy despite no change in in-distribution accuracy.

**Sequence length during training.** We trained separate models for addition involving up to 10,



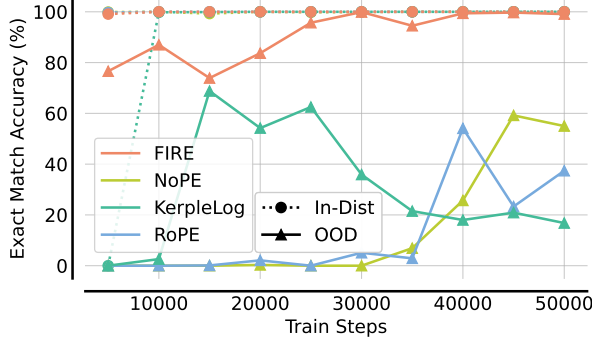


Figure 12 | Comparison of In-Distribution (30-digit addition) and Out-of-Distribution Generalization (90-digit addition, except for RoPE at 70-digit addition).

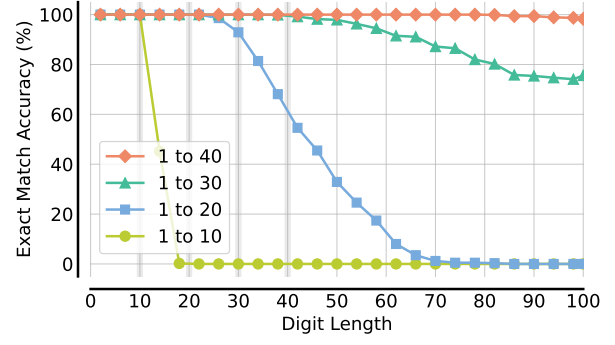


Figure 13 | Different training lengths: Increasing the training length significantly improves length generalization in FIRE, achieving near-perfect accuracy at length 100.

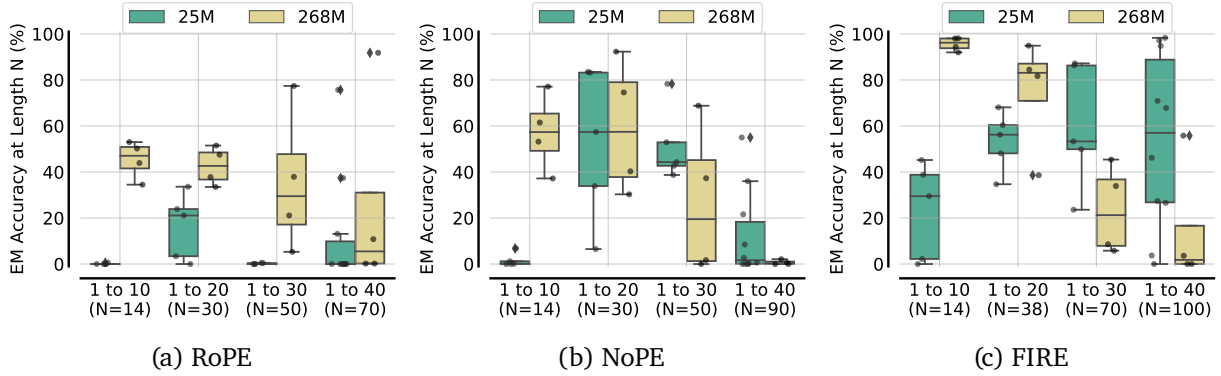


Figure 14 | Scaling model size inconsistently affects length generalization performance. While consistently enhancing performance in shorter length regimes (1-10, 1-20) across four PEs, this trend does not hold for larger regimes (1-30, 1-40). For instance, larger models outperform smaller ones with RoPE and KerpleLog (Figure C.14), but underperform with NoPE and FIRE. Moreover, increasing model size doesn't noticeably decrease performance variance, suggesting size scaling isn't vital for length generalization.

20, 30, and 40 digits, and evaluated them on addition of up to 100 digits. As depicted in Figures 13 and C.13, training length crucially improves performance in longer length generalizations across different PEs. Notably, not only that models that were trained on 40 digits generalize better than models that were trained on shorter sequences, the *generalization factor is also increasing*: the model that was trained on 40 digits generalizes to 100 digits (2.5 $\times$ ), while the model that was trained on up to 30 digits generalizes to 45 digits (1.5 $\times$ ), the model that was trained on up to 20 digits generalizes to 25 digits (1.25 $\times$ ), and the model that was trained on up to 10 digits does not generalize beyond training lengths (1.0 $\times$ ).

**Scaling model size.** The scaling of model size is crucial for improving large language models (Chowdhery et al., 2023; Thoppilan et al., 2022). To assess its effect on length generalization, we contrasted models with 25M and 268M parameters. We find that model size variation has a minor effect on length generalization. Figure 14 shows that larger models slightly improve generalization in short digit regimes (1 to 10 and 1 to 20 digit additions) but yield mixed results in longer regimes. While RoPE and KerpleLog show improvements, NoPE and FIRE experience performance degradation with a larger model, indicating model size may not be the primary factor in length generalization.

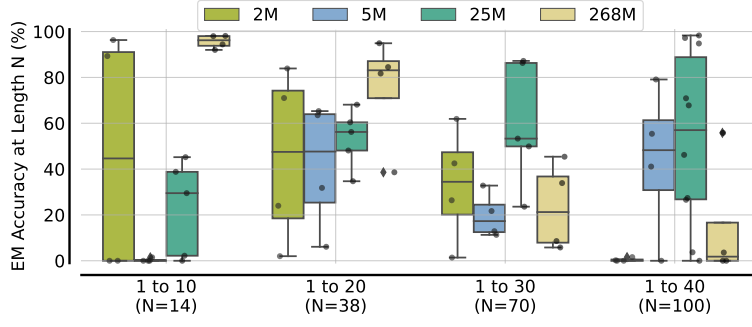


Figure 15 | Effect of different model sizes with FIRE as the position encoding.

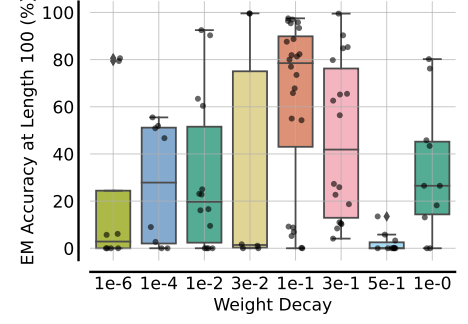


Figure 16 | Effect of weight decay with FIRE as the position encoding.

The efficacy of length generalization in the 25M model prompted us to explore the capabilities of smaller models. Specifically, we trained models with 2M and 5M parameters. As Figures 15 and C.15 illustrate, the 2M model’s performance deteriorates with longer sequences, indicating limited model capacity as a potential performance bottleneck. Intriguingly, this model outperforms its larger counterparts (5M and 25M models) in tasks involving 1 to 10 digit addition. Furthermore, the 5M model remarkably achieves 80% accuracy in 100 digit addition, trained only on 1 to 40 digit tasks, surpassing the 268M model’s performance.

**Does stronger regularization reduce variance?** To mitigate performance variance, we investigated standard regularization techniques, including weight decay and dropout. As depicted in Figure 16, higher weight decay values (e.g., 0.1, 0.3) slightly enhance the likelihood of achieving effective length generalization. Nonetheless, non-trivial length generalization remains attainable with either very low (e.g.,  $1e-6$ ) or high (e.g., 1.0) weight decay values, evidenced by approximately 80% accuracy in 100 digit addition trained on 40-digit sequences. Conversely, Figure C.17 shows that substantial dropout values (e.g., 0.2) severely impair length generalization. Dropout rates of 0.0 or 0.1, however, do not show statistically significant improvements over their counterparts. Overall, while regularization can modestly decrease performance variability, it falls short in ensuring robust length generalization. The variance in performance is still significantly influenced by the randomness of weights initialization and the training data order (Figures 10 and 11).

## 6. Related Work

Length generalization remains a significant challenge in neural networks, underscored by substantial research (Deletang et al., 2023; Dziri et al., 2023; Graves et al., 2016; Hupkes et al., 2020; Schwarzschild et al., 2021; Zhang et al., 2022). Despite their advanced reasoning capabilities, Transformer-based large language models (LLMs) (Chowdhery et al., 2023; Thoppilan et al., 2022) struggle with processing sequences beyond their training scope Anil et al. (2022). Enhancements in length generalization, especially in the addition task, primarily focus on two areas: refining positional encoding and optimizing data format.

**Position Encoding for Length Generalization** The inability of Transformers to extrapolate to longer sequences has been primarily attributed to Position Encoding (PE) Shaw et al. (2018). Various studies have suggested alternatives, such as relative positional encodings, which focus on the relative distances between tokens (Dai et al., 2019), the implementation of randomized position encoding (Ruoss et al., 2023), or the adoption of weighted attention mechanisms in place of position embeddings (Chi et al., 2022; Li et al., 2023; Press et al., 2022; Raffel et al., 2020). These approaches have shown promise in natural language processing (NLP). However, Kazemnejad et al. (2023) found

that omitting position encoding entirely yields better results for algorithmic tasks. In contrast, our experiments indicate that an effectively designed PE, such as the FIRE, is crucial for achieving optimal length generalization (Figure 3). Moreover, we show that a synergistic approach to consider both PE and data design markedly enhances length generalization capabilities.

**Data format for Length Generalization** A range of heuristic-based data formatting methods have been introduced, particularly for pretrained LLMs. These methods, including the use of scratchpads and the chain of thoughts approach, aim to facilitate arithmetic learning either through in-context learning or fine-tuning [Anil et al. \(2022\)](#); [Zhou et al. \(2022\)](#). Conversely, there is a body of research focused on Transformers trained from scratch. This research indicates that employing techniques such as reversed formatting and scratch pads can significantly boost length generalization performance [Lee et al. \(2023\)](#); [Shen et al. \(2023\)](#). Furthermore, it has been observed that both the data distribution and the sampling strategies can profoundly influence generalization [Lee et al. \(2023\)](#). [Awasthi and Gupta \(2023\)](#) further demonstrates the benefits of incorporating a simpler auxiliary task (e.g., identifying the successor element) in supporting the primary task (e.g., sorting). In contrast, [Jelassi et al. \(2023\)](#) finds that train set priming enables length generalization for an encoder-only Transformer model. In contrast, our good length generalization performance achieved with naive random sampling approach suggesting that sophisticated data sampling might be redundant.

## 7. Conclusion

Length generalization in Transformers has been a long-standing challenge. We evaluate the ability of Transformers to generalize to longer test sequences using the decimal addition task. Through extensive experiments, we find that there is no inherent limitation in Transformers’ design preventing effective length generalization. Instead, the missing ingredient is the right combination of data format and position encoding. We demonstrate that Transformers can achieve almost perfect generalization on sequences up to  $2.5\times$  the training length, given appropriate data formatting and position encoding.

Our thorough empirical analysis of common length generalization techniques reveals a significant dependency between the type of position encoding and the data format. This underscores the importance of synergizing data format with model architecture for optimal generalization. Despite these advancements, robust length generalization in Transformers remains elusive, even with meticulously finetuned regularization hyperparameters.

## References

- E. Abbe, S. Bengio, A. Lotfi, and K. Rizk. Generalization on the unseen, logic reasoning and degree curriculum. *arXiv preprint arXiv:2301.13105*, 2023.
- C. Anil, Y. Wu, A. Andreassen, A. Lewkowycz, V. Misra, V. Ramasesh, A. Slone, G. Gur-Ari, E. Dyer, and B. Neyshabur. Exploring length generalization in large language models. *Advances in Neural Information Processing Systems*, 35:38546–38556, 2022.
- P. Awasthi and A. Gupta. Improving length-generalization in transformers via task hinting. *arXiv preprint arXiv:2310.00726*, 2023.
- T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- S. Chen, S. Wong, L. Chen, and Y. Tian. Extending context window of large language models via positional interpolation. *arXiv preprint arXiv:2306.15595*, 2023.

