

Figure A.1: Histogram of instruction lengths in two instruction finetuning datasets: FLAN (CoT subset) (Longpre et al., 2023) and Super Natural Instructions (Wang et al., 2022). The dotted line indicates the median length of the instructions in each dataset.

532 A Number of instances decreases rapidly as sequence length grows

533 The recent trend of SFT-RLHF pipeline (Ouyang et al., 2022) relies on finetuning LLMs on the
 534 instruction following tasks. However, the training data of these datasets is often skewed towards
 535 shorter sequences. Figure A.1 shows the distribution of instruction lengths in two instruction
 536 finetuning datasets: FLAN (CoT subset) (Longpre et al., 2023) and Super Natural Instructions
 537 (Wang et al., 2022). The median length of instructions in these datasets is quite short compared
 538 to the maximum length that exists. Such distribution shape highlights the importance of length
 539 generalization in these tasks. In fact, the models are supposed to learn from short instructions and
 540 generalize to ones during inference that might be much longer.

541 B Background

542 B.1 Preliminaries

543 In this section, we lay the groundwork and introduce the notation we use throughout the paper. We
 544 will refer to this in Appendices C.1 and C.2.

545 Let f_θ be a decoder-only Transformer model, where θ denotes the full set of model parameters.
 546 f_θ processes the input sequence $\mathbf{x} = [x_0, x_1, \dots, x_T]$ and maps it to the output sequence $\mathbf{y} =$
 547 $[y_0, y_1, \dots, y_T]$ by applying a sequence of Transformer layers. Note that being decoder-only means
 548 the attention mechanism in each layer is causal, i.e. the attention weights are computed based on the
 549 previous positions only.

550 The layer $\text{TLayer}^{(l)}(\mathbf{H}^{(l-1)}; \theta_l)$, consisting of self-attention heads and a feed-forward sub-layer,
 551 reads the previous hidden state $\mathbf{H}^{(l-1)}$ and produces the hidden state at layer l : \mathbf{H}^l , where l is the
 552 layer index, and θ_l is the set of parameters of the l -th layer. Each hidden state $\mathbf{H}^{(l)} \in \mathbb{R}^{d \times (T+1)}$ is
 553 matrix where column t , denoted as $\mathbf{h}_t^{(l)}$, is the hidden state at position t .

554 A layer l is parameterized by a set of parameters $\theta_l = \{(\mathbf{W}_Q^m, \mathbf{W}_K^m, \mathbf{W}_V^m, \mathbf{W}_O^m)_m, \mathbf{W}_1, \mathbf{W}_2\}$, where
 555 $\mathbf{W}_Q^m, \mathbf{W}_K^m, \mathbf{W}_V^m \in \mathbb{R}^{h \times d}$ and $\mathbf{W}_O^m \in \mathbb{R}^{d \times h}$ are the query, key, value, and output matrices of the
 556 m -th head, respectively. $\mathbf{W}_1, \mathbf{W}_2 \in \mathbb{R}^{d \times k \cdot d}$ are the weight matrices of the feed-forward sub-layer.
 557 d denotes the model’s hidden state size, h is the attention dimension (where $h = \frac{d}{\# \text{heads}}$), and k is
 558 a multiplier of the hidden state size in the feed-forward sub-layer (it is usually set to 4 in common
 559 implementations of the Transformer). Note that we drop the layer index l and the attention head
 560 index m where it is clear from the context.

561 The Transformer layer $\text{TLayer}^{(l)}$ processes each column of $\mathbf{H}^{(l-1)}$ independently and in parallel to
 562 produce the output. The computation of the t -th column of $\mathbf{H}^{(l)}$ is as follows:

$$\mathbf{h}_t^{(l)} = \text{FF}(\lambda(\mathbf{a}_t + \mathbf{h}_t^{(l-1)})) + \mathbf{a}_t + \mathbf{h}_t^{(l-1)} \quad (3)$$

563 where FF is the feed-forward sub-layer, λ is layer normalization, and $\mathbf{a}_t \in \mathbb{R}^d$ is the output of the
 564 multi-head self-attention sub-layer at position t . Specifically, \mathbf{a}_t is computed as:

$$\mathbf{a}_t = \sum_m \text{Attn}^{(m)}(\mathbf{h}_t^{(l-1)}, \mathbf{H}^{(l-1)}) \quad (4)$$

565 where $\text{Attn}^{(m)}$ is the m -th attention head. Let $\mathbf{o}_t \in \mathbb{R}^d$ denote the output of an attention head at
 566 position t . Then, \mathbf{o}_t is computed as:

$$\mathbf{o}_t = \mathbf{W}_O \left(\sum_{i \leq t} \hat{\alpha}_i \mathbf{v}_i \right) \quad (5)$$

567 where $\hat{\alpha} = \text{softmax}(\boldsymbol{\alpha}) \in \mathbb{R}^{(t+1)}$, and $\boldsymbol{\alpha}$ is the attention weight vector such that:

$$\boldsymbol{\alpha} = [\langle \mathbf{q}_t, \mathbf{k}_0 \rangle, \langle \mathbf{q}_t, \mathbf{k}_1 \rangle, \dots, \langle \mathbf{q}_t, \mathbf{k}_t \rangle]^\top \quad (6)$$

568 where $\mathbf{q}_t = \mathbf{W}_Q \mathbf{h}_t^{(l-1)} \in \mathbb{R}^h$, $\mathbf{k}_i = \mathbf{W}_K \mathbf{h}_i^{(l-1)} \in \mathbb{R}^h$, and $\mathbf{v}_i = \mathbf{W}_V \mathbf{h}_i^{(l-1)} \in \mathbb{R}^h$. $\langle \cdot, \cdot \rangle$ denotes the
 569 dot product operation.

570 The feed-forward sub-layer $\text{FF}(\cdot) \in \mathbb{R}^d$ is a two-layer MLP:

$$\text{FF}(\mathbf{x}) = \mathbf{W}_2 \sigma(\mathbf{W}_1^\top \mathbf{x}) \quad (7)$$

571 where σ is a non-linear activation function (usually ReLU or GeLU (Hendrycks and Gimpel, 2020)).
 572 Additionally, $\lambda(\cdot) \in \mathbb{R}^d$ is layer normalization (Ba et al., 2016). Note that we take the *additive*
 573 (Elhage et al., 2021) view of attention heads in Equation (4) instead of *concatenate and multiple*
 574 view (Vaswani et al., 2017) as it is easier to understand and analyze. But, they are mathematically
 575 equivalent (Elhage et al., 2021).

576 The hidden state is initialized with a learned embedding of the input sequence $\mathbf{H}^{(0)} = \mathbf{W}_E \mathbf{X}$, where
 577 $\mathbf{W}_E \in \mathbb{R}^{d \times V}$ is the embedding matrix and $\mathbf{X} \in \mathbb{R}^{V \times (T+1)}$ is the one-hot encoded input sequence.
 578 V is the vocabulary size.

579 B.2 Positional Encoding

580 Almost all positional encoding methods can be explained and formulated as how they implement the
 581 dot product operation in Equation (6). So, in this section, we explain how the dot product $\langle \mathbf{q}_t, \mathbf{k}_i \rangle$ is
 582 implemented in different positional encoding schemes.

583 **Absolute Positional Encoding (APE)** The process of Absolute Positional Encoding (APE) involves
 584 assigning a position vector \mathbf{p}_i to each absolute position i and combining them with word embeddings
 585 before inputting them into the model. So, APE first modifies how the hidden state is initialized:

$$\mathbf{H}^{(0)} = \mathbf{W}_E \mathbf{X} + \mathbf{W}_P \mathbf{P} \quad (8)$$

586 where $\mathbf{W}_P \in \mathbb{R}^{d \times T}$ is the positional embedding matrix and $\mathbf{P} \in \mathbb{R}^{V_p \times (T+1)}$ is the one-hot encoded
 587 absolute position sequence. V_p is the maximum absolute position. Therefore, the hidden state at
 588 column j is:

$$\mathbf{h}_j^{(0)} = \mathbf{e}_j + \mathbf{p}_j \quad (9)$$

589 where $\mathbf{e}_j \in \mathbb{R}^d$ is the word embedding of token x_j and $\mathbf{p}_j \in \mathbb{R}^d$ is the positional embedding for
 590 position j . Then, the dot product for the first layer in Equation (6) is computed as:

$$\begin{aligned} \langle \mathbf{q}_t, \mathbf{k}_i \rangle &= \langle \mathbf{W}_Q \mathbf{h}_t^{(0)}, \mathbf{W}_K \mathbf{h}_i^{(0)} \rangle \\ &= \langle \mathbf{W}_Q (\mathbf{e}_t + \mathbf{p}_t), \mathbf{W}_K (\mathbf{e}_i + \mathbf{p}_i) \rangle \\ &= (\mathbf{W}_Q (\mathbf{e}_t + \mathbf{p}_t))^\top (\mathbf{W}_K (\mathbf{e}_i + \mathbf{p}_i)) \\ &= \mathbf{e}_t^\top \mathbf{W}_Q^\top \mathbf{W}_K \mathbf{e}_i + \mathbf{e}_t^\top \mathbf{W}_Q^\top \mathbf{W}_K \mathbf{p}_i \\ &\quad + \mathbf{p}_t^\top \mathbf{W}_Q^\top \mathbf{W}_K \mathbf{e}_i + \mathbf{p}_t^\top \mathbf{W}_Q^\top \mathbf{W}_K \mathbf{p}_i \end{aligned} \quad (10)$$

591 In the learned variant of APE, $\mathbf{p}_j \in \mathbb{R}^d$ is learned during training. In the sinusoidal variant, \mathbf{p}_j is
 592 calculated using a non-parametric function. Specifically, \mathbf{p}_j is computed as:

$$\mathbf{p}_j = [\sin(\omega_1 \cdot j), \cos(\omega_1 \cdot j), \sin(\omega_2 \cdot j), \cos(\omega_2 \cdot j), \dots, \sin(\omega_{d/2} \cdot j), \cos(\omega_{d/2} \cdot j)]^\top \quad (11)$$

593 where $\omega_i = \frac{1}{10000^{2i/d}}$.

