# ENTP: ENCODER-ONLY NEXT TOKEN PREDICTION

Anonymous authors	
-------------------	--

005 006 007

008 009

010

011

012

013

014

015

016

017

018

019 020 021 Paper under double-blind review

## ABSTRACT

Next-token prediction is conventionally done using decoder-only Transformers with causal attention, as this approach allows for efficient reuse of keys and values. What if we were not compute-limited, should we still use decoder-only Transformers? In this work, we introduce Encoder-only Next Token Prediction (ENTP). We use small scale experiments to explore the differences between ENTP and decoders, highlighting potential advantages of ENTP in setting with unbounded compute. We introduce the Count3 task and show, both theoretically and experimentally, that while ENTP can perform this task easily, a decoder-only Transformer cannot. Finally, we empirically demonstrate ENTP's superior performance across various synthetic tasks, such as length generalization and in-context learning. learning.

1 INTRODUCTION

Traditionally, auto-regressive language modeling has relied on decoder-only Transformers (Vaswani et al., 2017) with causal attention, trained using the next-token prediction objective. Causal attention ensures that each token can only attend to previous tokens, preventing future tokens from influencing past outputs. This mechanism makes training and inference more efficient, as past keys and values do not need to be recomputed for each token. This efficiency enables the scaling of decoder-only Transformers, such as GPT-4 (Achiam et al., 2023) and Llama-3 (Dubey et al., 2024), up to billions of parameters using current hardware.

However, causal attention also introduces artificial constraints. Given tokens  $x_1, x_2, ..., x_n$ , the contextual embedding of  $x_j$  (where j < n) can only attend to embeddings of earlier tokens, even when predicting  $x_{n+1}$ . While this constraint ensures a strict causal structure, it may not always be necessary or beneficial. We investigate what happens when we remove this constraint, while still maintaining causality externally.

Specifically, we look at Encoder-only Transformers, which are typically used for tasks like classification, and do not impose this causality constraint. Though traditionally not used for auto-regressive tasks, encoder-only architectures can be adapted for next-token prediction. When computing the output at the current time step, an encoder-only Transformer, or any sequence model, can be made causal by only providing inputs up to and including the current time step. Therefore, in this work, we investigate the idea of using encoder-only Transformers for next-token prediction. We summarize our findings below.

042

Functions expressible with Decoder-only and Encoder-only Transformers. We demonstrate that the sets of functions expressible by decoder-only and encoder-only Transformers are not comparable, which goes against intuition that the expressivity of encoders would subsume that of decoders. Rather, there exist functions expressible with decoder-only Transformers that are not expressible with encoder-only Transformers, and vice versa, as well as functions expressible by both architectures.

Complexity of Decoder-only and Encoder-only Transformers. Based on the minimum time and
 space complexities, we give a description of the functions that can be performed by decoder-only and
 encoder-only Transformers. We propose an auto-regressive task that can be performed by encoder only Transformers, and cannot be performed by decoder-only Transformers (given that some mild assumptions hold). We validate our hypothesis with small experiments using small decoder-only



Figure 1: An overview of decoder-only vs. encoder-only Transformers in next token prediction. Decoders use a causal attention, ensuring that each token attends only to the preceding tokens. In contrast, encoders allow all tokens to attend to each other by performing attention computation from scratch for each token prediction.

and encoder-only Transformers, as well as experiments fine-tuning GPT-40 (Achiam et al., 2023) and Llama3-8B (Dubey et al., 2024).

Additional Experiments. We compare the performance of decoder-only and encoder-only Transformers on a range of more realistic tasks. We test the sample complexity and length generalization capabilities of decoders and encoders using addition tasks (Lee et al., 2023). We also train both architectures to perform in-context learning (Garg et al., 2022) on various simple functions, such as linear functions and two-layer neural networks. Additionally, we train small decoder-only and encoder-only Transformers on a large text dataset (Gokaslan et al., 2019) to assess their performance in language modeling tasks.