- T.-C. Chi, T.-H. Fan, P. J. Ramadge, and A. Rudnicky. Kerple: Kernelized relative positional embedding for length extrapolation. *Advances in Neural Information Processing Systems*, 35:8386–8399, 2022.
- A. Chowdhery, S. Narang, J. Devlin, M. Bosma, G. Mishra, A. Roberts, P. Barham, H. W. Chung, C. Sutton, S. Gehrmann, et al. Palm: Scaling language modeling with pathways. *Journal of Machine Learning Research*, 24(240):1–113, 2023.
- Z. Dai, Z. Yang, Y. Yang, J. Carbonell, Q. V. Le, and R. Salakhutdinov. Transformer-xl: Attentive language models beyond a fixed-length context. *arXiv preprint arXiv:1901.02860*, 2019.
- G. Deletang, A. Ruoss, J. Grau-Moya, T. Genewein, L. K. Wenliang, E. Catt, C. Cundy, M. Hutter, S. Legg, J. Veness, and P. A. Ortega. Neural networks and the chomsky hierarchy. In *The Eleventh International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=WbxHAzkeQcn>.
- J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- S. Duan and Y. Shi. From interpolation to extrapolation: Complete length generalization for arithmetic transformers. *arXiv preprint arXiv:2310.11984*, 2023.
- Y. Dubois, G. Dagan, D. Hupkes, and E. Bruni. Location attention for extrapolation to longer sequences. *arXiv preprint arXiv:1911.03872*, 2019.
- N. Dziri, X. Lu, M. Sclar, X. L. Li, L. Jian, B. Y. Lin, P. West, C. Bhagavatula, R. L. Bras, J. D. Hwang, et al. Faith and fate: Limits of transformers on compositionality. *arXiv preprint arXiv:2305.18654*, 2023.
- J. Frankle and M. Carbin. The lottery ticket hypothesis: Finding sparse, trainable neural networks. *arXiv preprint arXiv:1803.03635*, 2018.
- Gemini, R. Anil, S. Borgeaud, Y. Wu, J.-B. Alayrac, J. Yu, R. Soricut, J. Schalkwyk, A. M. Dai, A. Hauth, et al. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.
- A. Graves, G. Wayne, M. Reynolds, T. Harley, I. Danihelka, A. Grabska-Barwinska, S. G. Colmenarejo, E. Grefenstette, T. Ramalho, J. P. Agapiou, A. P. Badia, K. M. Hermann, Y. Zwols, G. Ostrovski, A. Cain, H. King, C. Summerfield, P. Blunsom, K. Kavukcuoglu, and D. Hassabis. Hybrid computing using a neural network with dynamic external memory. *Nat.*, 538(7626):471–476, 2016. doi: 10.1038/NATURE20101. URL <https://doi.org/10.1038/nature20101>.
- A. Haviv, O. Ram, O. Press, P. Izsak, and O. Levy. Transformer language models without positional encodings still learn positional information. *arXiv preprint arXiv:2203.16634*, 2022.
- D. Hupkes, V. Dankers, M. Mul, and E. Bruni. Compositionality decomposed: How do neural networks generalise? *Journal of Artificial Intelligence Research*, 67:757–795, 2020.
- S. Jelassi, S. d’Ascoli, C. Domingo-Enrich, Y. Wu, Y. Li, and F. Charton. Length generalization in arithmetic transformers. *arXiv preprint arXiv:2306.15400*, 2023.
- A. Kazemnejad, I. Padhi, K. N. Ramamurthy, P. Das, and S. Reddy. The impact of positional encoding on length generalization in transformers. *arXiv preprint arXiv:2305.19466*, 2023.
- N. Lee, K. Sreenivasan, J. D. Lee, K. Lee, and D. Papailiopoulos. Teaching arithmetic to small transformers. *arXiv preprint arXiv:2307.03381*, 2023.

- A. Lewkowycz, A. Andreassen, D. Dohan, E. Dyer, H. Michalewski, V. Ramasesh, A. Slone, C. Anil, I. Schlag, T. Gutman-Solo, et al. Solving quantitative reasoning problems with language models. *Advances in Neural Information Processing Systems*, 35:3843–3857, 2022.
- S. Li, C. You, G. Guruganesh, J. Ainslie, S. Ontanon, M. Zaheer, S. Sanghai, Y. Yang, S. Kumar, and S. Bhojanapalli. Functional interpolation for relative positions improves long context transformers. *arXiv preprint arXiv:2310.04418*, 2023.
- Y. Li, D. Choi, J. Chung, N. Kushman, J. Schrittwieser, R. Leblond, T. Eccles, J. Keeling, F. Gimeno, A. Dal Lago, et al. Competition-level code generation with alphacode. *Science*, 378(6624):1092–1097, 2022.
- B. Liu, J. T. Ash, S. Goel, A. Krishnamurthy, and C. Zhang. Exposing attention glitches with flip-flop language modeling. *arXiv preprint arXiv:2306.00946*, 2023a.
- B. Liu, J. T. Ash, S. Goel, A. Krishnamurthy, and C. Zhang. Transformers learn shortcuts to automata. In *The Eleventh International Conference on Learning Representations*, 2023b. URL <https://openreview.net/forum?id=De4FYqjFueZ>.
- I. Loshchilov and F. Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017.
- M. Motamedi, N. Sakharlykh, and T. Kaldewey. A data-centric approach for training deep neural networks with less data. *arXiv preprint arXiv:2110.03613*, 2021.
- V. Nagarajan, A. Andreassen, and B. Neyshabur. Understanding the failure modes of out-of-distribution generalization. *arXiv preprint arXiv:2010.15775*, 2020.
- C. Olsson, N. Elhage, N. Nanda, N. Joseph, N. DasSarma, T. Henighan, B. Mann, A. Askell, Y. Bai, A. Chen, et al. In-context learning and induction heads. *arXiv preprint arXiv:2209.11895*, 2022.
- OpenAI. Gpt-4 technical report. *ArXiv*, abs/2303.08774, 2023. URL <https://api.semanticscholar.org/CorpusID:257532815>.
- B. Peng, J. Quesnelle, H. Fan, and E. Shippole. Yarn: Efficient context window extension of large language models. *arXiv preprint arXiv:2309.00071*, 2023.
- A. Power, Y. Burda, H. Edwards, I. Babuschkin, and V. Misra. Grokking: Generalization beyond overfitting on small algorithmic datasets. *arXiv preprint arXiv:2201.02177*, 2022.
- O. Press, N. Smith, and M. Lewis. Train short, test long: Attention with linear biases enables input length extrapolation. In *International Conference on Learning Representations*, 2022. URL <https://openreview.net/forum?id=R8sQPpGCv0>.
- C. Raffel, N. Shazeer, A. Roberts, K. Lee, S. Narang, M. Matena, Y. Zhou, W. Li, and P. J. Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *The Journal of Machine Learning Research*, 21(1):5485–5551, 2020.
- A. Ruoss, G. Delétang, T. Genewein, J. Grau-Moya, R. Csordás, M. Bennani, S. Legg, and J. Veness. Randomized positional encodings boost length generalization of transformers. *arXiv preprint arXiv:2305.16843*, 2023.
- A. Schwarzschild, E. Borgnia, A. Gupta, F. Huang, U. Vishkin, M. Goldblum, and T. Goldstein. Can you learn an algorithm? generalizing from easy to hard problems with recurrent networks. *Advances in Neural Information Processing Systems*, 34:6695–6706, 2021.

- P. Shaw, J. Uszkoreit, and A. Vaswani. Self-attention with relative position representations. *arXiv preprint arXiv:1803.02155*, 2018.
- N. Shazeer. Glu variants improve transformer. *arXiv preprint arXiv:2002.05202*, 2020.
- R. Shen, S. Bubeck, R. Eldan, Y. T. Lee, Y. Li, and Y. Zhang. Positional description matters for transformers arithmetic. *arXiv preprint arXiv:2311.14737*, 2023.
- D. R. So, W. Mańke, H. Liu, Z. Dai, N. Shazeer, and Q. V. Le. Primer: Searching for efficient transformers for language modeling. *arXiv preprint arXiv:2109.08668*, 2021.
- J. Su. Rectified rotary position embeddings. <https://github.com/bojone/rerope>, 2023.
- J. Su, M. Ahmed, Y. Lu, S. Pan, W. Bo, and Y. Liu. Roformer: Enhanced transformer with rotary position embedding. *Neurocomputing*, 568:127063, 2024.
- R. Thoppilan, D. De Freitas, J. Hall, N. Shazeer, A. Kulshreshtha, H.-T. Cheng, A. Jin, T. Bos, L. Baker, Y. Du, et al. Lamda: Language models for dialog applications. *arXiv preprint arXiv:2201.08239*, 2022.
- A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- P. Veličković, A. P. Badia, D. Budden, R. Pascanu, A. Banino, M. Dashevskiy, R. Hadsell, and C. Blundell. The clrs algorithmic reasoning benchmark. In *International Conference on Machine Learning*, pages 22084–22102. PMLR, 2022.
- Y. Wu, A. Q. Jiang, W. Li, M. Rabe, C. Staats, M. Jamnik, and C. Szegedy. Autoformalization with large language models. *Advances in Neural Information Processing Systems*, 35:32353–32368, 2022.
- Y. Zhang, A. Backurs, S. Bubeck, R. Eldan, S. Gunasekar, and T. Wagner. Unveiling transformers with lego: a synthetic reasoning task. *arXiv preprint arXiv:2206.04301*, 2022.
- H. Zhou, A. Nova, H. Larochelle, A. Courville, B. Neyshabur, and H. Sedghi. Teaching algorithmic reasoning via in-context learning. *arXiv preprint arXiv:2211.09066*, 2022.
- H. Zhou, A. Bradley, E. Littwin, N. Razin, O. Saremi, J. Susskind, S. Bengio, and P. Nakkiran. What algorithms can transformers learn? a study in length generalization. *arXiv preprint arXiv:2310.16028*, 2023.



## A. Positional Encoding

### A.1. Additive Relative Positional Encoding (RPE)

Shaw et al. (2018) pioneered additive RPE by integrating position encodings into the attention layer’s key, and optionally the value, rather than the input. This concept was further simplified in T5 (Raffel et al., 2020), where the vector representations of relative positions are simplified to scalar biases added to pre-softmax attention logits. Subsequent advancements in additive RPE, aimed at enhancing length generalization and computational efficiency, include notable methods like Alibi (Press et al., 2022), Kerple (Chi et al., 2022), and FIRE (Li et al., 2023). A commonality among these methods is the unified computation formula for pre-softmax attention logits, as outlined by Li et al. (2023):

$$A_{\text{RPE}}(X) = XW_Q(XW_K)^\top + B, \quad (\text{A.1})$$

where the bias matrix  $B \in \mathbb{R}^{n \times n}$  is induced by the **position encoding function**  $b : \mathbb{N}^{*2} \rightarrow \mathbb{R}$ , has its  $(i, j)$ -th entry defined as  $b(i, j)$ . Variations in  $b$ ’s formulations and parameterizations give rise to diverse RPE variants.

- T5 (Raffel et al., 2020): T5’s RPE segments relative distances into distinct buckets with a logarithmic scale, each associated with a unique parameter. With  $K + 1$  buckets and a pre-defined distance  $L_1$ , the attention bias is calculated as (assuming  $K + 1$  is even)

$$b(i, j) = \begin{cases} r_{i-j} & 0 \leq i - j < \frac{K+1}{2} \\ r_{\frac{K+1}{2} + \lfloor \frac{K+1}{2} \log\left(\frac{2(i-j)}{K+1}\right) / \log\left(\frac{2L_1}{K+1}\right) \rfloor} & \frac{K+1}{2} \leq i - j < L_1 \\ r_K & i - j \geq L_1 \end{cases}. \quad (\text{A.2})$$

- Alibi (Press et al., 2022):  $b(i, j) = -r|i - j|$ , where  $r > 0$  is a hyper-parameter.
- Kerple (Chi et al., 2022):  $b(i, j) = -r_1 \log(1 + r_2|i - j|)$  (logarithmic variant) or  $-r_1|i - j|^{r_2}$  (power variant), where  $r_1, r_2 > 0$  are learnable scalars.
- FIRE (Li et al., 2023):  $b(i, j) = f_\theta\left(\frac{\psi(i-j)}{\psi(\max\{L, i\})}\right)$ , where  $f_\theta : \mathbb{R} \rightarrow \mathbb{R}$  is a learnable MLP parameterized by  $\theta$ ,  $\psi : \mathbb{N} \rightarrow \mathbb{R}_+$  is monotonically increasing and  $L > 0$  is a learnable scalar.

## B. Implementation Details

### B.1. Data Generation

As shown in Figure 2, we adopt the reversed format with index hints as our default data format. During training, we randomly sample a consecutive index hints from a pre-defined ordered index set with 102 distinct symbols, thereby enhancing the learning of hint sequences and their order. At inference, the same hint sampling strategy is applied to questions, prompting the model for answers.

To generate addition examples, we opt for a naive random sampling approach instead of structured data sampling Lee et al. (2023), as our analysis indicates that carry operations are not a major hindrance to length generalization (See Figure 9). Our approach involves uniformly selecting the number’s length from 1 to the maximum training length, followed by independent sampling of two operands based on this length, with an additional zero padding to accommodate potential carry-induced extra digits. For training, datasets comprising 30M, 40M, 60M, and 120M examples are generated for number lengths 1-40, 1-30, 1-20, and 1-10, respectively. In contrast, the test set consists of 1,000 examples per digit length.