594 **T5’s Relative PE** The Relative bias in T5 is a type of relative positional encoding that initially
 595 calculates the relative distance $(t - i)$ between tokens at positions t and i . This distance is then
 596 transformed into a scalar bias value b and is incorporated into the dot product between the query and
 597 key. b is learned during training. Thus, the dot product in every layer can be written as:

$$\langle \mathbf{q}_t, \mathbf{k}_i \rangle = \mathbf{q}_t^\top \mathbf{k}_i + b_{\text{bucket}(n-m)} \quad (12)$$

598 where

$$\text{bucket}(n) = \begin{cases} n & \text{if } n < \frac{\mathcal{B}}{2} \\ \frac{\mathcal{B}}{2} + \left\lfloor \frac{\log\left(\frac{n}{\mathcal{B}/2}\right)}{\log\left(\frac{\mathcal{D}}{\mathcal{B}/2}\right)} \times \frac{\mathcal{B}}{2} \right\rfloor & \text{if } \frac{\mathcal{B}}{2} \leq n < \mathcal{D} \\ \mathcal{B} - 1 & \text{if } n \geq \mathcal{D} \end{cases}$$

599 This function maps the relative distance d to a bucket index, which will be used to look up the weight
 600 corresponding to that bucket. \mathcal{B} is the number of buckets, \mathcal{D} is the maximum distance. It assigns
 601 half of the buckets to distances smaller than $\frac{\mathcal{D}}{2}$ with linear spacing and the other half to distances
 602 larger than $\frac{\mathcal{D}}{2}$ with logarithmic spacing. The weight for distances larger than \mathcal{D} is the same. This is
 603 to facilitate generalization to unseen distances. In the original implementation of T5, $\mathcal{B} = 32$ and
 604 $\mathcal{D} = 128$. Following shows an example of the bucket function with $\mathcal{B} = 5$ and $\mathcal{D} = 6$:

$$\text{bucket}\left(\begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 2 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 3 & 2 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 4 & 3 & 2 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 5 & 4 & 3 & 2 & 1 & 0 & 0 & 0 & 0 & 0 \\ 6 & 5 & 4 & 3 & 2 & 1 & 0 & 0 & 0 & 0 \\ 7 & 6 & 5 & 4 & 3 & 2 & 1 & 0 & 0 & 0 \\ 8 & 7 & 6 & 5 & 4 & 3 & 2 & 1 & 0 & 0 \\ 9 & 8 & 7 & 6 & 5 & 4 & 3 & 2 & 1 & 0 \end{bmatrix}\right) = \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 2 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 3 & 2 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 4 & 3 & 3 & 2 & 1 & 0 & 0 & 0 & 0 & 0 \\ 5 & 4 & 4 & 3 & 3 & 2 & 1 & 0 & 0 & 0 \\ 6 & 4 & 4 & 4 & 3 & 3 & 2 & 1 & 0 & 0 \\ 7 & 4 & 4 & 4 & 4 & 3 & 3 & 2 & 1 & 0 \\ 8 & 4 & 4 & 4 & 4 & 4 & 3 & 3 & 2 & 1 \\ 9 & 4 & 4 & 4 & 4 & 4 & 4 & 3 & 3 & 2 \end{bmatrix}$$

605 **ALiBi** Similar to T5’s Relative PE, ALiBi subtracts a scalar bias from the attention score. As the
 606 distance between the query and key tokens increases, the bias grows linearly. Specifically, the dot
 607 product in every layer can be written as:

$$\langle \mathbf{q}_t, \mathbf{k}_i \rangle = \mathbf{q}_t^\top \mathbf{k}_i - (t - i) \cdot C^{(m+1)} \quad (13)$$

where m is head index and C is a constant defined as:

$$C = 2^{-2^{-\log_2(\#\text{heads}+3)}}$$

608 For example, if the number of heads is 8, then we have $\frac{1}{2}, \frac{1}{2^2}, \dots, \frac{1}{2^8}$ (Press et al., 2022).

609 **Rotary** The Rotary is a relative PE that applies a rotation to the query and key representations
 610 based on their absolute positions before dot product attention. Due to this rotation, the attention dot
 611 product relies solely on the relative distance between tokens.

612 First, we formulate Rotary for model dimension $d = 2$. Rotary positional encoding defines the dot
 613 product as:

$$\begin{aligned}
\langle \mathbf{q}_t, \mathbf{k}_i \rangle &= \langle \text{Rot}(\mathbf{q}_t, t), \text{Rot}(\mathbf{k}_i, i) \rangle \\
&= \langle \mathbf{R}^{t\theta} \mathbf{q}_t, \mathbf{R}^{i\theta} \mathbf{k}_i \rangle \\
&= (\mathbf{R}^{t\theta} \mathbf{q}_t)^\top (\mathbf{R}^{i\theta} \mathbf{k}_i) \\
&= \mathbf{q}_t^\top (\mathbf{R}^{t\theta})^\top \mathbf{R}^{i\theta} \mathbf{k}_i \\
&= \mathbf{q}_t^\top \mathbf{R}^{(i-t)\theta} \mathbf{k}_i
\end{aligned} \tag{14}$$

614 where $\mathbf{R}^{t\theta}$ is a rotation matrix that rotates \mathbf{x} by $t\theta$ radians:

$$\mathbf{R}^{n\theta} = \begin{bmatrix} \cos(n\theta) & -\sin(n\theta) \\ \sin(n\theta) & \cos(n\theta) \end{bmatrix} \tag{15}$$

615 for $d > 2$, Rotary applies the same approach on every two consecutive dimensions of \mathbf{q}_t and \mathbf{k}_i , but
616 with different θ angles. Refer to [Su et al. \(2021\)](#) for the exact formulation.

617 **NoPE** NoPE does not explicitly encode positional encodings. So, the dot product in every layer
618 can be written as:

$$\langle \mathbf{q}_t, \mathbf{k}_i \rangle = \mathbf{q}_t^\top \mathbf{k}_i \tag{16}$$

619 C Proofs

620 In this section, we provide proof of why NoPE can implicitly learn both absolute and relative positions.
621 We refer the readers to [Appendix B.1](#) for the notation and definitions used in this section.

622 C.1 Absolute Positional Encoding in NoPE

623 This section discusses how NoPE can recover absolute positions in the hidden state. Our proof is
624 inspired by [Weiss et al. \(2021\)](#); [Lindner et al. \(2023\)](#) and relies on the causal attention mask in the
625 decoder-only Transformer and the softmax function to recover absolute positions.

Theorem 1 (Absolute Encoding). *Let $\mathbf{x} = [\langle \text{bos} \rangle, x_1, \dots, x_T]$ be an input sequence of length $T + 1$ to the model. Then, the first layer of f_θ can recover absolute positions $[1, \dots, T + 1]$ in the hidden state $\mathbf{H}^{(1)}$. That is, there exist $\mathbf{W}_Q, \mathbf{W}_K, \mathbf{W}_V, \mathbf{W}_O, \mathbf{W}_1,$ and \mathbf{W}_2 such that the self-attention and feedforward operations in the first layer compute absolute positions and write it to the next hidden state.*

626 *Proof.*

627 Our proof only specifies the weights of a single attention head in the first layer (and additionally the
628 parameterization of feedforward sub-layer). In this parameterization, we only require the first three
629 dimensions of the hidden states. The rest of the heads, as long as they do not override the first three
630 dimensions, can be arbitrary. This does not impose any challenges as Transformers used in practice
631 usually have a very large model dimension d . In the rest, we provide the construction of the weights
632 and then verify that they can recover absolute positions.