079 080 081

066

067

068

069

071

072 073

074

075

076

077

# 2 RELATED WORK

082 083

**Expressive Power of Transformers.** There have been various literature exploring the expressive 084 power of Transformers. From the lens of universal approximation, Yun et al. (2020) showed that any 085 continuous sequence-to-sequence function over a compact set can be approximated arbitrary close by a Transformer (of finite albeit very large size). Other works approach expressiveness from the 087 perspective of computability and complexity such as Pérez et al. (2021) which showed Transformers 880 are Turing complete and Merrill et al. (2022); Merrill & Sabharwal (2024); Li et al. (2024) which 089 use circuit complexity to characterize the languages recognizable by Transformers of fixed depth. Giannou et al. (2023) presents a framework for Transformers as universal computers by placing 090 them in a loop. Communication complexity has also been used to show the impossibility of one-layer 091 Transformers of expressing certain functions e.g. induction head, without the model size being linear 092 in the length of the input (Sanford et al., 2024b;a). We modify functions introduced in (Sanford et al., 093 2024b) to compare the expressive power decoder-only Transformers against ENTP. We note that the 094 existing bounds (Sanford et al., 2024b;a) are highly related, but they do not directly imply the relative 095 expressive power of encoder and decoder (and specifically of the same fixed model size). 096

097

Transformer Architectures for Next Token Prediction. Transformers have become the de facto 098 backbones for next-token prediction tasks, leading to several variants such as encoder-decoder, causal decoder-only, and prefix decoder-only models. In the encoder-decoder model (Lewis et al., 100 2019; Chung et al., 2024), similar to the vanilla Transformer (Vaswani et al., 2017), the encoder 101 transforms the input tokens into conditioning features, and the decoder auto-regressively predicts 102 the target tokens by using cross-attention over the encoded representation and causal attention over 103 the output tokens. In contrast, the causal decoder-only model (Brown et al., 2020; Chowdhery et al., 104 2023) uses only the Transformer decoder and applies causal attention to all tokens to perform next-105 token prediction, ensuring that each token attends only to previous tokens. The prefix decoder-only model (Raffel et al., 2020; Wu et al., 2021) is similar to the causal decoder-only model but differs 106 in that it applies non-causal attention (i.e., full self-attention) to the input sequence (see Figure 8 for 107 visualizations of the attention patterns in these variants).

With the development of these models, recent studies have investigated the performance of each variant across various tasks. Notably, Wang et al. (2022) examined the zero-shot generalization performance of each model along with various objectives, and Ding et al. (2024) analyzed the performance of causal decoder-only and prefix decoder-only models in in-context learning. However, despite these diverse studies, there is a lack of research on encoder-only models that do not impose the constraint of causal attention for every next-token prediction. Therefore, in this work, we analyze the characteristics of encoder-only next token prediction (ENTP), comparing them with decoder-only models.

116 117

118

128 129

132

133

136

140

146

147

148

154

157

#### **3** PRELIMINARIES

**Sequence-to-Token Functions and Autoregression.** Given a vocabulary V (which we traditionally think of as some finite set, but could in general be any arbitrary perhaps uncountable set e.g.  $\mathbb{R}^d$ ), we can define a sequence over V as  $(x_1, \ldots, x_n)$  where  $x_1, \ldots, x_n \in V$ . Let  $V^* = \{(x_1, \ldots, x_n) : n \in \mathbb{N}; x_i \in V\}$  be the set of all sequences generated by V. Then, we say that  $f: V^* \to V$  is a sequence-to-token function.