### B.2. Training Details

Our base model, following Zhou et al. (2023), is a 25M parameter Transformer featuring 6 blocks, a 512 hidden size, a feedforward layer with a hidden dimension of 2048 using GeGLU activation (Shazeer, 2020), and an 8-head attention mechanism. We also adopt RMSNorm, integrating both PreNorm and PostNorm layers, following the Primer architecture (So et al., 2021). Additionally, our preliminary investigations underscore the significance of employing causal language modeling when applying the index hint technique. Conversely, attempts to leverage prefix language modeling paired with bidirectional attention in model inputs consistently falter in length generalization. Our three other model variants with size [2M, 5M, 268M] consist of [2, 4, 16] blocks, a [256, 256, 1024] hidden size, a feedforward layer with a hidden dimension of [1024, 1024, 4096], and a [4, 4, 16]-head attention mechanism, respectively.

In our implementation of FIRE Li et al. (2023), we employ layerwise sharing of attention bias across all attention blocks to enhance training efficiency. The paraterization of FIRE consists of a 2-layer MLP with a 32-unit hidden layer, utilizing ReLU activation.

We use the AdamW optimizer (Loshchilov and Hutter, 2017) to train the model with a weight decay value of 0.1 and dropout rate of 0.0. The learning rate schedule incorporates an initial 500-step linear warm-up, followed by a cosine decay, starting at  $3e-4$ . We train the model with sequence packing, a batch size of 128, and a sequence length of 2048, over 50,000 steps. We use greedy decoding to generate the model output during evaluation. We summarize the hyperparameters in Table B.1.

Table B.1 | Hyperparameters Summary for Length Generalization

Hyperparameter	Value
Language Model Type	Causal
Activation Functions	GeGLU
Normalization Layer	RMSNorm
Normalization Type	PreNorm and PostNorm
Optimizer	AdamW
Training Steps	50,000
Batch size	128
Weight Decay	0.1
Dropout	0.0
Learning Rate (LR)	0.0003
LR Warmup Steps	500
LR Cooldown (Begin, End)	(500, 50,000)
Warmup Schedule	Linear (from 0 to LR)
Cooldown Schedule	Cosine Decay (from LR to 0.1LR)
Training Sequence Length	2048
Evaluation	Greedy

## C. Additional Results

### C.1. Training Loss and Sequence Exact Match Accuracy of Reverse Format with Index Hint trained up to 40

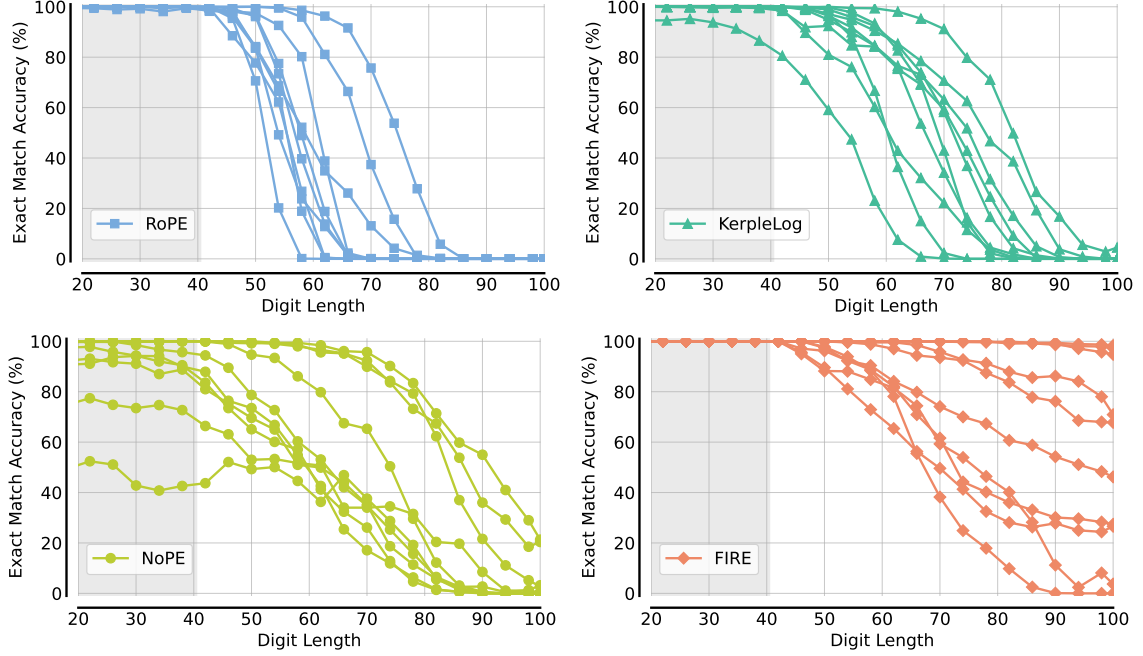


Figure C.1 | Exact match accuracy on 20 to 100 digit addition of all 10 trials trained on up to 40-digit addition with index hint and reverse format using four different position encodings.

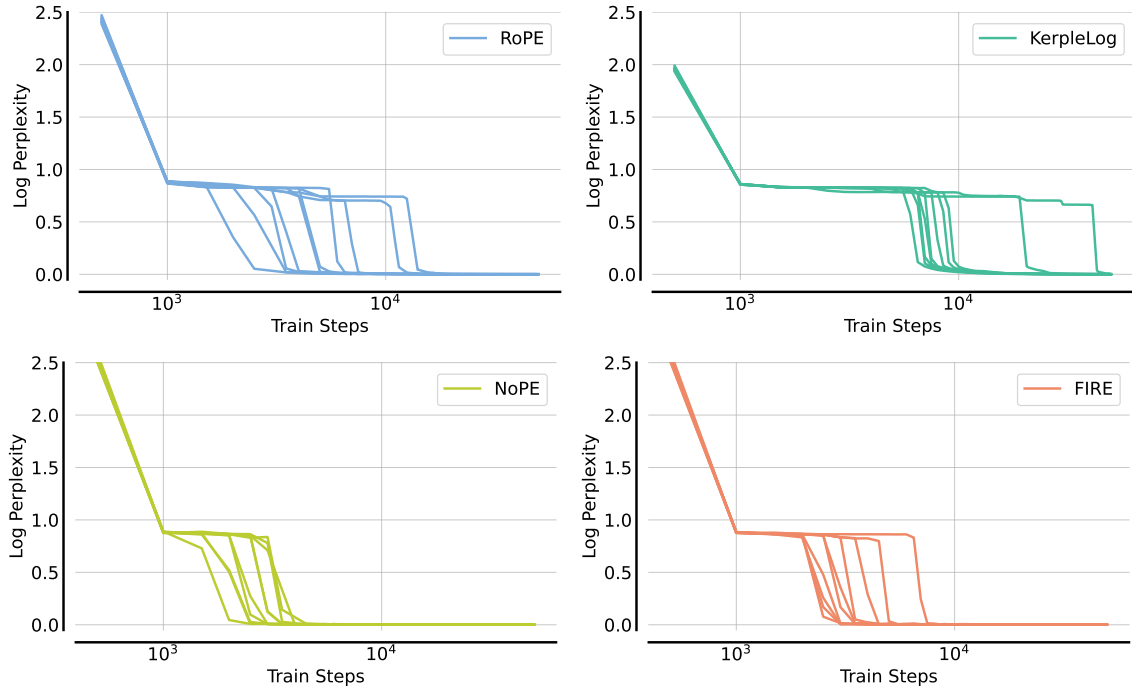


Figure C.2 | Training loss over 10 trials in reverse formats. Despite similar nearly 0 log perplexity losses across runs after 10K training steps, different runs exhibit very different length generalization.

## C.2. Training Loss and Next-token Prediction Accuracy of Standard and Reverse Format

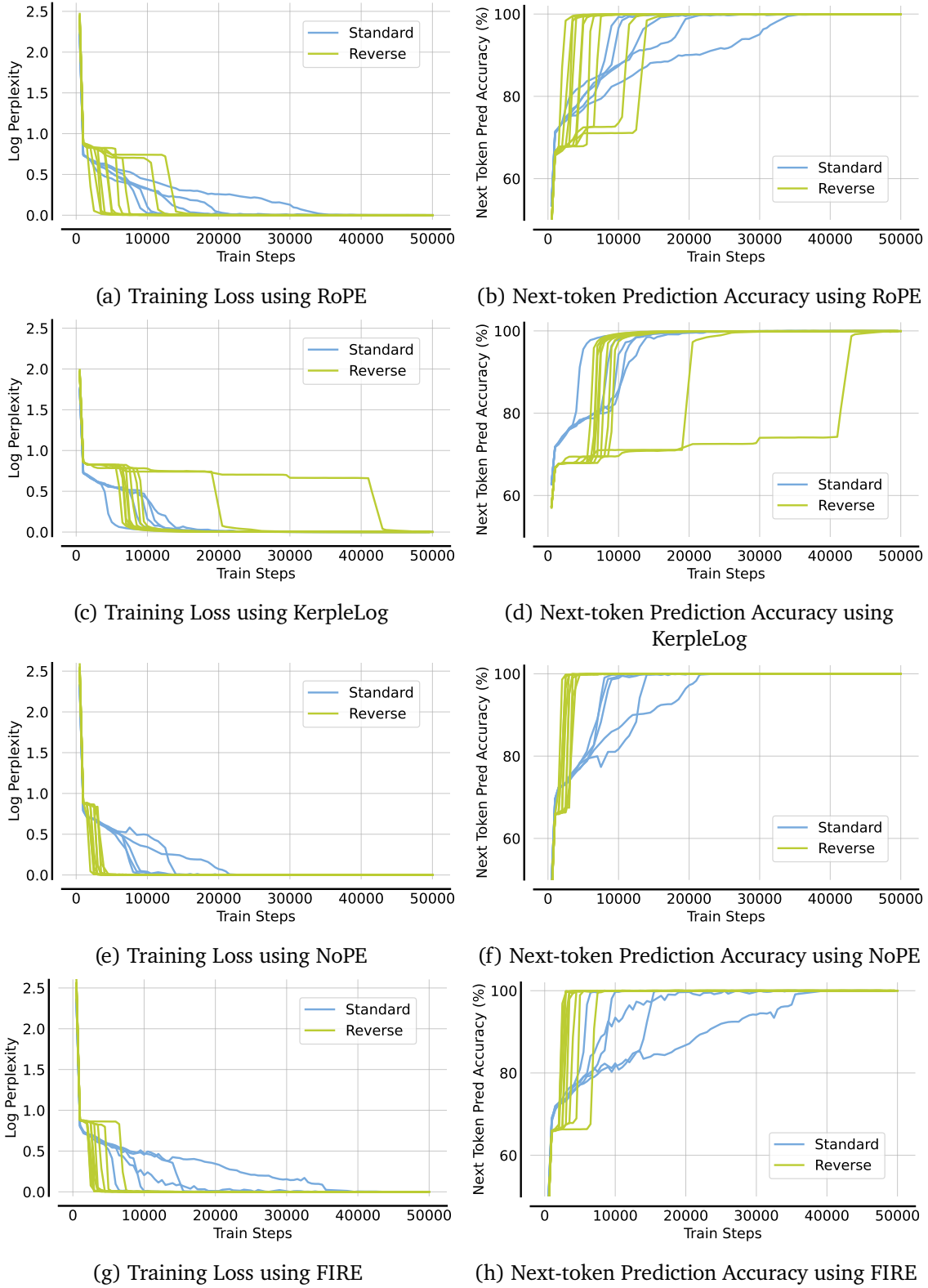


Figure C.3 | Training log perplexity and next-token prediction accuracy over 10 trials in standard versus reverse formats using RoPE, KerpleLog, NoPE and FIRE. Reverse format shows distinct grokking during training, unlike the gradual enhancement in standard format.