633 First, we construct the word embedding matrix $\mathbf{W}_E \in \mathbb{R}^{d \times V}$, where each column is the embedding
634 of a token in the vocabulary. We construct \mathbf{W}_E such that it always sets the first dimension of every
635 embedding vector to be 1. Additionally, it sets the second dimension to 1 if and only if the token
636 is $\langle \text{bos} \rangle$. Otherwise, it sets it to zero. The third dimension of all embedding vectors is set to zero.
637 Other dimensions can take any arbitrary values. Without loss of generality, assume $\langle \text{bos} \rangle$ is the first
638 token in the vocabulary, i.e. The first column. Then, we have:

$$\mathbf{W}_E = \begin{bmatrix} 1 & 1 & 1 & \dots & 1 \\ 1 & 0 & 0 & \dots & 0 \\ 0 & 0 & 0 & \dots & 0 \\ e_{4,1} & e_{4,2} & e_{4,3} & \dots & e_{4,V} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ e_{d,1} & e_{d,2} & e_{d,2} & \dots & e_{d,V} \end{bmatrix}_{d \times V} \tag{17}$$

639 where $e_{d,i} \in \mathbb{R}$.

640 Secondly, for head dimensions $h \geq 1$, we construct the weights $\mathbf{W}_Q, \mathbf{W}_K, \mathbf{W}_V, \mathbf{W}_O$ of the first
641 attention head in the first layer. Specifically,

$$\mathbf{W}_K = \begin{bmatrix} 1 & 0 & \dots & 0 \\ 1 & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 1 & 0 & \dots & 0 \end{bmatrix}_{h \times d} \quad \mathbf{W}_V = \begin{bmatrix} 0 & 1 & 0 & \dots & 0 \\ 0 & 0 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & 0 \end{bmatrix}_{h \times d} \quad (18)$$

642 \mathbf{W}_K reads from the first dimension of the hidden state, which is initialized with 1 using the embedding
643 matrix. Since all word embeddings have one in their first dimension, this parameterization will result
644 all key vectors to be the same. Moreover, \mathbf{W}_V reads from the second dimension of the hidden state,
645 which is initialized with 1 if the token is $\langle \text{bos} \rangle$. So, the value vector will have 1 in its first dimension
646 only if the corresponding token is $\langle \text{bos} \rangle$.

647 \mathbf{W}_Q can be any arbitrary matrix. \mathbf{W}_O will write the result of the attention to the third dimension of
648 the hidden state and can be constructed as:

$$\mathbf{W}_O = \begin{bmatrix} 0 & 0 & 0 & 0 & \dots & 0 \\ 0 & 0 & 0 & 0 & \dots & 0 \\ 1 & 0 & 0 & 0 & \dots & 0 \\ 0 & 0 & 0 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & 0 & \dots & 0 \end{bmatrix}_{d \times h} \quad (19)$$

649 Now, we verify that for any input sequence $\mathbf{x} = [\langle \text{bos} \rangle, x_1, \dots, x_T]$, the first layer can recover
650 absolute positions $[1, \dots, T+1]$ in the hidden state $\mathbf{H}^{(1)}$. We verify this for column t of $\mathbf{H}^{(1)}$. That
651 is, we show that absolute position information is available in the third dimension of $\mathbf{h}_t^{(1)}$.

652 First, we use the word embedding matrix \mathbf{W}_E to compute the embedding $\mathbf{H}^{(0)}$:

$$\mathbf{H}^{(0)} = \mathbf{W}_E \mathbf{X} = \begin{bmatrix} 1 & 1 & 1 & \dots & 1 \\ 1 & 0 & 0 & \dots & 0 \\ 0 & 0 & 0 & \dots & 0 \\ e_{4,1} & e_{4,2} & e_{4,3} & \dots & e_{4,V} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ e_{d,1} & e_{d,2} & e_{d,2} & \dots & e_{d,V} \end{bmatrix}_{d \times (T+1)} \quad (20)$$

653 We now provide the attention computation at position $1 \leq t \leq T+1$. First, we use \mathbf{W}_Q to compute
654 the query vector \mathbf{q}_t by applying $\mathbf{q}_t = \mathbf{W}_Q \mathbf{h}_t^{(0)}$:

$$\mathbf{q}_t = [q_1, q_2, q_3, \dots, q_h]^\top \quad (21)$$

655 Recall that \mathbf{W}_Q can be any arbitrary matrix. So, $q_j \in \mathbb{R}$ can take any arbitrary value. Next, we
656 compute the key vectors by applying $\mathbf{k}_i = \mathbf{W}_K \mathbf{h}_i^{(0)}$:

$$\mathbf{k}_1 = \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix} \quad \mathbf{k}_2 = \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix} \quad \dots \quad \mathbf{k}_t = \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix} \quad (22)$$

657 Note that all key vectors are the same and we only need to compute them up to position t as the
658 attention mask is causal, i.e query can only look at positions $\leq t$. Next, we compute the attention
659 weight vectors $\boldsymbol{\alpha}$:

$$\boldsymbol{\alpha} = [\langle \mathbf{q}_t, \mathbf{k}_1 \rangle, \langle \mathbf{q}_t, \mathbf{k}_2 \rangle, \dots, \langle \mathbf{q}_t, \mathbf{k}_t \rangle]^\top \quad (23)$$

$$= [\alpha^*, \alpha^*, \dots, \alpha^*]^\top \quad (24)$$

660 where $\alpha^* = q_1 + q_2 + \dots + q_h$. Next, we apply softmax to compute the attention probabilities.
 661 Since all α^i 's are the same, we have:

$$\hat{\alpha} = \text{softmax}(\alpha) = \left[\frac{1}{t}, \frac{1}{t}, \dots, \frac{1}{t} \right]^\top \quad (25)$$

662 Now, we compute the value vectors by applying $v_i = \mathbf{W}_V \mathbf{h}_i^{(0)}$:

$$\mathbf{v}_1 = \begin{pmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix} \quad \mathbf{v}_2 = \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{pmatrix} \quad \dots \quad \mathbf{v}_t = \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{pmatrix} \quad (26)$$

663 Finally, we compute the output of the attention head by applying \mathbf{W}_O :

$$\mathbf{o}_t = \mathbf{W}_O \left(\sum_{i \leq t} \hat{\alpha}_i \mathbf{v}_i \right) = \mathbf{W}_O \left(\frac{1}{t} \sum_{i \leq t} \mathbf{v}_i \right) = \mathbf{W}_O \begin{pmatrix} 1/t \\ 0 \\ \vdots \\ 0 \end{pmatrix}_h = \begin{pmatrix} 0 \\ 0 \\ 1/t \\ 0 \\ \vdots \\ 0 \end{pmatrix}_d \quad (27)$$

664 Thus, the output of our constructed attention head recovers the absolute position information and
 665 writes it to the third dimension of output.

666 We used the decoder-only property of Transformer implicitly in Equation (23), which helped us to
 667 only attend to position $\leq t$. So, the lengths of the attended sequence are always t . Moreover, the
 668 presence of <bos> token in the input sequence helped us to anchor the absolute position information.
 669 This is not a problem as in practice models are often prompted with some instructions which can act
 670 as <bos> token.

671 With this information available to the rest of the network, the feedforward sub-layer, with sufficient
 672 hidden width, can recover the absolute positions $[1, 2, \dots, T+1]$ from the third dimension of attention
 673 output. This is because the feedforward sub-layer is MLP with ReLU activation. So, it can learn
 674 any arbitrary function (Park et al., 2020). Note that the layer-norm operation can be bypassed as
 675 explained by Akyurek et al. (2023). \square

676 C.2 Relative Positional Encoding in NoPE

677 In this section, we show if the hidden state contains absolute positional information as explained in
 678 the previous section, then the attention mechanism in all subsequent layers can implement a relative
 679 positional encoding. We refer the readers to Appendices B.1 and C.1 for the notation and definitions
 680 used in this section.

Theorem 2 (Relative Encoding). *Suppose that the hidden state $\mathbf{H}^{(1)}$ contains absolute positional information, as stated in Theorem 1, and assume that it is not overwritten by any subsequent layers. Then, the self-attention in all subsequent layers can implement a relative positional encoding: there exists a parameterization of f_θ such that, for $l \geq 2$, the attention dot product between query \mathbf{q}_t and key \mathbf{k}_i at positions t and i ($t \geq i$) can be expressed as:*

$$\langle \mathbf{q}_t, \mathbf{k}_i \rangle = f_{\text{cnt}}(\mathbf{q}_t, \mathbf{k}_i) + f_{\text{rel}}(t - i) \quad (1)$$

where f_{cnt} is a function of their content, and f_{rel} is a function of their relative distance.