We can view a causal model as a map from an input sequence to an output sequence with the causality constraint that the *i*'th ouput token depends only on the first *i* input tokens. Mathematically, we enforce this causality contraint by characterizing our causal model,  $\mathcal{T}_f : V^* \to V^*$ , with some sequence-to-token function *f* where on input sequence  $(x_1, \ldots, x_n)$  we have that

$$\mathcal{T}_f(x_1, x_2, \dots, x_n) \coloneqq (f(x_1), f(x_1, x_2), \dots, f(x_1, x_2, \dots, x_n)). \tag{1}$$

Observe that a sequence-to-token function f can be used auto-regressively to generate tokens from some initial sequence  $(x_1, \ldots, x_n)$  via the following update rule:

$$x_{n+i} \coloneqq f(x_1, \dots, x_{n+i-1}). \tag{2}$$

This can also be viewed as a special case of the causal model, where the input sequence is chosen so that

$$\mathcal{T}_f(x_1, x_2, \dots, x_n) = (x_2, x_3, \dots, x_{n+1}).$$
(3)

Hence, if we are trying to learn a causal or auto-regressive model, it suffices to learn the sequence function that generates it. Thus in this paper, we focus on the type of sequence functions that encoders versus decoders can learn and express.

**Encoders and Decoders.** We will use the letters  $\mathcal{E}$  and  $\mathcal{D}$  respectively to refer to encoders and decoders. In this paper, both models refer to variants of the Transformer architecture introduced in Vaswani et al. (2017), where the only difference lies in the masking used on the attention scores (decoder uses a causal mask while encoder allows full attention, as illustrated in Figure 1). The model size of a Transformer is determined by two parameters:

- *L*: number of Transformer blocks.
- *D*: embedding dimension.

As Transformers are sequence-to-sequence maps, we will use subscript notation where  $\mathcal{T}(x_1, \ldots, x_n)_i$  denotes the *i*'th value in the output sequence. We also allow our models the option to use positional embeddings. Tilde-notation i.e.  $\tilde{\mathcal{E}}$  and  $\tilde{\mathcal{D}}$  will denote models that do not use positional embeddings. For token embeddings  $x_1, \ldots, x_n$  and positional embeddings  $p_1, \ldots, p_n$ :

$$\mathcal{E}(x_1,\ldots,x_n) \coloneqq \mathcal{E}(x_1+p_1,\ldots,x_n+p_n).$$

In our experiments, we will use encoder and decoder models that have access to trainable positionalembeddings.

**Encoder and Decoders as Causal Models.** Given encoder  $\mathcal{E}$  and decoder  $\mathcal{D}$ , we can associate them with sequence-to-token functions as follows:

160 161  $f_{\mathcal{E}}: (x_1, \dots, x_n) \mapsto \mathcal{E}(x_1, \dots, x_n)_n$  $f_{\mathcal{D}}: (x_1, \dots, x_n) \mapsto \mathcal{D}(x_1, \dots, x_n)_n.$  We then have  $\mathcal{T}_{\mathcal{E}}$  and  $\mathcal{T}_{\mathcal{D}}$  as the causal models of  $\mathcal{E}$  and  $\mathcal{D}$  when used as sequence functions  $f_{\mathcal{E}}$ and  $f_{\mathcal{D}}$  respectively.<sup>1</sup> Under this characterization, we make two observations. Firstly, we can view  $\mathcal{T}_{\mathcal{E}}$  as an explicit and necessary way to introduce causality to the encoder  $\mathcal{E}$  since there is nothing implicit to the encoder that forces causality. Secondly, in juxtaposition to the previous statement, the causal model  $\mathcal{T}_{\mathcal{D}}$  is exactly equivalent to just using  $\mathcal{D}$ , that is  $\mathcal{T}_{\mathcal{D}} = \mathcal{D}$ . This is because  $\mathcal{D}$ enforces causality implicitly via the attention mask (see Appendix A.2 for formal proof). Therefore the explicit enforcement becomes redundant.