### C.3. The evolution of EM Accuracy during training in reverse format using 4 PEs

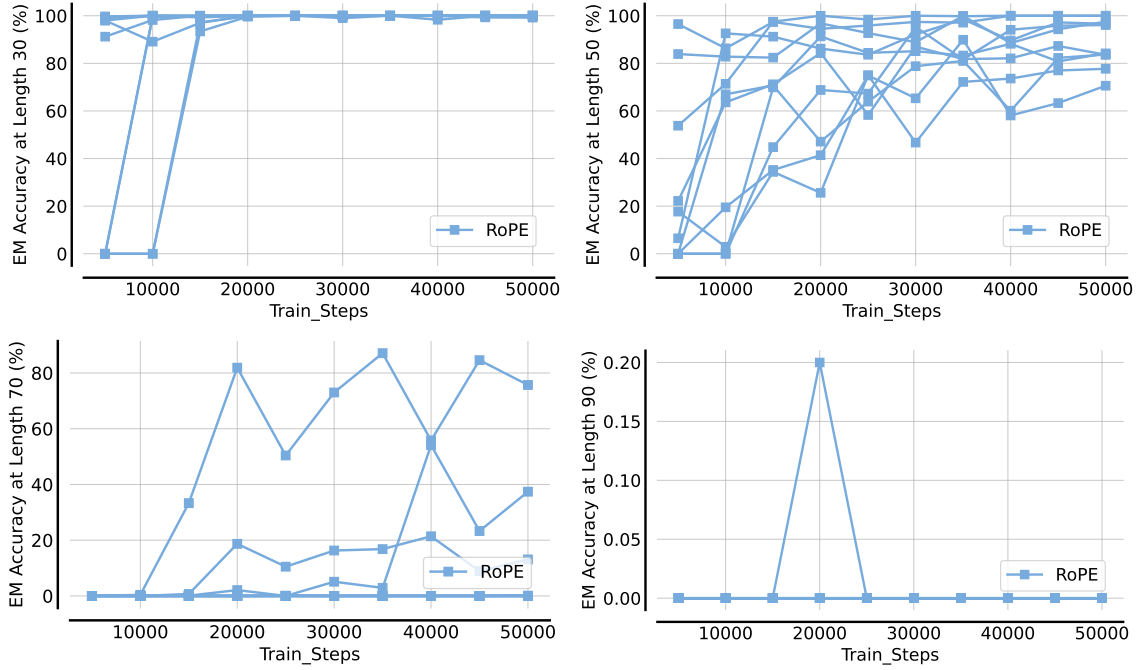


Figure C.4 | Exact match accuracy on [30, 50, 70, 90] digit addition of all 10 trials trained on up to 40-digit addition with index hint and reverse format using RoPE.

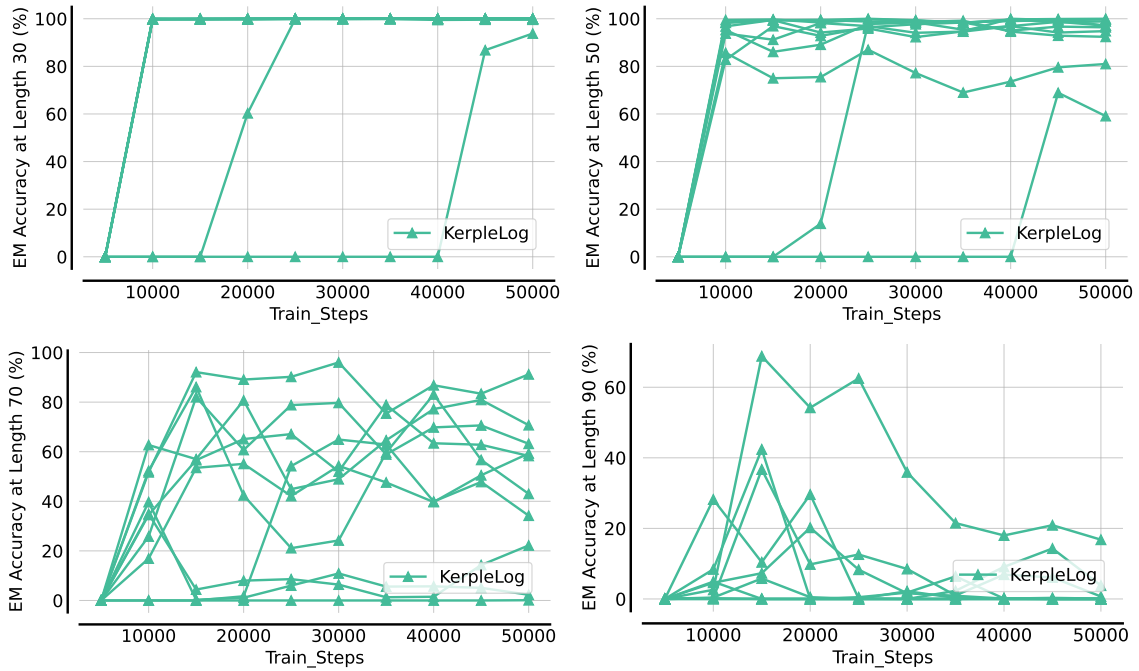


Figure C.5 | Exact match accuracy on [30, 50, 70, 90] digit addition of all 10 trials trained on up to 40-digit addition with index hint and reverse format using KerpleLog.



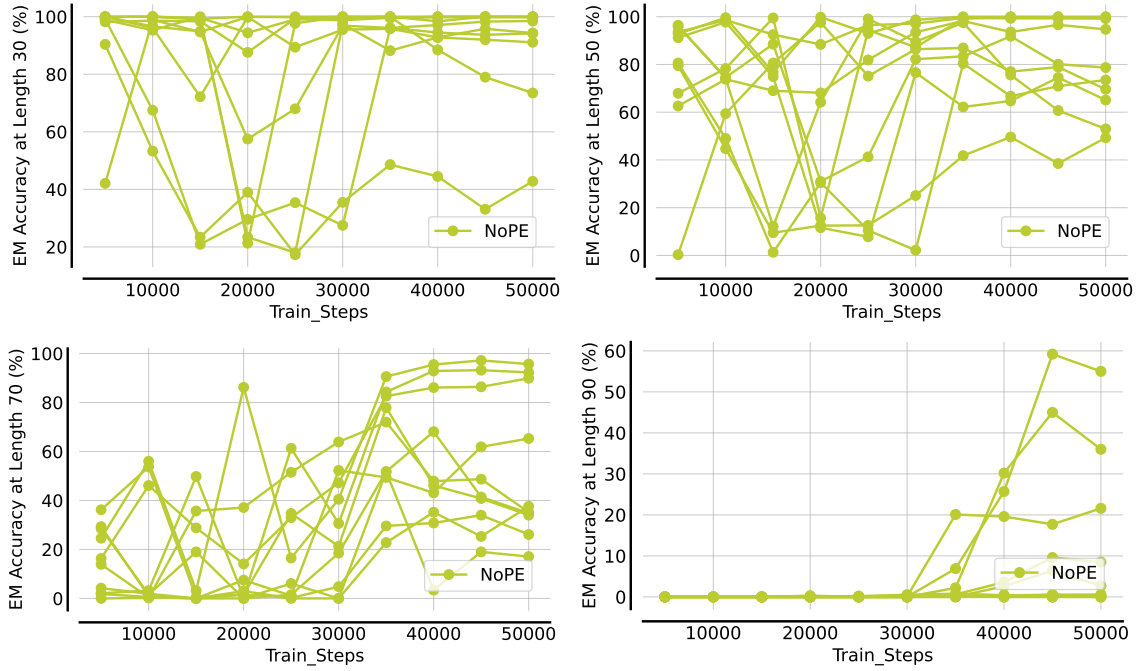


Figure C.6 | Exact match accuracy on [30, 50, 70, 90] digit addition of all 10 trials trained on up to 40-digit addition with index hint and reverse format using NoPE.

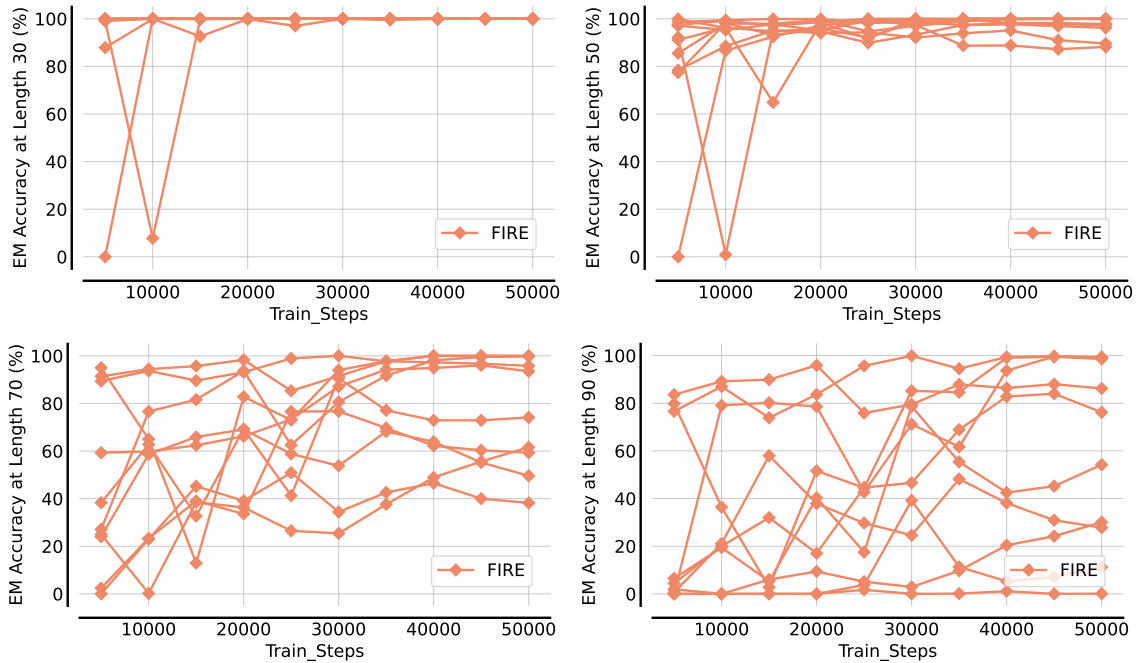
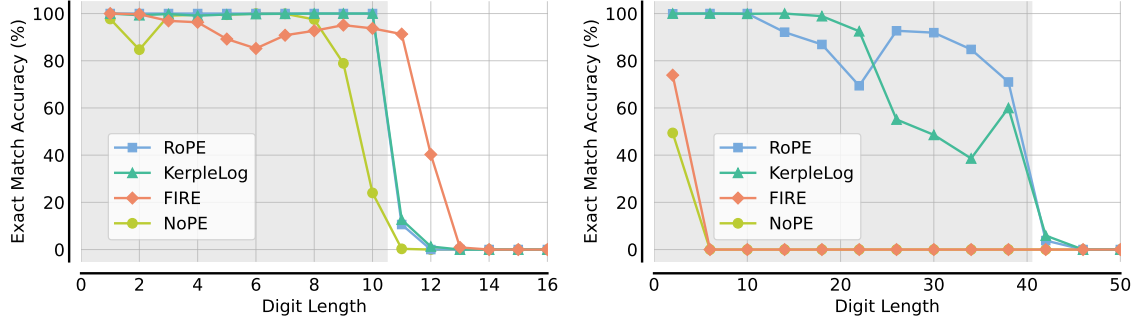


Figure C.7 | Exact match accuracy on [30, 50, 70, 90] digit addition of all 10 trials trained on up to 40-digit addition with index hint and reverse format using FIRE.

#### C.4. Effect of Index Hint



(a) Models trained up to 10-digit addition

(b) Models trained up to 40-digit addition

Figure C.8 | Best sequence exact match accuracy over five trials without index hint, trained up to length 10. All position encoding methods fail to generalize beyond trivial lengths and struggle with in-distribution generalization, highlighting the crucial role of index hints in length generalization. See the performance of each run in Appendix D.2 and Appendix D.1 for trained up to 10-digit and 40-digit addition.

#### C.5. Source of Variance

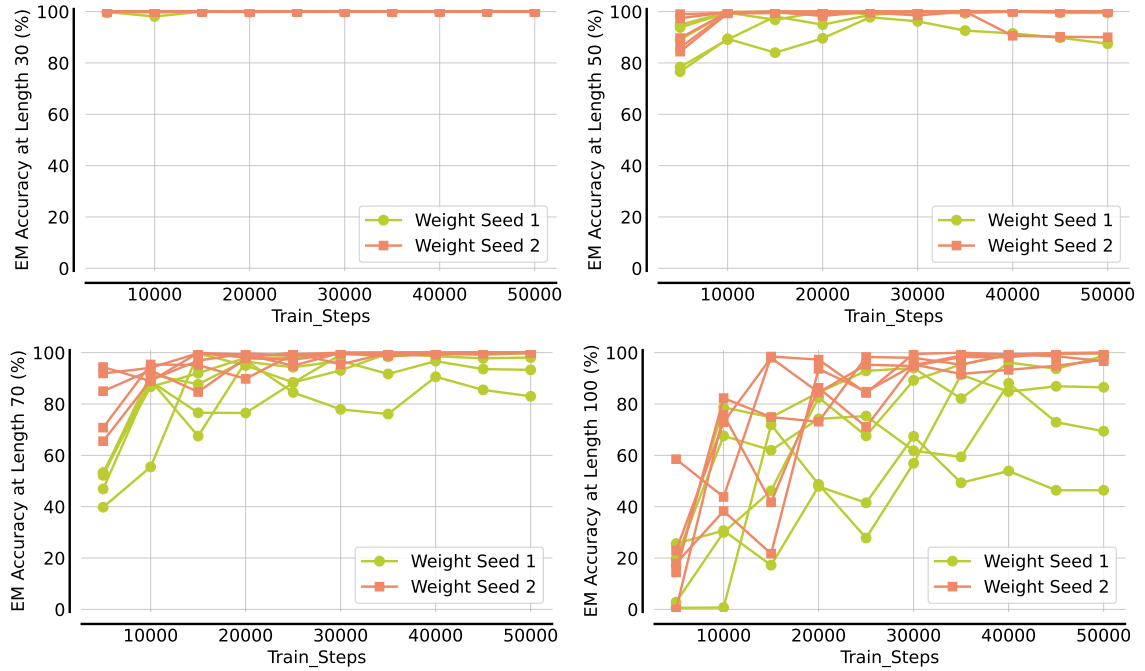


Figure C.9 | Exact match accuracy on [30, 50, 70, 100] digit addition of all 10 trials (2 weight seeds x 5 data seeds) trained on up to 40-digit addition with index hint and reverse format using FIRE.