681 *Proof.*

682 Our proof only specifies a few entries of weight matrices for attention heads in layers $l \geq 2$, which
 683 does not impose any challenges for Transformers used in practice as they usually have a very large
 684 model dimension d . Moreover, we require to have absolute positions in the third dimension of the
 685 hidden state as explained in Theorem 1. To show NoPE can implement relative encoding, we only
 686 need to prove that its attention dot product depends on the relative distance between tokens (See

687 Appendix B.1 for an overview of relative encoding methods). In the rest, we provide the construction
 688 of the weights and then verify that they can implement relative position encoding.

689 For head dimension $h \geq 2$, we construct the weights $\mathbf{W}_Q, \mathbf{W}_K$ of the attention heads in the second
 690 layers and above. Specifically,

$$\mathbf{W}_Q = \begin{bmatrix} 1 & 0 & 0 & 0 & \dots & 0 \\ 0 & 0 & -1 & 0 & \dots & 0 \\ w_{3,1} & w_{3,2} & w_{3,3} & w_{3,4} & \dots & w_{3,d} \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ w_{h,1} & w_{h,2} & w_{h,3} & w_{h,4} & \dots & w_{h,d} \end{bmatrix}_{h \times d} \quad (28)$$

691

$$\mathbf{W}_V = \begin{bmatrix} 0 & 0 & 1 & 0 & \dots & 0 \\ 1 & 0 & 0 & 0 & \dots & 0 \\ w'_{3,1} & w'_{3,2} & w'_{3,3} & w'_{3,4} & \dots & w'_{3,d} \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ w'_{h,1} & w'_{h,2} & w'_{h,3} & w'_{h,4} & \dots & w'_{h,d} \end{bmatrix}_{h \times d} \quad (29)$$

692 where $w_{i,j}, w'_{i,j} \in \mathbb{R}$ can take any arbitrary value. Their corresponding \mathbf{W}_V and \mathbf{W}_O can take any
 693 arbitrary values as long as they do not override the first three dimensions of the residual stream.

694 Now we verify that for any input sequence $\mathbf{x} = [\langle \text{bos} \rangle, x_1, \dots, x_T]$, the attention dot product
 695 between query \mathbf{q}_t and key \mathbf{k}_i at positions t and i ($t \geq i$) will depend the relative distance between
 696 tokens.

697 First, assume that absolute positions are computed in the hidden state $\mathbf{H}^{(l)}$ for $l \geq 1$, as stated in
 698 Theorem 1. Specifically,

$$\mathbf{H}^{(l)} = \begin{bmatrix} 1 & 1 & 1 & 1 & \dots & 1 \\ 1 & 0 & 0 & 0 & \dots & 0 \\ 1 & 2 & 3 & 4 & \dots & T+1 \\ h_{4,1} & h_{4,2} & h_{4,3} & h_{4,4} & \dots & h_{4,T+1} \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ h_{d,1} & h_{d,2} & h_{d,3} & h_{d,4} & \dots & h_{d,T+1} \end{bmatrix}_{d \times (T+1)} \quad (30)$$

699 where $h_{i,j} \in \mathbb{R}$ can be any arbitrary value as the first three dimensions of the hidden state are
 700 reserved for PE computation. The rest of the dimensions can take any arbitrary values as in regular
 701 computation of Transformers.

702 We now present the attention computations at position $1 \leq t \leq T+1$. We use \mathbf{W}_Q to compute the
 703 query vector \mathbf{q}_t by applying $\mathbf{q}_t = \mathbf{W}_Q \mathbf{h}_t^{(l)}$:

$$\mathbf{q}_t = [1, -t, q_3, \dots, q_h]^\top \quad (31)$$

704 where $q_j \in \mathbb{R}$ can take any arbitrary value. Next, we compute the key vectors by applying $\mathbf{k}_i =$
 705 $\mathbf{W}_K \mathbf{h}_i^{(l)}$:

$$\mathbf{k}_1 = \begin{pmatrix} 1 \\ 1 \\ k_{3,1} \\ \vdots \\ k_{h,1} \end{pmatrix} \quad \mathbf{k}_2 = \begin{pmatrix} 2 \\ 1 \\ k_{3,2} \\ \vdots \\ k_{h,2} \end{pmatrix} \quad \mathbf{k}_3 = \begin{pmatrix} 3 \\ 1 \\ k_{3,3} \\ \vdots \\ k_{h,3} \end{pmatrix} \quad \dots \quad \mathbf{k}_t = \begin{pmatrix} t \\ 1 \\ k_{3,t} \\ \vdots \\ k_{h,t} \end{pmatrix} \quad (32)$$

706 where $k_{(\cdot, \cdot)} \in \mathbb{R}$ can have any arbitrary value. So, for \mathbf{k}_i we have:

$$\mathbf{k}_i = [i, 1, k_{3,i}, \dots, k_{h,i}]^\top \quad (33)$$

707 Next, we let us present the attention dot product between \mathbf{q}_t and \mathbf{k}_i :

$$\langle \mathbf{q}_t, \mathbf{k}_i \rangle = 1 \cdot i + (-t) \cdot 1 + q_3 \cdot k_{3,i} + \dots + q_h \cdot k_{h,i} \quad (34)$$

$$= i - t + \sum_{j=3}^h q_j \cdot k_{j,i} \quad (35)$$

$$= \left(\sum_{j=3}^h q_j \cdot k_{j,i} \right) - (t - i) \quad (36)$$

$$= f_{\text{cnt}}(\mathbf{q}_t, \mathbf{k}_i) + f_{\text{rel}}(t - i) \quad (37)$$

708 Thus, the dot product between \mathbf{q}_t and \mathbf{k}_i depends on the relative distance between tokens (assuming
 709 the rest of the terms do not cancel out which can be easily avoided by setting the respective weights
 710 in Equations (28) and (29)). Note that our proof uses the linear spacing between tokens, but the MLP
 711 the first layer can write any arbitrary function of absolute positions to the third dimension of the
 712 hidden state, which enables more complex relative encoding schemes. \square

713 D Experimental Details

714 D.1 Tasks

715 Here we provide the details and more examples of the tasks and datasets we used in our evaluation.
 716 For each task, we sample 100K examples for the training set and 10K for the test. Also, we use 15%
 717 of the train as the validation set.

718 **Addition** The addition task (Nye et al., 2021) asks the model to compute the sum of two numbers.
 719 Each number is represented as a sequence of digits that are separated by space. So, the model has
 720 access to the exact digits.

```
Input
Compute: 5 3 7 2 6 + 1 9 1 7 ?
Output
The answer is 5 5 6 4 3.
```

721 we create each length bucket based on the number of digits in each number, e.g. 6-by-3, 6-by-4,
 722 etc. For the training set, we use buckets where one of the numbers has at most L digits. For the test
 723 set, we use buckets where any of the numbers have at most L digits. The model is evaluated on the
 724 correctness of its predicted result.

725 **Polynomial Evaluation** The polynomial evaluation task (Nye et al., 2021) asks the model to
 726 evaluate a polynomial expression at a given value of x . The polynomial terms and digits are separated
 727 to make just the tokenizer does not glue symbols together.

```
Input
Evaluate x = 3 in ( 3 x ** 0 + 1 x ** 1 + 1 x ** 2 ) % 10 ?
Output
The answer is 5.
```

728 The length bucket is created based on the number of terms in the polynomial expression. We sample x
 729 from $\mathcal{U}(-2, 2)$, the degree of each term from $\mathcal{U}(0, 3)$, and the coefficient of each term from $\mathcal{U}(-3, 3)$.
 730 We take the modulo of the result by 10 to make the task easier for the model and make sure we only
 731 measure the generalization of the length of the problem instance not the value of the polynomial. The
 732 model is evaluated on the correctness of its predicted result.

733 **Sorting** The sorting task (Saxton et al., 2019) asks the model to sort a sequence of input numbers.
 734 We use this task in two variants: Single Token and Multi Digit. In the Single Token variant, we create
 735 an alphabet of 50 tokens from the model’s vocabulary and fix some canonical ordering among them

736 through task. Each instance is a sequence of tokens from the alphabet in a random order, and the
737 model is asked to sort them in the canonical order.

```
Input
Sort the following numbers: 3 1 4 1 5 ?
Output
The answer is 1 1 3 4 5.
```

738 In the Multi Digit variant, we simply present a sequence of multi digit/tokens numbers to the model,
739 and ask it to sort them in ascending order. Each number is represented by its digits and they are
740 separated by a space.

```
Input
Sort the following numbers: 3 1, 4 1, 5 9, 1 2 6, 5 3 3 ?
Output
The answer is 3 1, 4 1, 5 9, 1 2 6, 5 3 3.
```

741 In this case, we sample each number from $\mathcal{U}(0, 10000)$. In both cases, the length bucket is created
742 based on the length of the input sequence. The model is evaluated on the correctness of its predicted
743 result.

744 **Summation** In this task (Saxton et al., 2019), we ask the model to compute the sum of a sequence
745 of input numbers modulo 10 as we want to specifically measure how the model generalizes to longer
746 sequences not the value of summation result:

```
Input
Compute: ( 1 + 2 + 3 + 4 + 7 ) % 10 ?
Output
The answer is 7.
```

747 Each digit is randomly sampled from $\mathcal{U}(1, 9)$. The length bucket is created based on the length of the
748 input sequence. The model is evaluated on the correctness of its predicted result.