**EXPRESSIVE POWER OF ENCODER-ONLY VS. DECODER-ONLY** 

- 169
- 170

# 171

4

**TRANSFORMERS** 

172

Given that we can learn causal functions (defined as causal model in preliminary) using either encoders and decoders, the natural question to ask is how the expressive power of each model is related, i.e. can the encoder express more causal functions than a decoder of the *same* model size? Or perhaps they express the exact same class of causal functions? Towards answering this question, one trivial observation is that one-layer decoders and encoders are equivalent (formal proof in Appendix A.3) which directly implies the existence of causal functions<sup>2</sup> that both architecture can model exactly.

Now, what about causal functions that a decoder can model but not encoder, or vice versa — that
 encoder can model but not decoder? These functions exist too, as we show in the following two
 theorems:

**Theorem 1.** For any  $L \ge 2$  and  $D \ge 1$ , there exists a position-free decoder  $\widetilde{D}$  that has L-layers and embedding dimension D, such that for any encoder  $\mathcal{E}$ , there exists some input sequence  $(x_1, x_2, ...)$ with  $x_1, x_2, \dots \in \mathbb{R}^D$ , and  $\mathcal{T}_{\widetilde{D}}(x_1, x_2, \dots) \neq \mathcal{T}_{\mathcal{E}}(x_1, x_2, \dots)$ .

**Theorem 2.** For any  $L \ge 2$  and  $D \ge 1$ , there exists a position-free encoder  $\widetilde{\mathcal{E}}$  that has L-layers and embedding dimension D, such that for any decoder D with positional embeddings satisfying  $p_1 \ne p_2$ , there exists some input sequence  $(x_1, x_2, \ldots)$  with  $x_1, x_2, \cdots \in \mathbb{R}^D$ , and  $\mathcal{T}_{\widetilde{\mathcal{E}}}(x_1, x_2, \ldots) \ne \mathcal{T}_D(x_1, x_2, \ldots)$ .

The above two theorems are existential in nature. Informally, Theorem 1 says that if we consider causal model defined over the entirety of  $\mathbb{R}^D$  as its vocabulary, we can find some decoder, for which any encoder will differ from it on some input sequence. Theorem 2 makes a similar (albeit weaker statement) in the other direction; namely the existence of a causal function computable by an encoder, but not by any decoder that uses "non-trivial" positional embeddings (e.g. embeddings for different positions are unique). Detailed proof of both theorems are deferred to Appendix A.

197 Of course, the setting and assumptions of the above two statements are not necessarily very realistic. For one, they focus on general class of causal models rather than only auto-regressive ones. Furthermore, the results only pertain to exact realization and say nothing about approximation. The 199 assumption of unbounded domain is also not realistic as in practice decoders are trained and used 200 over a finite domain of tokens, each with some fixed embeddings. And specific to Theorem 2, no 201 claim is made about decoders that do not use positional embeddings. But despite the limitations, 202 these theorems give an indication that the expressive power of encoder and decoder model are dif-203 ferent — despite the almost identical description modulo the attention mask. Changing the mask on 204 the attention scores causes significant changes to the properties of the model. Thus, in the following 205 sections we propose an auto-regressive tasks and run experiments comparing encoders and decoders 206 that corroborates this view. 207

208 209

210

211

212

215

# 5 TIME AND SPACE COMPLEXITY COMPARISONS

Inspired by the different computational models of encoder-only and decoder-only Transformers, we characterize the causal sequence functions learnable by encoders and decoders based on their

<sup>&</sup>lt;sup>1</sup>To be fully consistent with notation in equation 1, we should denote them as  $\mathcal{T}_{f_{\mathcal{E}}}$  and  $\mathcal{T}_{f_{\mathcal{D}}}$  respectively. However we abuse notation and use  $\mathcal{T}_{\mathcal{E}}$  and  $\mathcal{T}_{\mathcal{D}}$  for sake of simplicity.