## C.6. Error Analysis

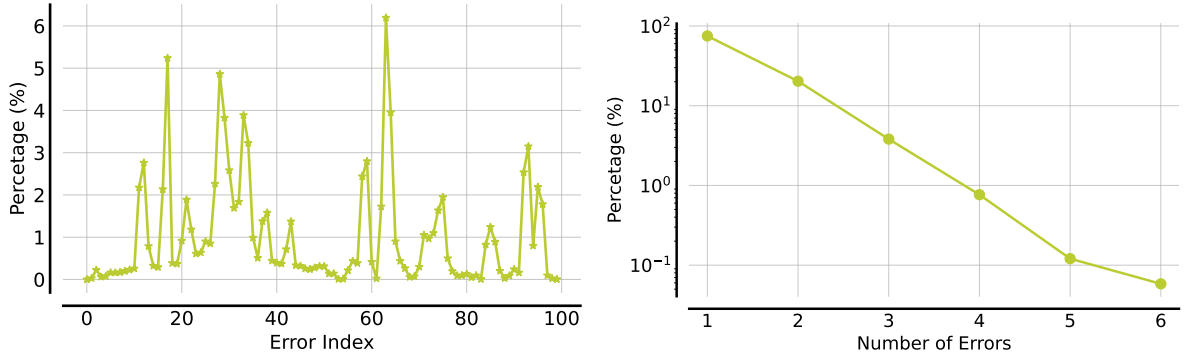


Figure C.10 | (Left) Average error position distribution over 10 runs, showing a broad error spread across all positions. Specific checkpoints exhibit a propensity for errors at certain positions (refer to Figure C.12). (Right) Notably, in successful generalizations, more than 90% of errors are confined to single-digit inaccuracies, exhibiting an exponential distribution.

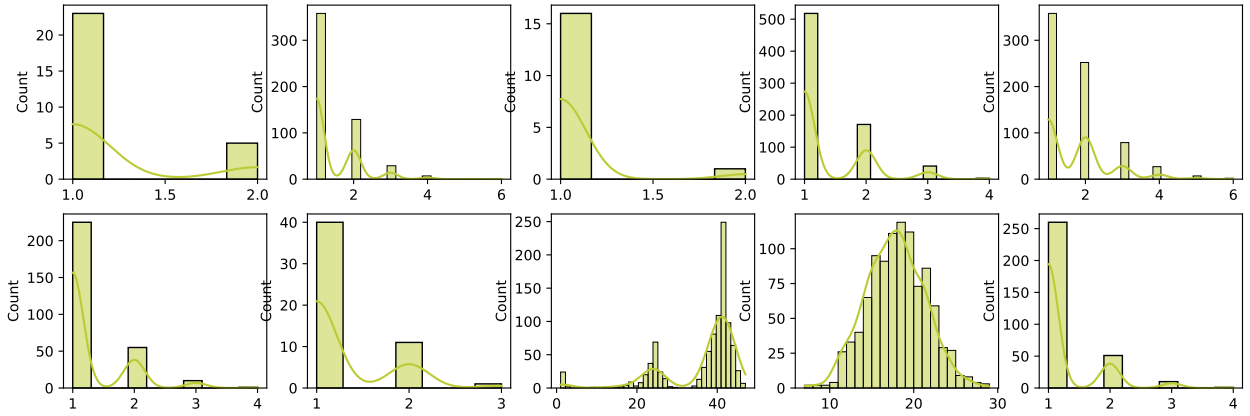


Figure C.11 | Error count distribution

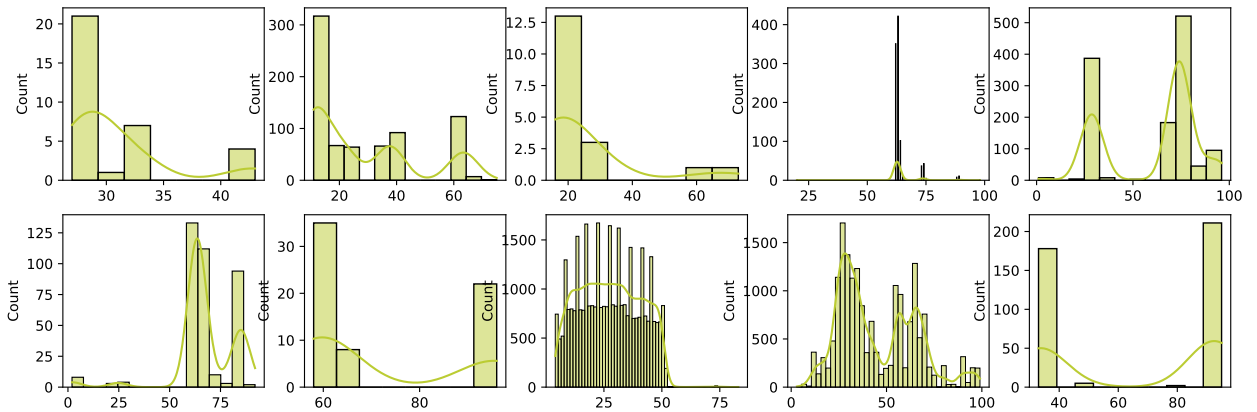


Figure C.12 | Error position distribution (FIRE)

## C.7. Training Digit Length Effect

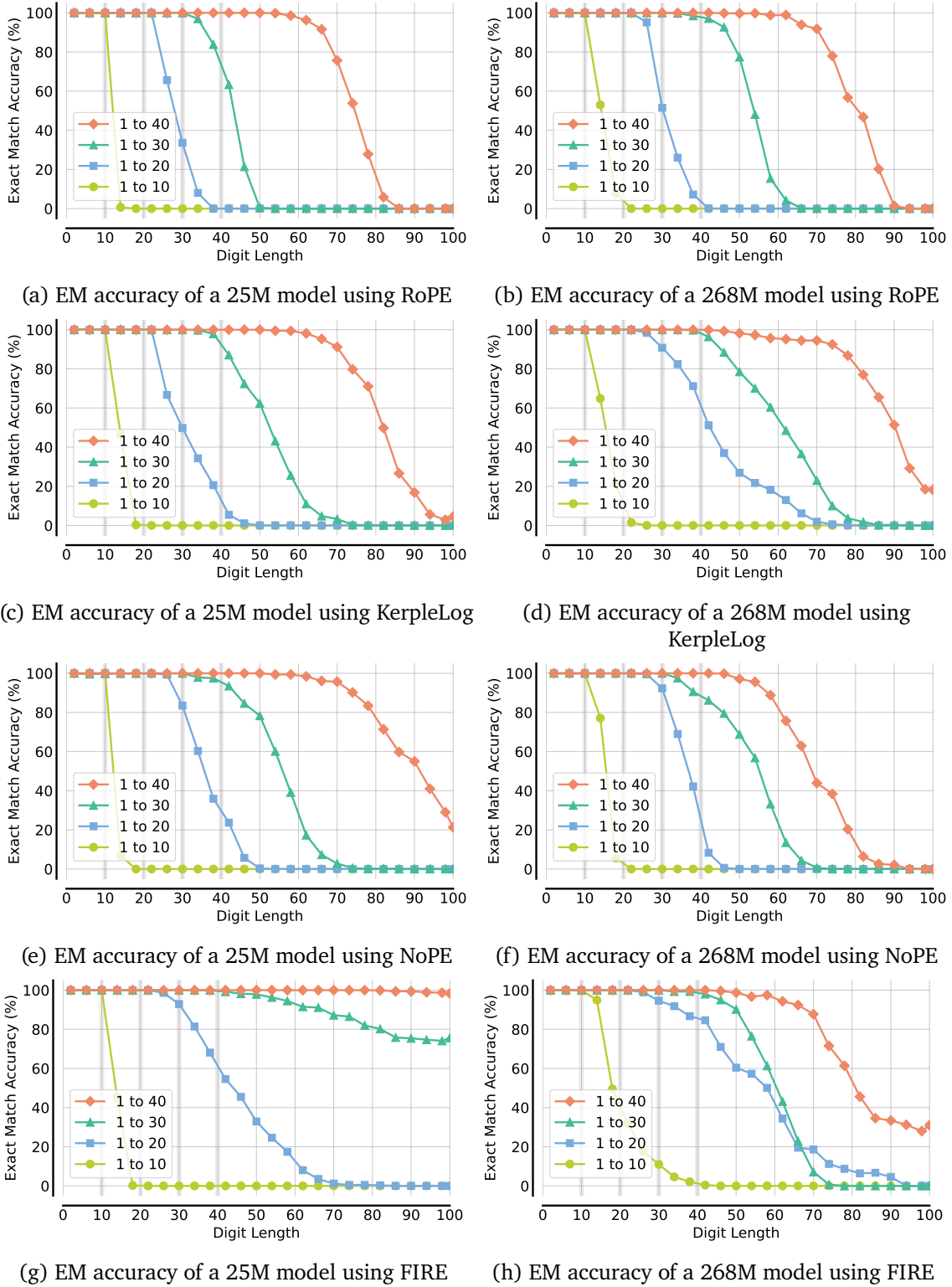


Figure C.13 | Best sequence exact match accuracy of 5 trials with two model sizes (i.e., 25M and 268M), trained on up to 10, 20, 30 and 40 digit length using 4 PEs.