749 **Parity** In the parity task (Anil et al., 2022), we ask the model to compute the parity of a binary
750 sequence.

```
Input
Is the number of 1's even in [ 1 0 0 1 1 ] ?
Output
The answer is No.
```

751 **LEGO** In the LEGO task (Zhang et al., 2023), the model is provided with a simple computation
752 graph (DAG), where each node represents a variable, and variables are connected by simple operations
753 which created the edges in the computation graph. We refer to Zhang et al. (2023) for a detailed
754 description.

```
Input
If a = -1; b = -a; c = +b; d = +c. Then what is c?
Output
The answer is +1.
```

755 To sample each example, we first sample the list of variables based on the length of the example, and
756 then we uniformly sample the value of each variable to make sure all variables are represented with
757 both -1 and +1. Finally, given the value of variables, we deterministically compute the operation on
758 each edge. For each example, we query all variables from the middle of the computation graph to the
759 end. The model is evaluated on the correctness of its predicted result.

760 **Copy** The copy task is straightforward. The model has to repeat the input sequence in the output.

Input

Copy the following words: <w1> <w2> <w3> <w4> <w5> .

Output

<w1> <w2> <w3> <w4> <w5>

761 We create multiple variants of this task to better understand the models' generalization behavior. In
762 the first variant, the input tokens are the same, so the model has to basically count the number of input
763 tokens. In the second variant, the model has to replace the input tokens with another token sampled
764 from the vocabulary. In the third variant, we sample the input tokens from the model's vocabulary,
765 and the model has to predict them in the same order. We also create 2x versions of variants 1 and 3 to
766 make the tasks more challenging.

767 **Reverse** In this task the model, the model has to reverse the order of input tokens in its output.

Input

Reverse the following words: <w1> <w2> <w3> <w4> <w5> .

Output

<w5> <w4> <w3> <w2> <w1> .

768 As in the copy task, we create multiple variants of this task. In the first variant, the model has to
769 reverse the order of input tokens, where the tokens are randomly sampled from the model's vocabulary.
770 In the second variant, the model has to reverse the order of input tokens, as in the first variant, but
771 also it has to reverse it one more time, recreating the original input.

772 D.2 Hyperparameters

773 Table 2 shows the hyperparameters we used in our experiments. We use the same hyperparameters
774 for all models and positional encoding schemes. In our initial experiment, we tried a few more
775 hyperparameters such as $lr \in \{0.00001, 0.00003, 0.00005\}$ and $WeightDecay \in \{0, 0.05, 0.1\}$,
776 but we did not observe any significant difference in the results. So, we decided to use the same
777 hyperparameters throughout our experiments.

778 D.3 Compute

779 In our experiments, we used single-GPU training setup for the models. Specifically, we ran our
780 experiments on a mix of NVIDIA V100 32G, NVIDIA RTX8000 48G, NVIDIA A100 40G, and
781 NVIDIA A100 80G GPUs. Depending on the GPU type and the positional encoding, each of our
782 training runs took 6 to 15 hours, per each seed, on average to complete. Considering all the datasets,
783 and positional encoding schemes, in addition to the scratchpad experiments, and three seeds, we ran
784 about 870 individual training runs for results in this paper.

785 D.4 Reproducibility

786 In this study, all experiments employed open-source libraries, specifically HuggingFace (Wolf et al.,
787 2020) from which we utilized their implementation as a foundation for the training loop, optimizer,
788 and the Transformer architecture. To ensure the reproducibility, we will also release a singularity
789 binary with all dependencies and libraries to enable running our experiments on any machine with
790 NVIDIA GPUs and at anytime in the future. Moreover, every reported number in this paper is linked
791 to source code package that deterministically (up to GPU stochasticity) reproduces the results, which
792 we release them publicly on GitHub at <https://www.omitted.link>.

793 E Full Results

794 E.1 Detailed Model Accuracy

795 We report the detailed results of our experiments in Figures E.2 to E.4. We refer the readers to
796 Appendix D.1 for the description of each task.

Table 2: Summary of hyperparameters used in the experiments.

Parameter	Value
Optimizer	AdamW
Learning rate	0.00003
Weight Decay	0.05
Batch size	64
Learning Rate Scheduler	Polynomial
Warm Up	6% of training steps
# Train Steps	40K steps
Dropout (<i>taken from HuggingFace</i>)	0.1
Model dimension (<i>taken from HuggingFace</i>)	768
# Layers (<i>taken from HuggingFace</i>)	12
# Attention Heads (<i>taken from HuggingFace</i>)	12

797 **E.2 Detailed Head Distance**

798 [Figure E.5](#) shows the layer-wise distance of No PE’s attention patterns with other positional encoding
 799 schemes measured across instances of the SCAN dataset. We refer the readers to [Section 5.2](#) for the
 800 details and analysis of these results.

801 **E.3 Detailed Model Accuracy On Various Scratchpad Formats**

802 [Figure E.6](#) shows the generalization of various scratchpad formats for each model aggregated across
 803 all datasets. [Figures E.7 to E.13](#) show the generalization of various scratchpad formats for each model
 804 on each dataset. We refer the readers to [Section 6](#) for the details and analysis of these results.

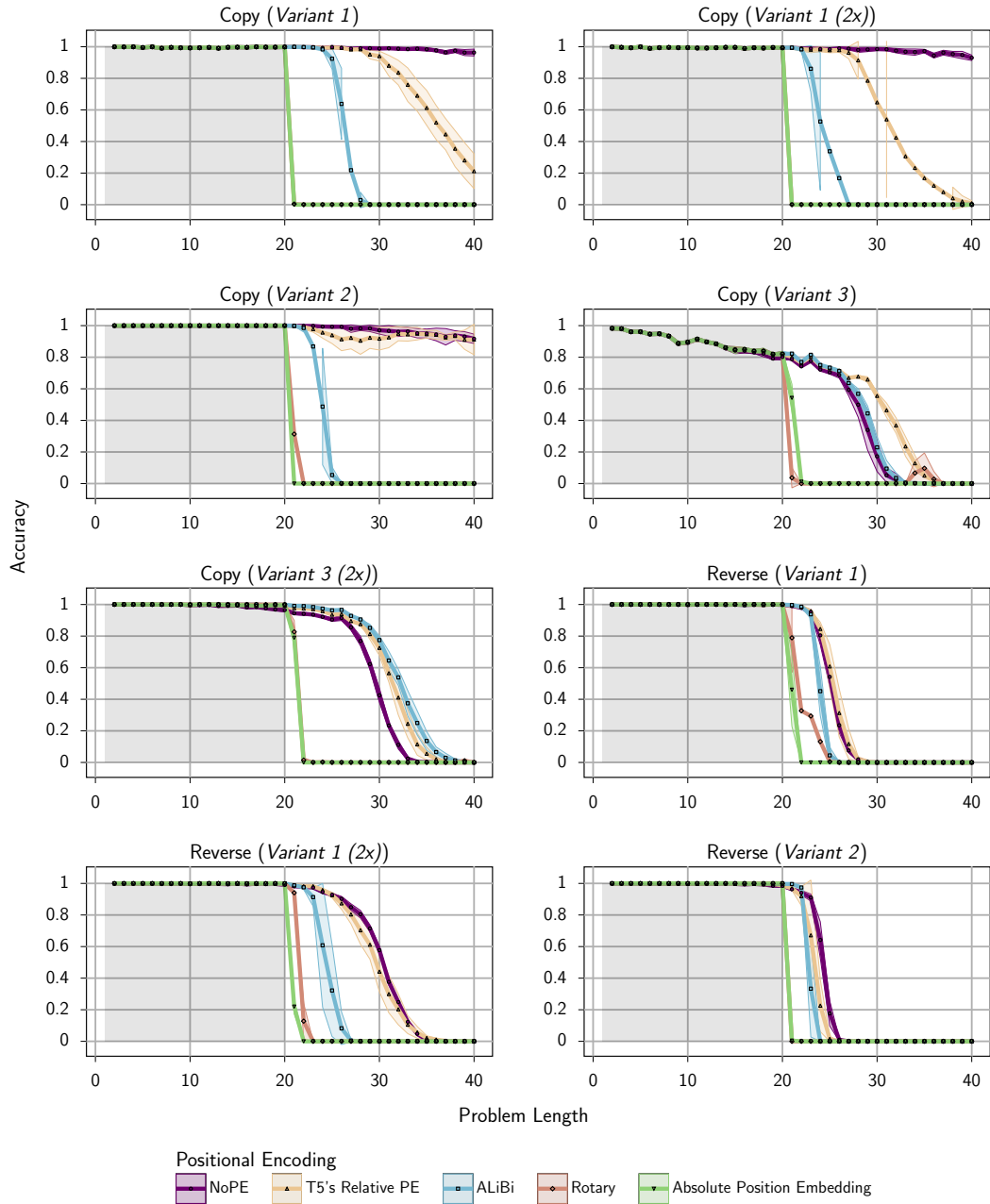


Figure E.2: Generalization behavior of positional encoding schemes on Primitive tasks.