<sup>&</sup>lt;sup>2</sup>Here causal function refers to sequence-to-sequence function where the outputs are only determined by current and previous inputs, i.e.  $\mathcal{F}(x_1, x_2, ..., x_n)_i$  only depends on  $x_1, x_2, ..., x_i$ , and not  $x_j$  for any j > i.

required computational complexity. We give an informal comparison of encoders and decoders in terms of their required time and space complexities — both over the entire sequence and for each additional token. We propose Count3, which is closely related to Match3 (Sanford et al., 2024b), to highlight the "gap" between the complexity of encoders and decoders. Count3 is feasible for an encoder but challenging for a decoder due to its limited computation complexity.

222 **Time Complexity Comparison.** Decoder-only Transformers, using KV-Cache, take O(n) time to generate each token, and  $O(n^2)$  time to generate an entire sequence. Because an ENTP has to 223 compute the entire attention matrix for every token, it takes  $O(n^2)$  time to generate each token. Thus, 224 it takes  $O(n^3)$  time to generate the entire sequence. While this implies that ENTP is more compute-225 intensive (i.e., ENTP will be slower than decoder-only Transformer), this also implies that ENTP 226 can express more compute-intensive functions than decoders. Specifically, since the total amount of 227 compute that decoders use for generating n tokens is  $O(n^2)$ , they cannot run any algorithm whose 228 runtime is  $\omega(n^2)$  (strictly greater than quadratic time). 229

**Space Complexity Comparison.** Both encoder-only and decoder-only use O(n) space complexity to generate an entire sequence. Although the standard implementation of attention uses  $O(n^2)$ space, attention can be implemented using only O(n) space. For details of algorithmic implementations of attention using O(D) memory, refer to Algorithm 3 in the appendix. Thus, we need a more detailed approach to find a difference between the space complexity of the two models.

Towards this end, we classify memory used for the computation over the current token as either precomputed or additional. Precomputed memory stores values from computation over past tokens. Values stored in precomputed memory *persist*, and are used for computation over current and future tokens, e.g. the keys and values of previous tokens for a decoder. Additional memory stores values that depend on the current token, e.g. keys and values of the current token.

When generating the *n*th token, a decoder uses O(n) precomputed memory to store keys and values of previous tokens and O(1) additional memory to compute results over the current token. An encoder computes everything from scratch for each token, so it uses O(n) additional memory and no precomputed memory. Under this view, there is a space complexity gap between encoder and decoder.

Table 1: Complexity for next token inference.

Complexity	Encoder-only	Decoder-only
Additional Time Complexity	$O(n^2 DL)$	O(nDL)
Precomputed Space Complexity	N/A	O(nDL)
Additional Space Complexity	O(nD)	O(D)

250 251 252

221

Most of our complexity analysis focuses on Transformers with fixed sizes, so we primarily consider 253 complexity with respect to the sequence length n. However, we also account for the embedding 254 dimension D and the number of layers L in Table 1. Both encoder-only and decoder-only Trans-255 formers use O(LD) time because the attention operation is performed O(L) times, and computing 256 each query, key, and value vector is O(D). In the case of multi-head attention, we assume D = hd, 257 where h is the number of heads and d is the dimension of the query, key, and value vectors. A 258 decoder uses O(nhdL) = O(nDL) precomputed space because it stores nhL query, key, and value 259 d-dimensional vectors. Both encoders and decoders use O(D) additional space for current token's 260 embedding vector — and we specifically note that there is no dependence on L as Transformer do 261 computation sequentially on the layer (i.e. all the additional computation required for layer  $\ell$  is done 262 before all the additional computation required for layer  $\ell + 1$ ). Thus, we can do the computation over L layers using O(D) space by overwriting computation over previous layers. 263

5.1 Count3 DEFINITION AND COMPLEXITY

Consider the sequence function that maps an input sequence of positive integers  $x_1, x_2, \ldots, x_n$  to the number of pairs  $x_i, x_j$  where the modulo-*n* sum of  $x_i, x_j$  and  $x_n$  is equal to 0. More formally,