## C.8. Model Size Effect

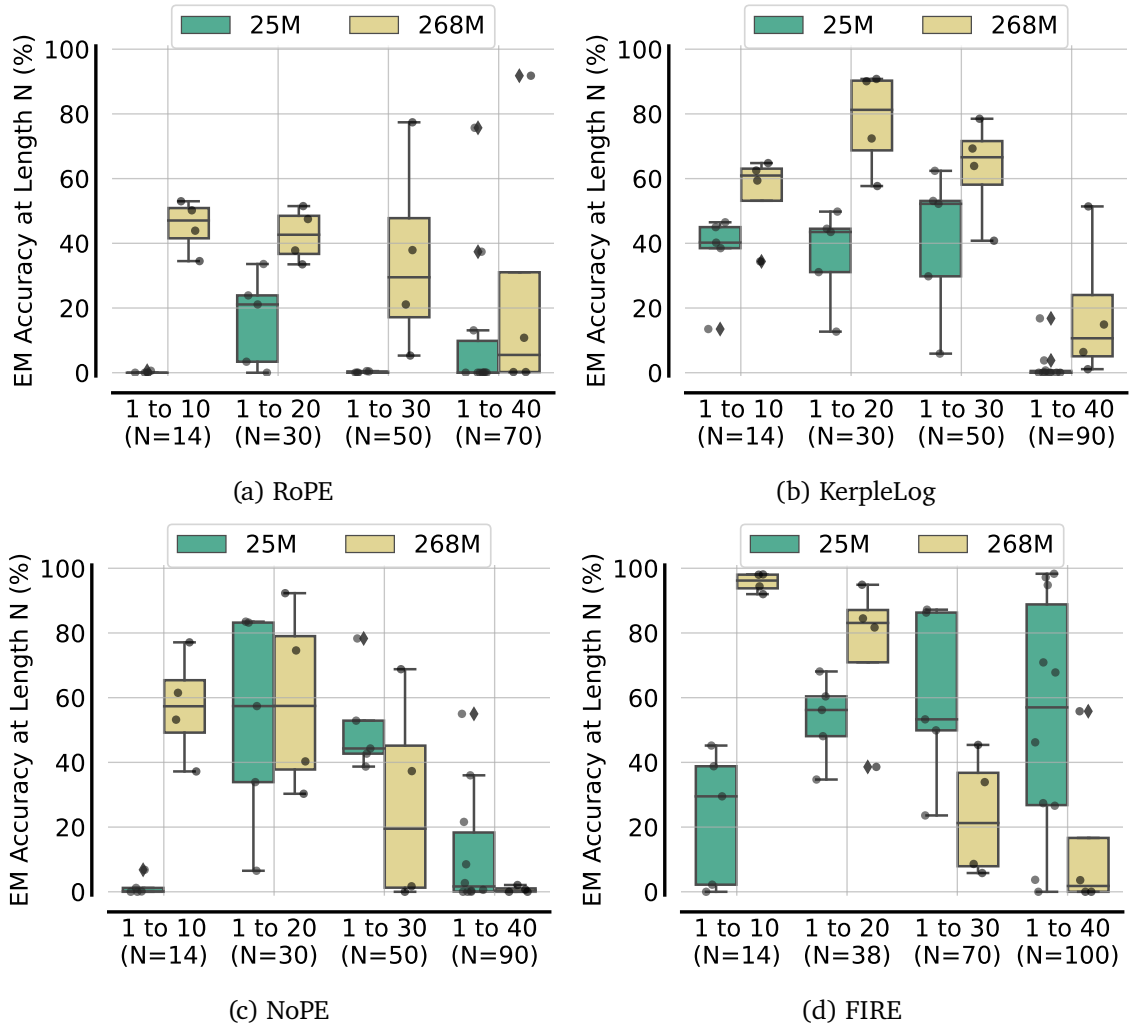


Figure C.14 | Scaling model size inconsistently affects length generalization performance. While consistently enhancing performance in shorter length regimes (1-10, 1-20) across four position encodings, this trend does not hold for larger regimes (1-30, 1-40). For instance, larger models outperform smaller ones with RoPE and KerpleLog encodings, but underperform with NoPE and FIRE. Moreover, increasing model size doesn't noticeably decrease performance variance, suggesting size scaling isn't vital for length generalization.

## C.9. FIRE Related Scaling

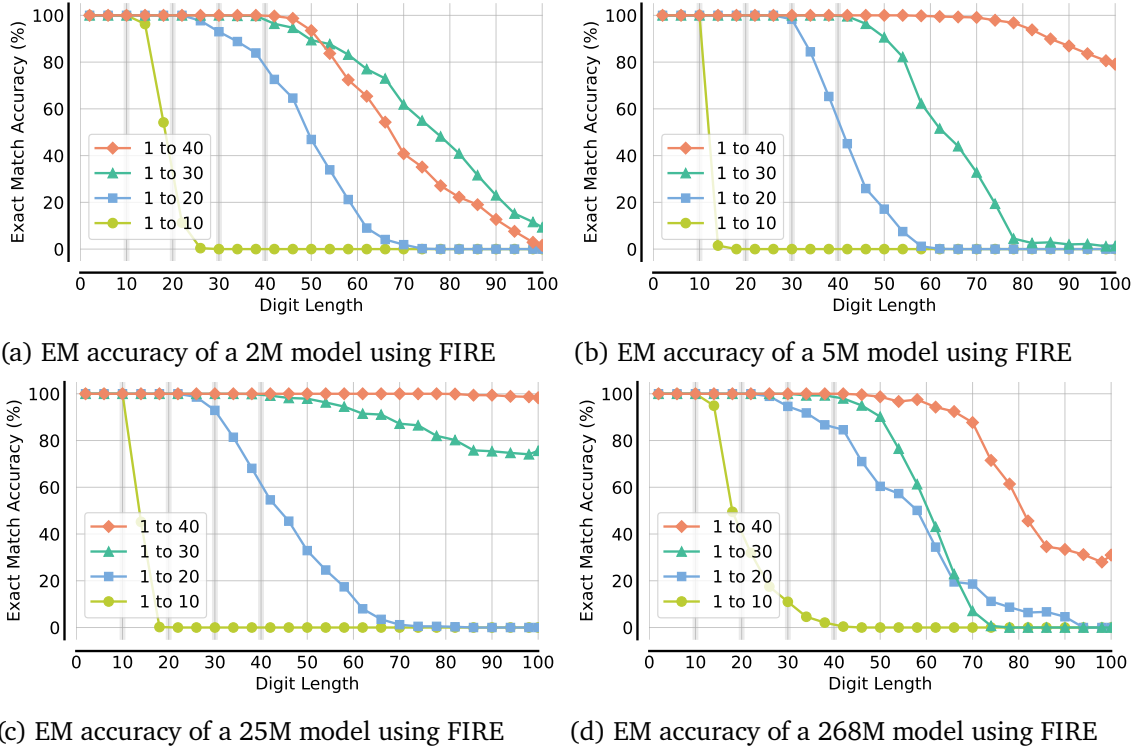


Figure C.15 | Best sequence exact match accuracy of 5 trials with four model sizes (i.e., 2M, 5M, 25M and 268M), trained on up to 10, 20, 30 and 40 digit length using **FIRE**.



## C.10. Hyperparameter Study

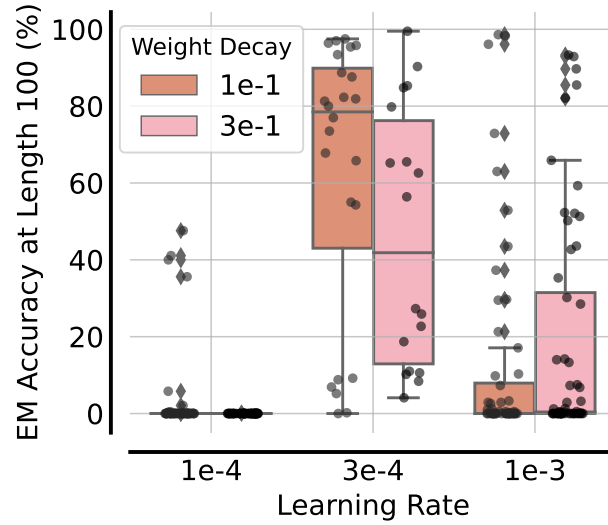


Figure C.16 | Sequence exact match accuracy for test digit length 100, trained on digit lengths 1-40.  $3e-4$  seems to be the optimal learning rate.

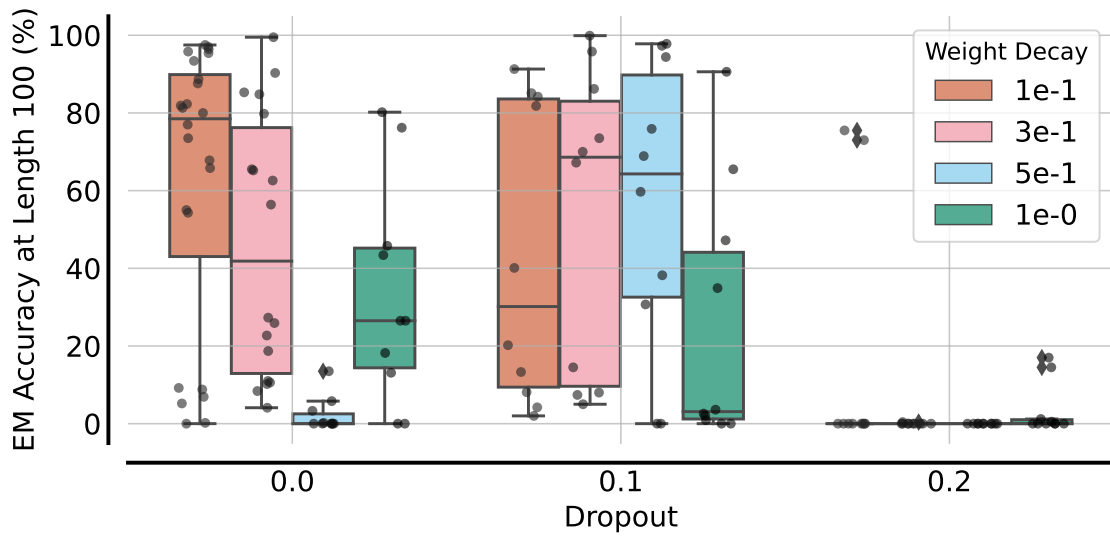


Figure C.17 | Sequence exact match accuracy for test digit length 100, trained on digit lengths 1-40. A higher dropout rate markedly impedes length generalization, whereas a lower rate shows negligible impact.

## D. Training Loss and Sequence Exact Match Accuracy

### D.1. Reverse Format without Index Hint trained up to 40-digit addition

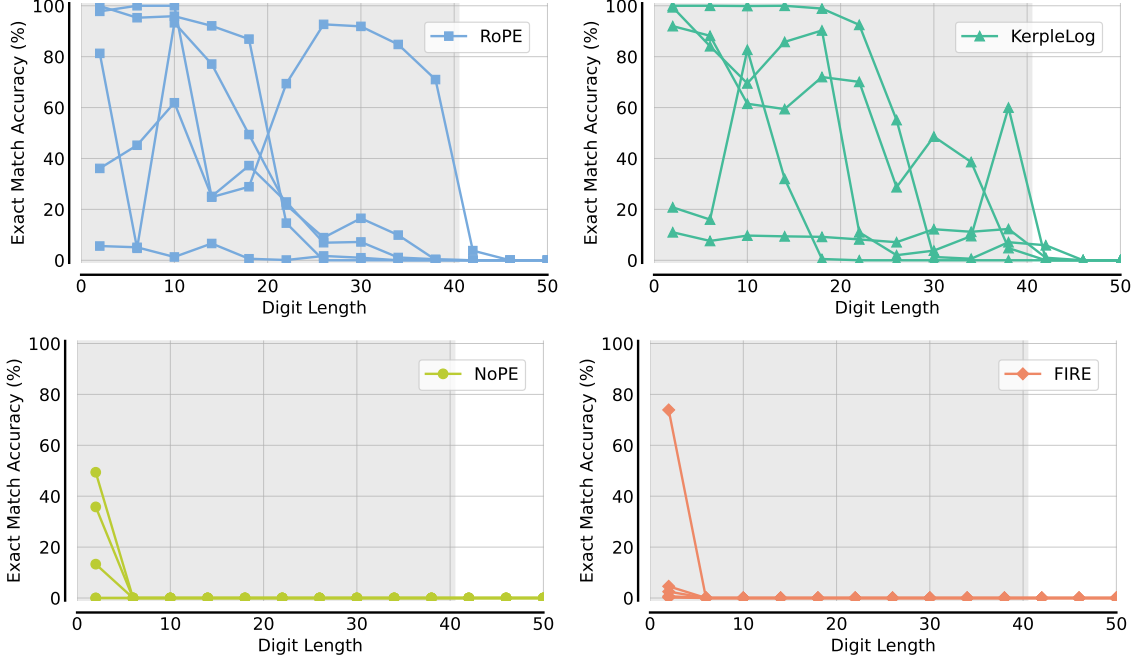


Figure D.1 | Exact match accuracy on 20 to 100 digit addition of all 10 trials trained on up to 40-digit addition with index hint and reverse format using four different position encodings.

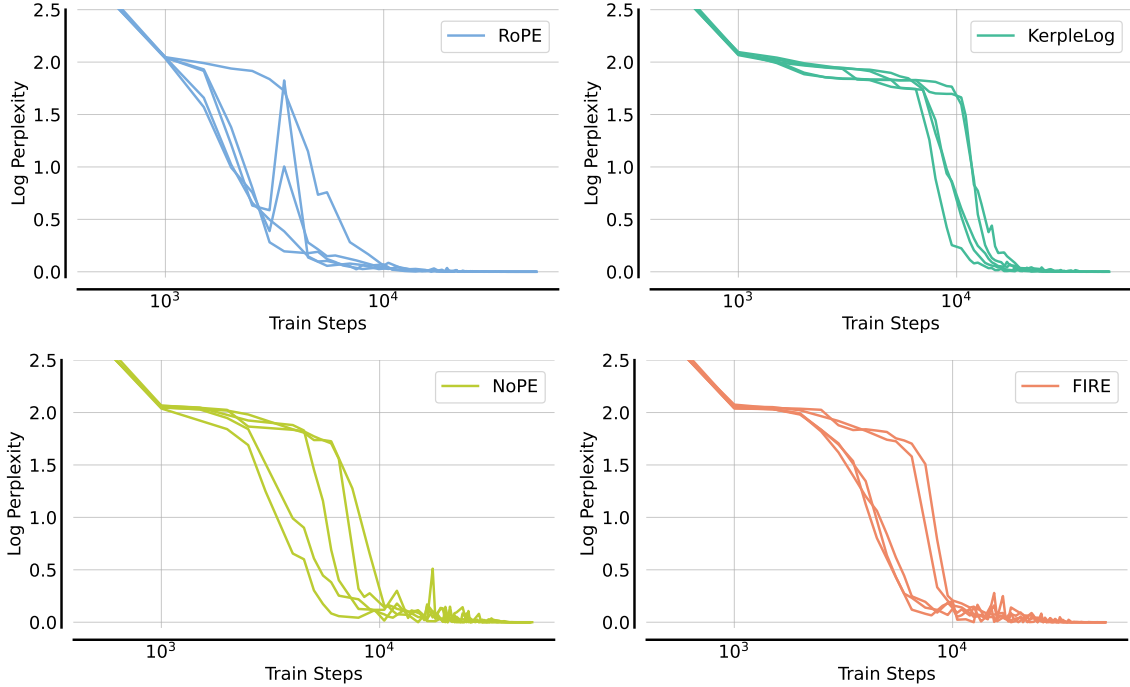


Figure D.2 | Training loss over 10 trials in reverse formats. Despite similar nearly 0 log perplexity losses across runs after 10K training steps, different runs exhibit very different length generalization.