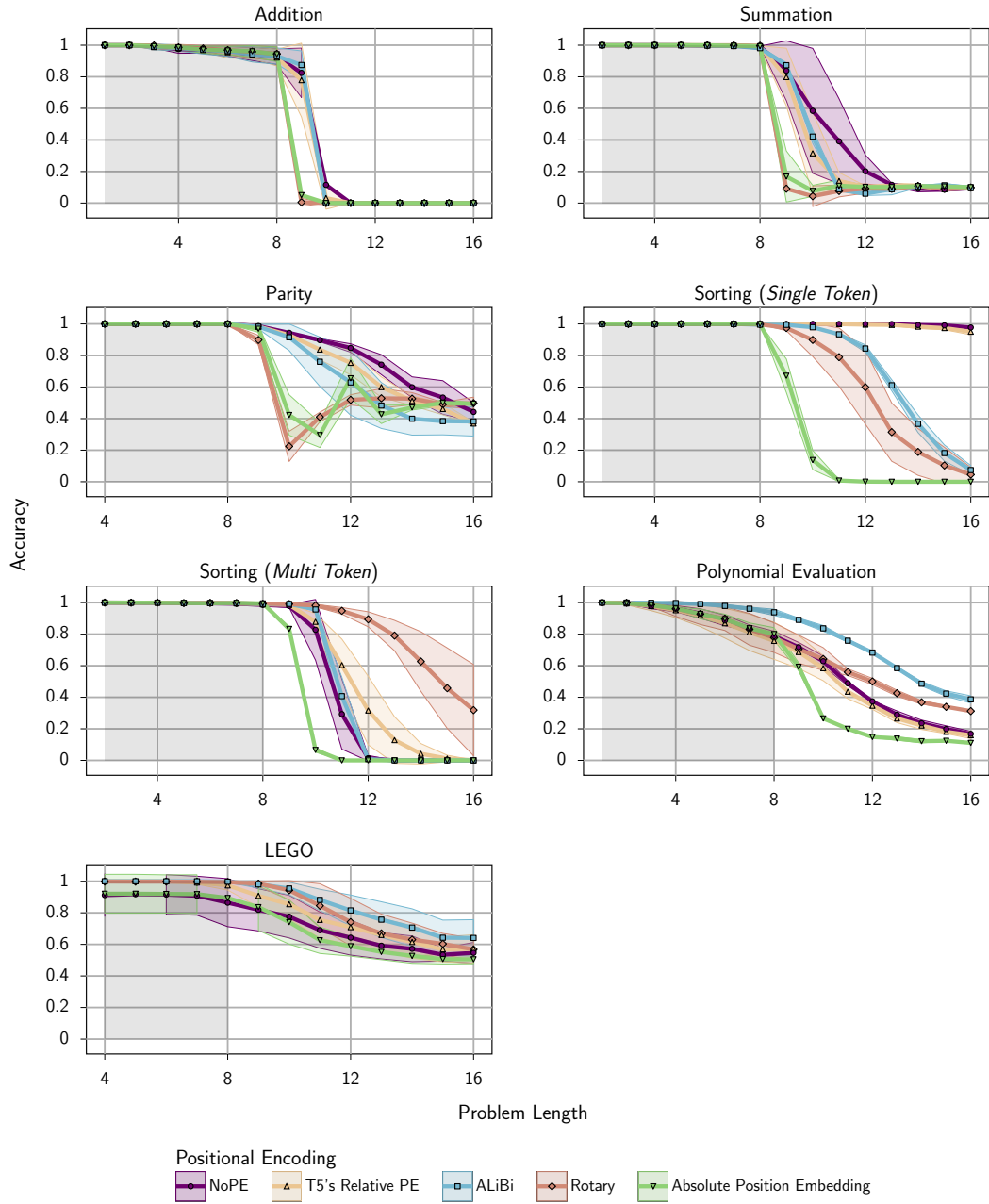


Figure E.3: Generalization behavior of positional encoding schemes on Mathematical & Reasoning tasks.

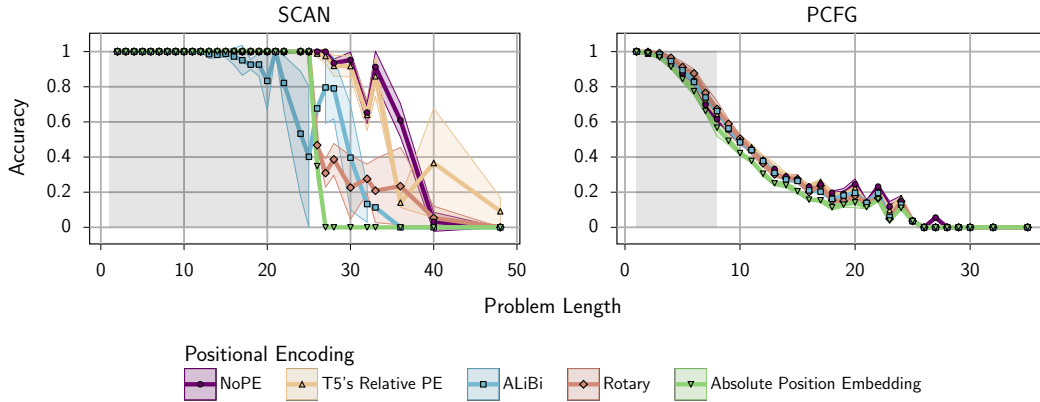


Figure E.4: Generalization behavior of positional encoding schemes on Classic Length Generalization tasks.

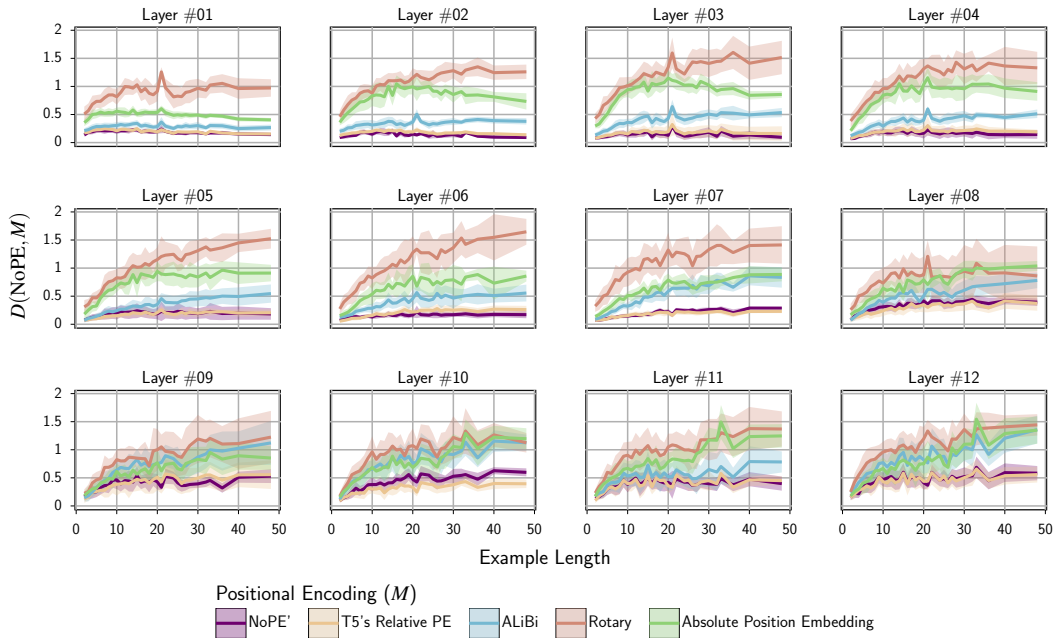


Figure E.5: Layer-wise distance of No PE's attention patterns with other positional encoding schemes measured across instances of SCAN dataset.

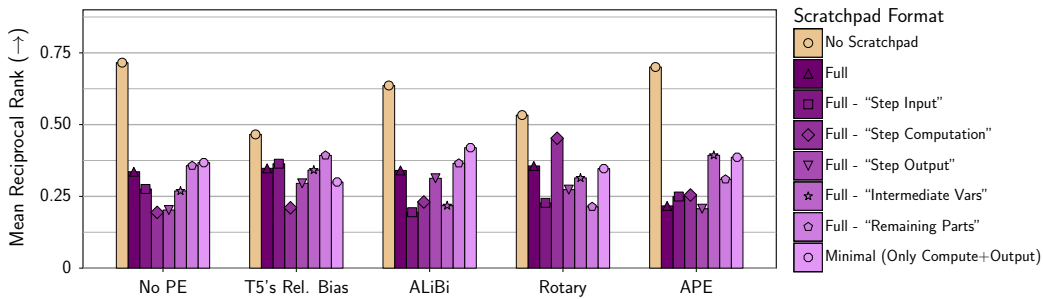


Figure E.6: The optimal scratchpad format is aggregated across all datasets per each model. The optimal scratchpad format is different for each model.