268 269

264 265

266

Count3
$$(x_1, x_2, \dots, x_n) \coloneqq |\{(i, j) \in [n]^2 : x_i + x_j + x_n \equiv 0 \pmod{n}\}| \pmod{n}.$$
 (4)

281

282

283

284

286

287

289

291 292

293

294

295

306

307 308

270

271

272

273

274

Algorithm 1 Algorithm to compute Count3 in  $O(n^2)$  time and O(1) space

**Require:** length n sequence of integers  $(x_1, \ldots, x_n)$ 1: count  $\leftarrow 0$ 2: for  $i = 1, \ldots, n$  do 3: for  $j = 1, \ldots, n$  do 4: if  $(x_i + x_j + x_n) \equiv 0 \pmod{n}$  then 5: count  $\leftarrow$  count + 16: return count (mod n)

Algorithm 2 Algorithm to compute Count3 in O(n) time and O(n) space

**Require:** length *n* sequence of integers  $(x_1, \ldots, x_n)$ 1: count  $\leftarrow 0$ 2: table  $\leftarrow$  zero-indexed length-*n* array of 0's 3: for  $i = 1, \ldots, n$  do 4:  $k \leftarrow -x_i \mod n$ 5: table $[k] \leftarrow$  table[k] + 16: for  $i = 1, \ldots, n$  do 7:  $k \leftarrow (x_i + x_n) \mod n$ 8: count  $\leftarrow$  count + table[k]9: return count (mod *n*)

Count3 is an augmented version of Match3 (Sanford et al., 2024b). As shown in (Sanford et al., 2024b), the triple-wise relationships, used in both Match3 and Count3, are difficult for Transformers to represent, because of the pairwise nature of self-attention.

Note that there exists algorithms that can compute Count3 on some length-n input sequence 296  $x_1, x_2, \ldots, x_n$  in either  $O(n^2)$  time and O(1) space or in O(n) time and O(n) space. See Al-297 gorithm 1 and Algorithm 2 for exact pseudocode implementations. In brief, Algorithm 1 iterates 298 through all  $n^2$  pairs checking if they meet the modulo-n sum requirement. Algorithm 2 uses the 299 fact that  $x_i + x_j + x_n \equiv 0 \pmod{n}$  is equivalent to  $x_i + x_n \pmod{n} \equiv -x_j \pmod{n}$ . In two 300 linear passes, it counts each  $-x_i \pmod{n}$  and stores each count in a table, then sums the values 301 in the table for each  $x_i + x_n \pmod{n}$ . Now, given these two algorithms, we make the following 302 conjecture: 303

**Conjecture 1.** Given an algorithm A that computes  $\text{Count}3(x_1, x_2, \ldots, x_n)$ , at least one of the following must hold true:

(i) A requires  $\Omega(n^2)$  time with access to  $x_n$ 

(ii) A requires  $\Omega(n)$  space storing values unique to n.

This conjecture seems plausible given that both Algorithm 1 and Algorithm 2, which we consider to be optimal, adhere to it. Algorithm 1 uses  $O(n^2)$  time after accessing  $x_n$  (line 4 of Algorithm 1) and Algorithm 2 requires O(n) memory, where the stored values are a function of n (line 4 of Algorithm 2). Given this conjecture, we can show the following lemma:

**Lemma 1.** Given that Conjecture 1 holds and assuming  $O(\log n)$  precision, any decoderonly Transformer with fixed embedding dimension D satisfying  $\mathcal{D}(x_1,\ldots,x_m)_m =$ Count $3(x_1,\ldots,x_m)$  for all sequences of length  $m \leq n$  must have  $L = \Omega(n)$ .<sup>3</sup>

318 319 *Proof.* Let  $\mathcal{D}$  be a decoder with L layers and embedding dimension D satisfying  $\mathcal{D}(x_1, \ldots, x_m)_m =$ 320 Count $3(x_1, \ldots, x_m)$  for all sequences of length  $m \le n$ . We can use  $\mathcal{D}$  as an algorithm to compute 321 Count $3(x_1, \ldots, x_n)$ , by outputting  $\mathcal{D}(x_1, \ldots, x_n)_n$  on input  $(x_1, \ldots, x_n)$ . Thus either condition 322 (i) or (ii) of Conjecture 1 must hold when  $\mathcal{D}$  computes the output sequence over  $(x_1, \ldots, x_n)$ .