## D.2. Reverse Format without Index Hint trained up to 10-digit addition

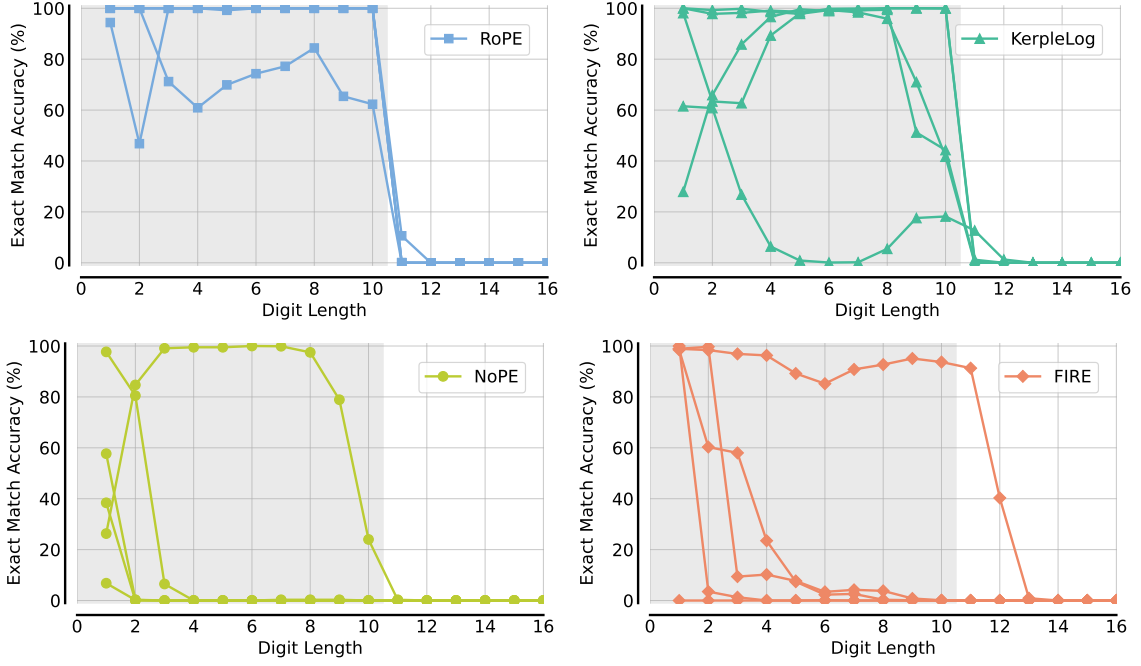


Figure D.3 | Exact match accuracy on 20 to 100 digit addition of all 10 trials trained on up to 10-digit addition with index hint and reverse format using four different position encodings.

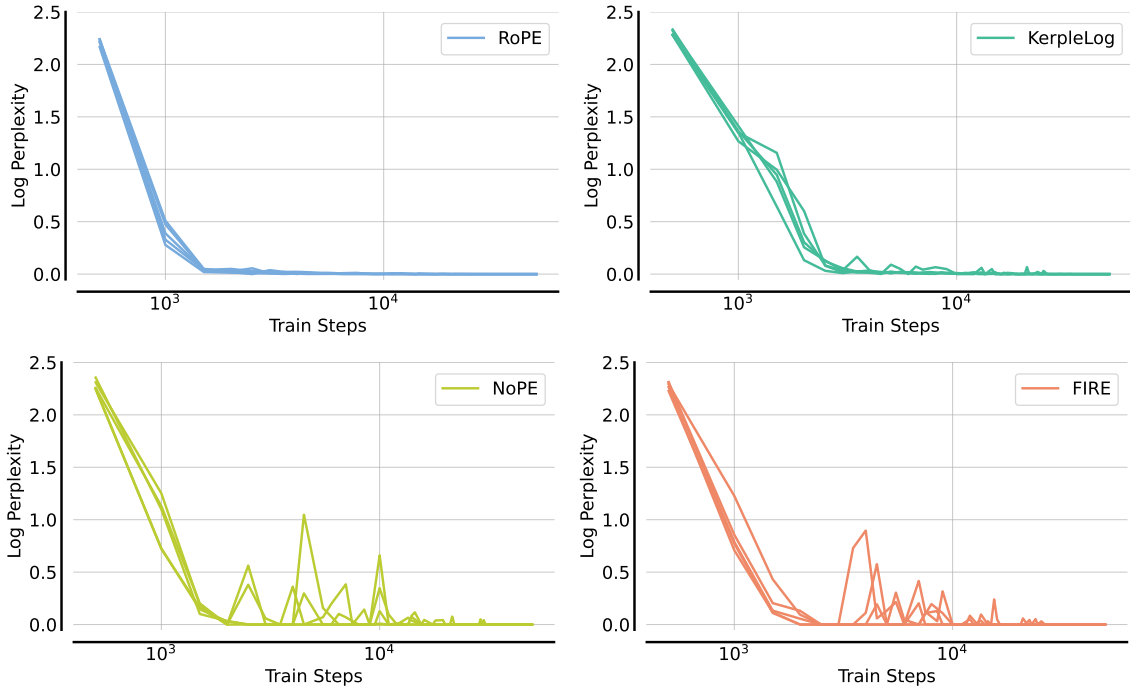


Figure D.4 | Training loss over 10 trials in reverse formats. Despite similar nearly 0 log perplexity losses across runs after 10K training steps, different runs exhibit very different length generalization.

### D.3. Standard Format with Index Hint trained up to 40

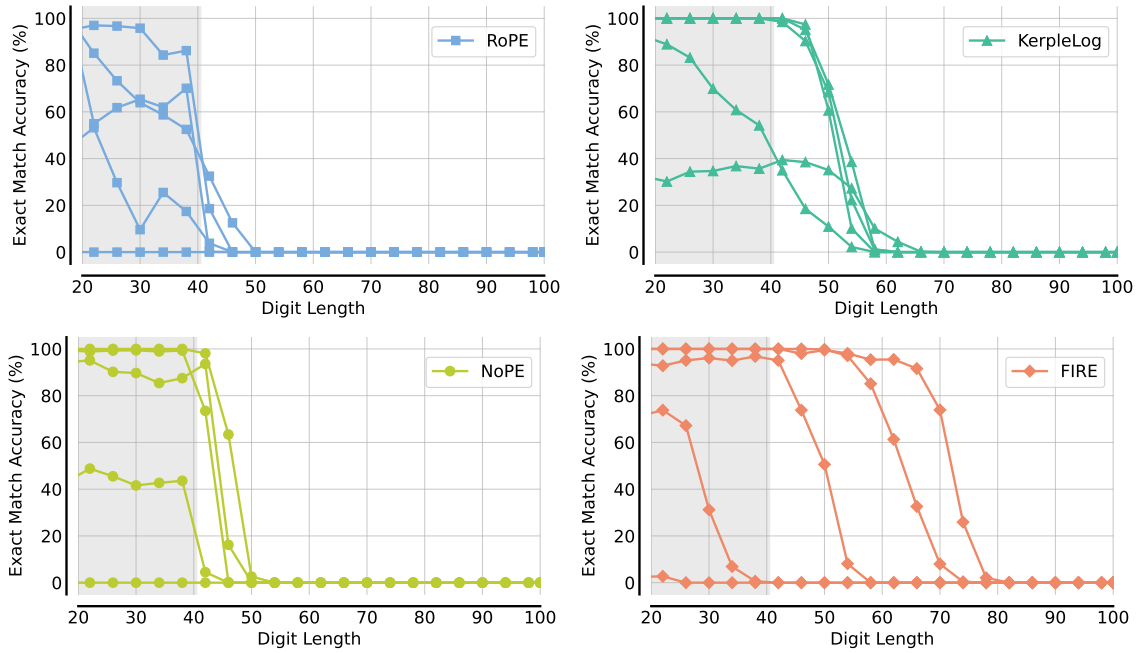


Figure D.5 | Exact match accuracy on 20 to 100 digit addition of all 10 trials trained on up to 40-digit addition with index hint and reverse format using four different position encodings.

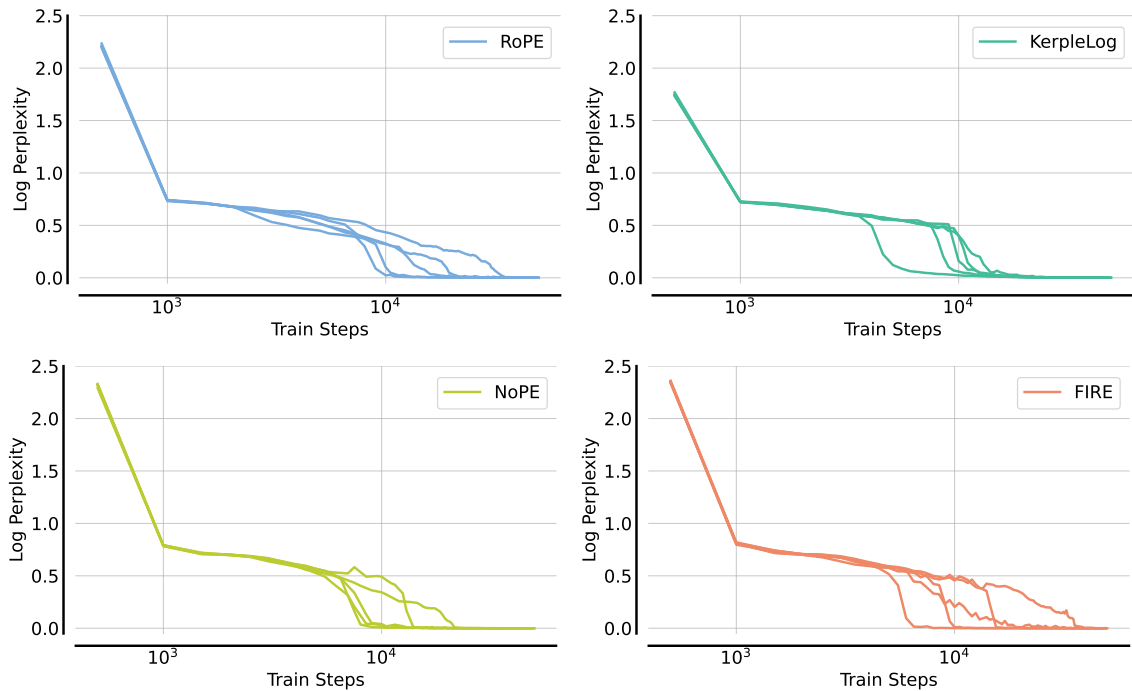


Figure D.6 | Training loss over 10 trials in reverse formats. Despite similar nearly 0 log perplexity losses across runs after 10K training steps, different runs exhibit very different length generalization.