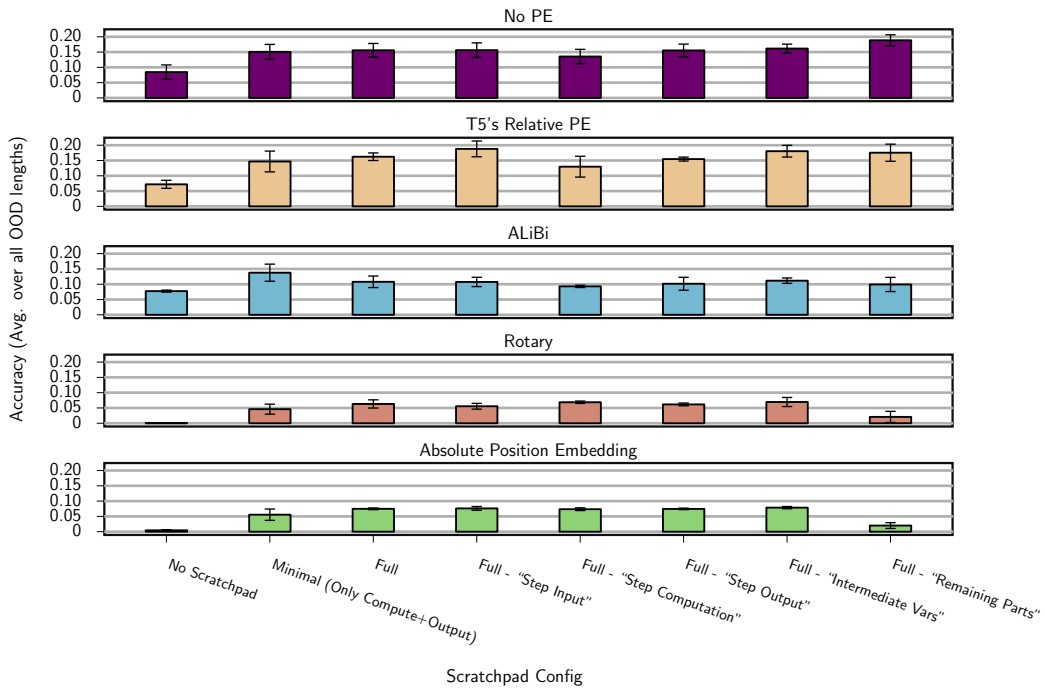


Figure E.7: Generalization of various scratchpad formats for each model on the **Addition** task.

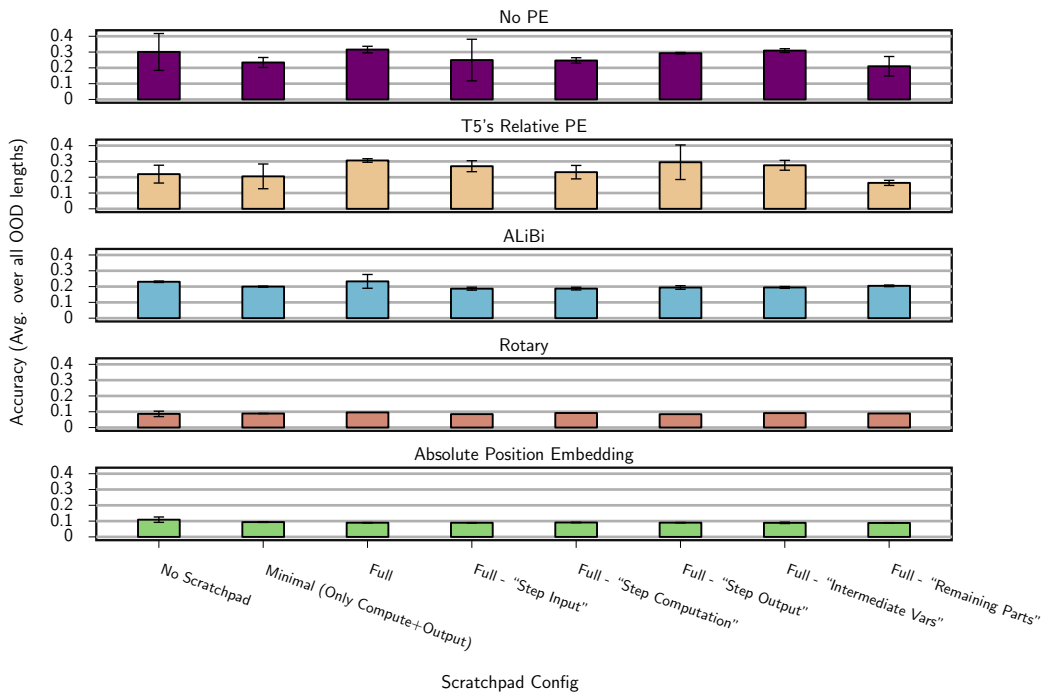


Figure E.8: Generalization of various scratchpad formats for each model on the **Summation** task.

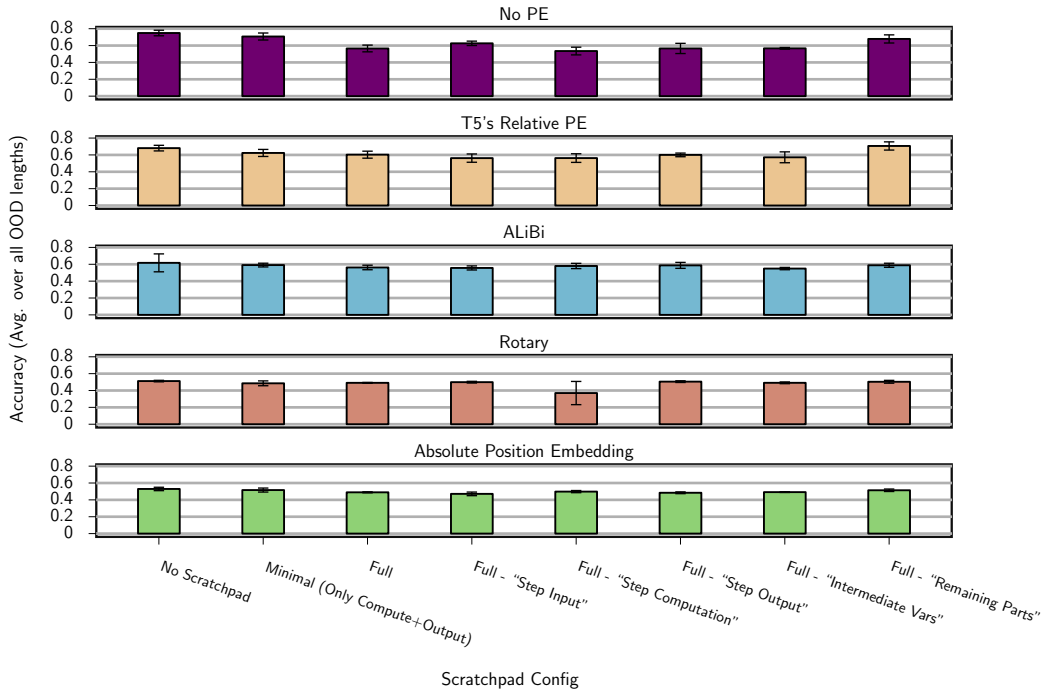


Figure E.9: Generalization of various scratchpad formats for each model on the **Parity** task.