323

 $<sup>{}^{3}</sup>O(\log n)$  precision is not required for Lemma 1, but it is included to be consistent with Lemma 2.

327

328

330

331 332

333 334

335

336

337



Figure 2: An example of a sequence used in a Count3 experiment.

**Case 1:** (i) is true. In this case  $\mathcal{D}$  requires  $\Omega(n^2)$  time with access to  $x_n$ , i.e.  $\mathcal{D}$  uses  $\Omega(n^2)$  time for the last token. From Table 1, we know that  $\mathcal{D}$  uses O(nLD) time for each token. Thus condition (i) is true only if we have that  $nLD = \Omega(n^2)$ . Since D is fixed, it follows that  $L = \Omega(n)$ .

338 **Case 2:** (ii) is true. In this case,  $\mathcal{D}$  requires  $\Omega(n)$  space storing values unique to n. Because 339 decoders are causal, we have  $\mathcal{D}(x_1,\ldots,x_n)_i = \text{Count3}(x_1,\ldots,x_i)$  for all  $i \in [n]$ . Then since 340 we assume (ii) is true, computing  $\mathcal{D}(x_1,\ldots,x_n)_i$  requires  $\Omega(i)$  space for each  $i \in [n]$ . Fur-341 thermore by the uniqueness assumption of (ii), for  $i \neq j$ , the values stored when computing Count3 $(x_1, \ldots, x_i)$  are different from the values stored when computing Count3 $(x_1, \ldots, x_i)$ . 342 Since decoders are causal, the space used to compute  $\mathcal{D}(x_1,\ldots,x_n)_i$  cannot be overwritten when 343 computing  $\mathcal{D}(x_1,\ldots,x_n)_i$ , for j > i. Hence, when  $\mathcal{D}$  computes  $\operatorname{Count3}(x_1,\ldots,x_n)$ , it uses  $\Omega(i)$ 344 space to compute Count3( $x_1, \ldots, x_i$ ), for each  $i \in [n]$ . Thus,  $\mathcal{D}$  uses  $\Omega\left(\sum_{i \in [n]} i\right) = \Omega(n^2)$  space 345 346 to compute  $Count3(x_1, ..., x_n)$ . From Table 1, we know that  $\mathcal{D}$  uses O(nLD) space for the entire 347 sequence. Then  $nLD = \Omega(n^2)$ . Since D is fixed, it follows that  $L = \Omega(n)$ . 348

Finally, as  $L = \Omega(n)$  is a necessary condition for both conditions (i) and (ii), Lemma 1 follows.  $\Box$ 

**Lemma 2.** Assuming  $O(\log n)$  precision, there exists an encoder  $\mathcal{E}$  with L = O(1) and D = O(1) such that  $\mathcal{E}(x_1, \ldots, x_m)_m = \text{Count3}(x_1, \ldots, x_m)$  for all sequences of length  $m \le n$ .<sup>4</sup>

353 Proof of Lemma 2 is in Appendix A.6.

**Remark 1.** With linear chain-of-thought (generating O(n) tokens before answering), a decoder would be able to perform Count3. We provide a RASP<sup>5</sup> (Weiss et al., 2021) program, Algorithm 6, to demonstrate this.