#### D.4. Random Space Augmentation with Reverse Format with Index Hint trained up to 40-digit addition

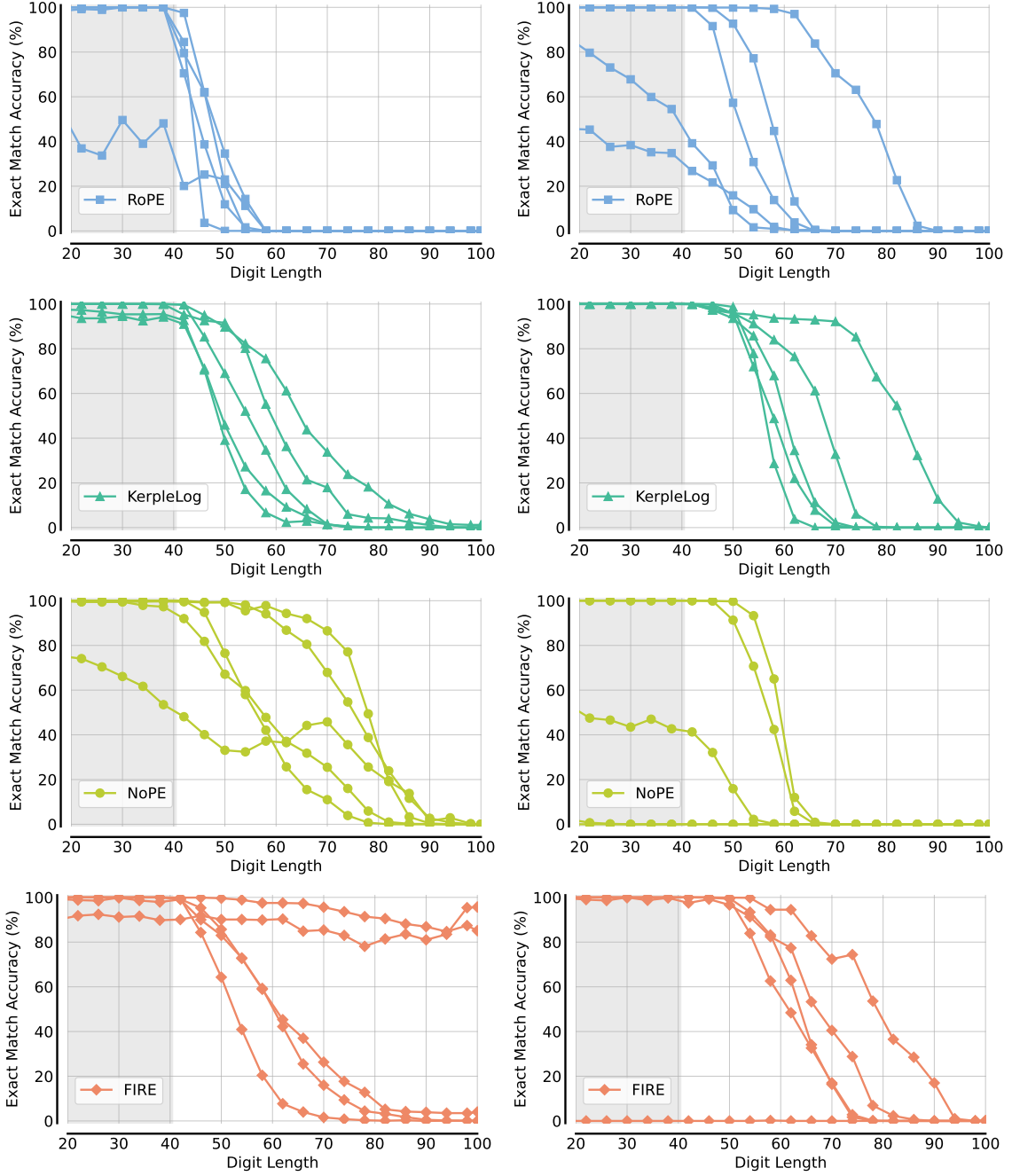


Figure D.7 | (Left) With Random Space Augmentation. (Right) Without Random Space Augmentation. Exact match accuracy on 20 to 100 digit addition of all 10 trials trained on up to 40-digit addition with index hint and reverse format using four different position encodings.

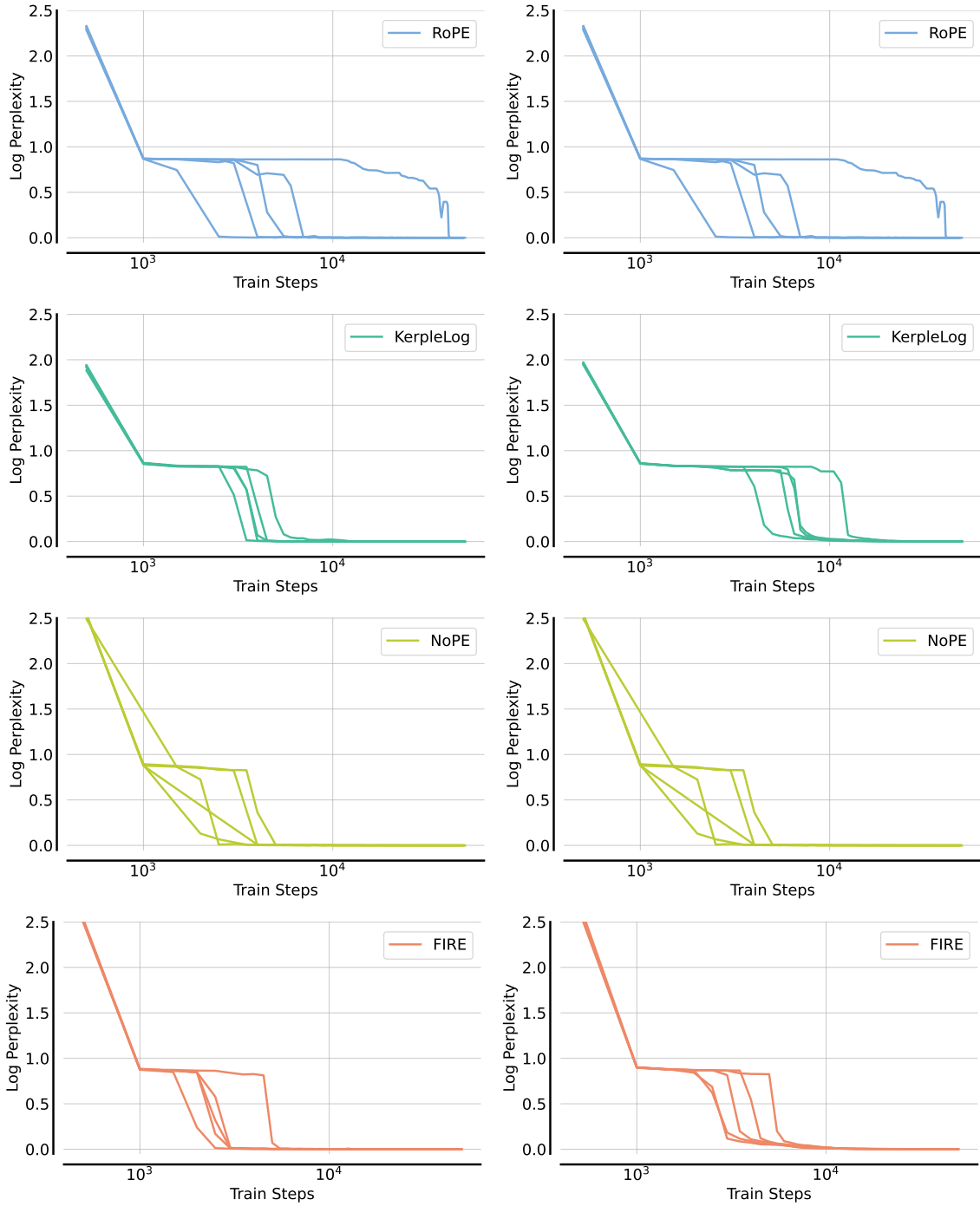


Figure D.8 | (Left) Without Random Space Augmentation. (Right) With Random Space Augmentation. Training loss over 10 trials in reverse formats. Despite similar nearly 0 log perplexity losses across runs after 10K training steps, different runs exhibit very different length generalization.



### D.5. Randomized Position Encoding with Reverse Format with Index Hint trained up to 40-digit addition

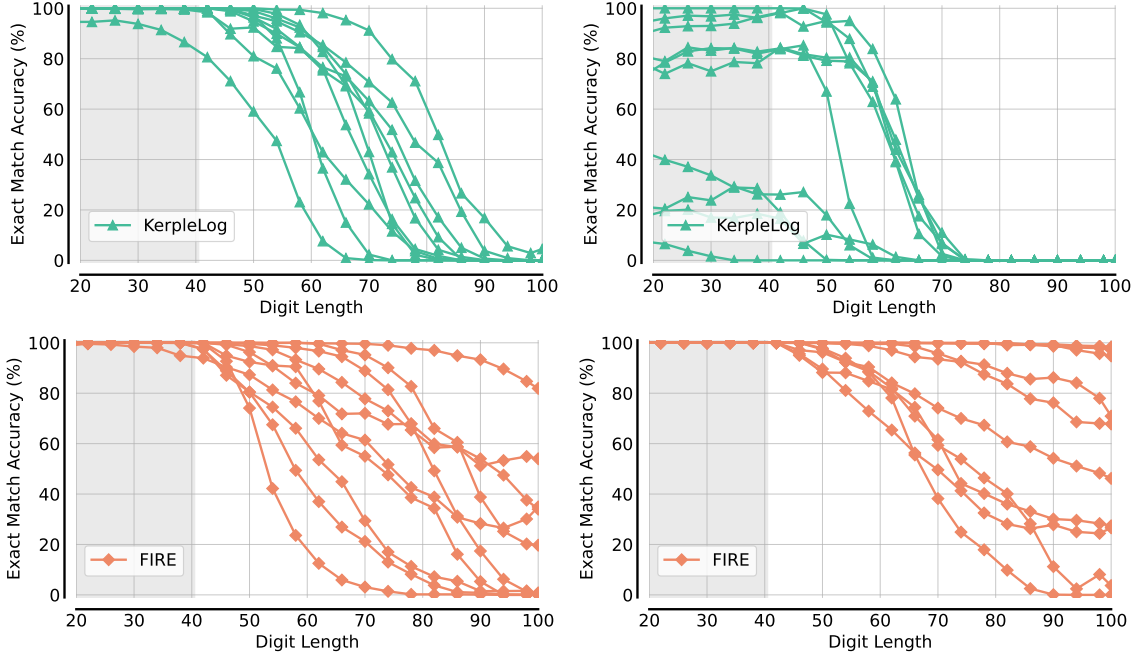


Figure D.9 | (Left) Without Randomized Position Encoding (Right) With Randomized Position Encoding. Exact match accuracy on 20 to 100 digit addition of all 10 trials trained on up to 40-digit addition with index hint and reverse format using four different position encodings.

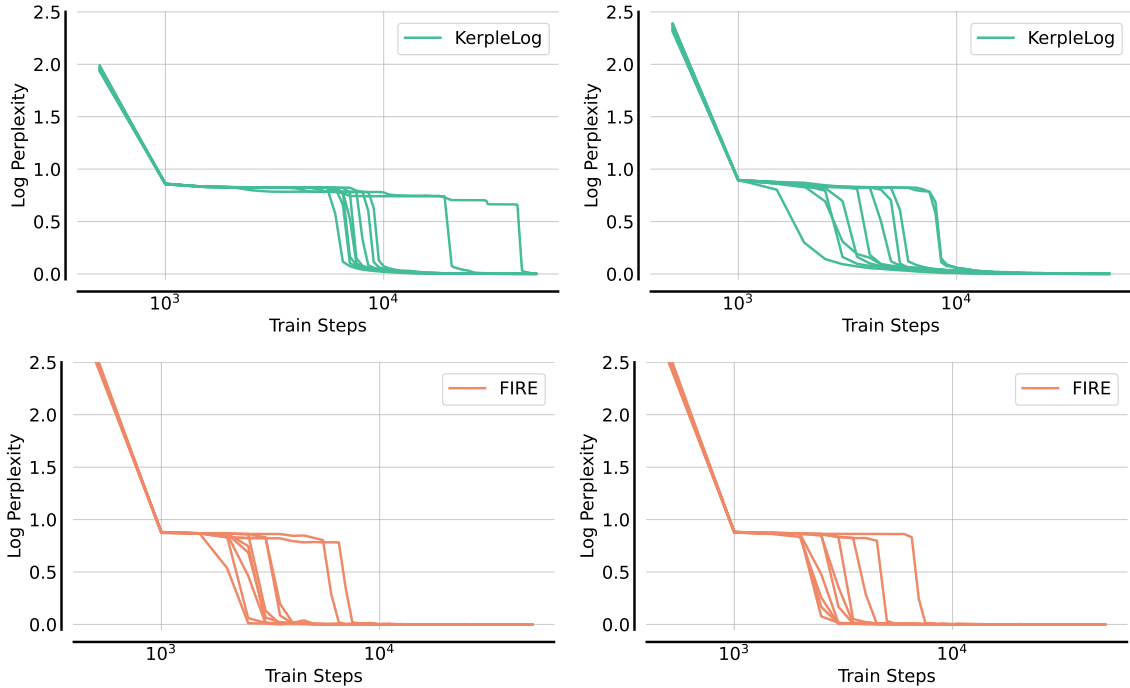


Figure D.10 | (Left) Without Randomized Position Encoding (Right) With Randomized Position Encoding. Training loss over 10 trials in reverse formats.