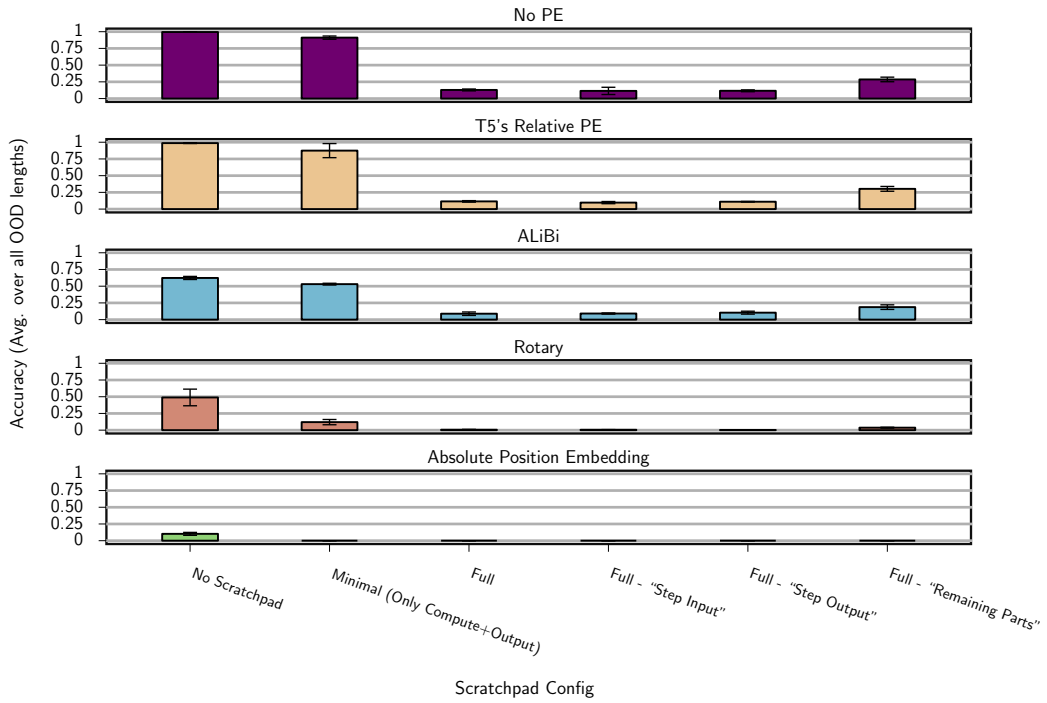


Figure E.10: Generalization of various scratchpad formats for each model on the **Sorting** task (Single Digit).

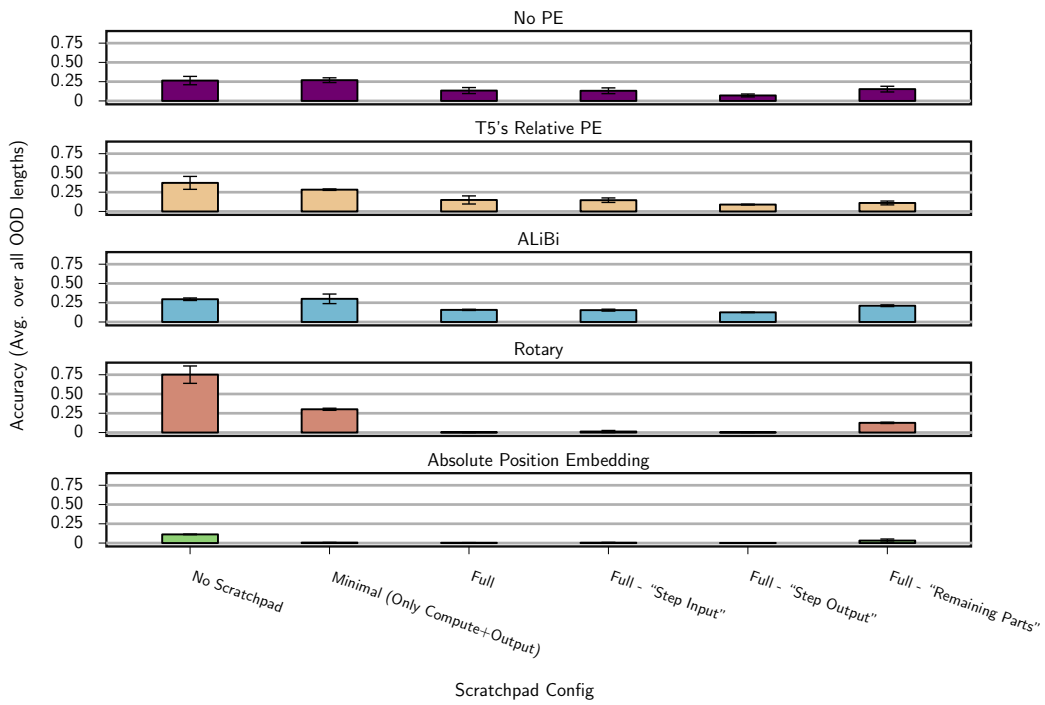


Figure E.11: Generalization of various scratchpad formats for each model on the **Sorting** task (Multi Digit).

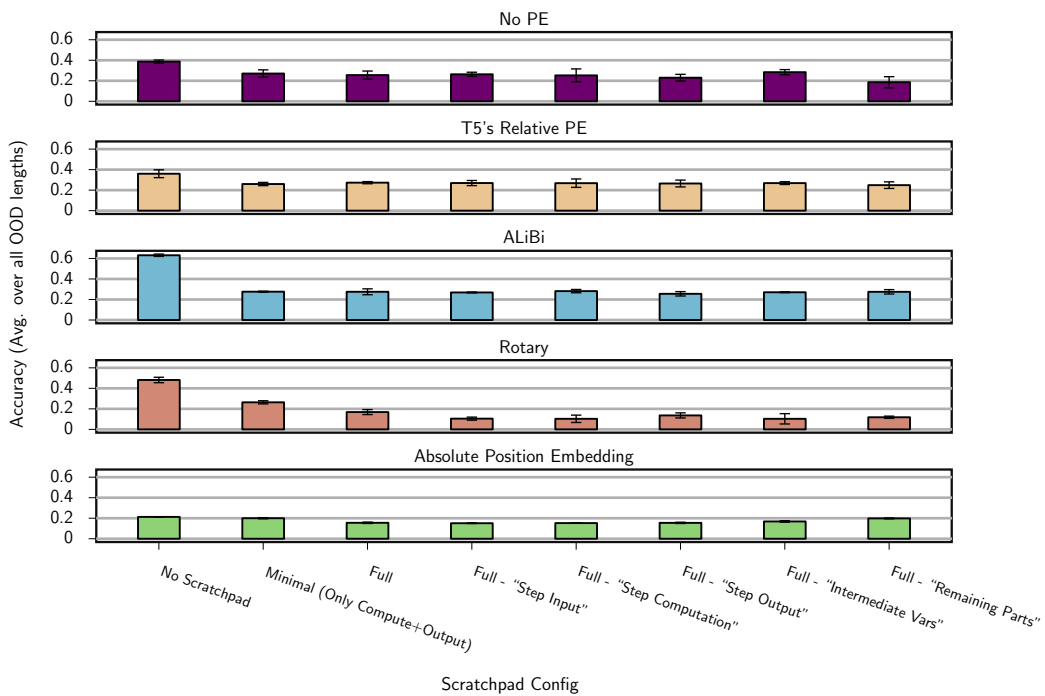


Figure E.12: Generalization of various scratchpad formats for each model on the **Polynomial Evaluation** task.

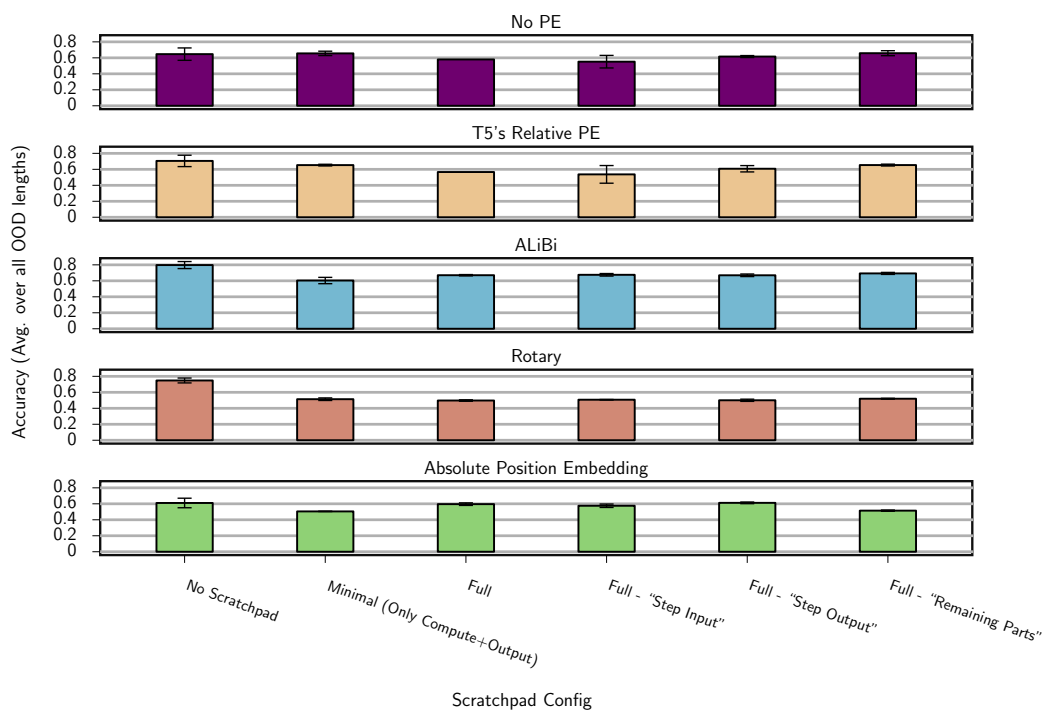


Figure E.13: Generalization of various scratchpad formats for each model on the **LEGO** task.