350

351

352

354

355

5.2 SMALL-SCALE EXPERIMENT

We train small decoder-only and encoder-only Transformers on auto-regressive sequences generated from Equation (4). To generate unique sequences, we start each sequence with a seed containing 16 random integers between 0 and 63. Then we extend the sequence to 64 integers using Equation (4) (see Figure 2). Seeds are generated randomly during training. The seed portion of the sequence was not used to compute the loss during training, so the model was only trained on the deterministic part of the sequence. As shown in Figure 3, the decoder-only Transformers demonstrate some ability to learn patterns related to the distribution of numbers in Count3 sequences, but they completely fail to learn the task. In contrast, the encoder successfully learns the task with near-perfect accuracy.

368 369

376

377

5.3 Count3 with Large Language Model

Furthermore, we investigate the performance of decoder-only large language models (LLMs) on the Count3 task. We fine-tune Llama-3 (Dubey et al., 2024) and GPT-40 (Achiam et al., 2023) using sequences of 64 integers, including 16 seeds, as introduced in Section 5.2. To enable the LLMs to leverage their knowledge, we also include the code for Algorithm 1 in the prompt, asking the models to provide the result after executing the code (see Table 4 for the full prompt). As shown in Figure 4,

<sup>&</sup>lt;sup>4</sup>We assume  $x_1, ..., x_m < m$ .

<sup>&</sup>lt;sup>5</sup>RASP (Weiss et al., 2021) is a programming language that describes Transformer computations, by mapping attention and feed-forward computation into simple primitives.



Figure 3: **Training loss and sequence accuracy curve for the** Count3. The ENTP successfully learns to perform the Count3 task, but the decoder-only models struggle to learn it.

the LLMs struggle with the Count3 task, which is consistent with our small-scale experiment. It demonstrates that the suggested characteristics of causal decoder-only models hold true even at large scales. We provide the validation of the prompt design used, as well as details about the LLM fine-tuning, in the Appendix C.1.



Figure 4: Results of LLM fine-tuning on Count3.

Count3 sequences have very small Kolmogorov complexity i.e., the size of the smallest program
that can generate this data, as Algorithm 1 (which uses only six lines of pseudocode) serves as an
upper bound. However, as shown by Lemma 1 and the experiments in Section 5.3, a decoder, no
matter how large, cannot efficiently learn the task.

We fine-tune BERT (Devlin et al., 2019) as an ENTP model, and it successfully learns Count3 faster
 than small-scale ENTP models. See Figure 11 for details.

#### 5.4 SIMILAR FUNCTION LEARNABLE BY DECODER

Motivated by the question of how we need to change Count3 so that it can be learned by a decoder,
 we examine a modified version of Match3 (Sanford et al., 2024b). <sup>6</sup>

$$Match3'(x_1, x_2, \dots, x_n) \coloneqq \begin{cases} 1 & \exists \ (i, j) : x_1 + x_i + x_j = 0 \pmod{128} \\ 0 & \text{otherwise} \end{cases}$$
(5)

There are several key differences between Count3 and Match3':

- (i) Match3' uses a fixed modulus, whereas Count3 employs the sequence length as the modulus. This simplifies the decoder's task, as the modulus remains constant across all tokens, enabling reuse of intermediate values from previous tokens.
- (ii) Match3' operates on triplets  $(x_1, x_i, x_j)$  rather than  $(x_i, x_j, x_n)$ . By using  $x_1$  instead of  $x_n$ , it becomes easier for the decoder since  $x_1$  remains unchanged for different tokens, facilitating the reuse of intermediate values across tokens.
- (iii) Match3' checks for the existence of a condition rather than counting occurrences. Counting is challenging to implement within the attention mechanism without scaling values by the

<sup>6</sup>The main difference between Match3 and Match3' is that Match3' operates on triplets  $(x_1, x_i, x_j)$ .

sequence length. Due to the causal mask, scaling value vectors by sequence length is not straightforward.

We provide RASP (Weiss et al., 2021) program, Algorithm 5 that satisfies  $\mathcal{D}(x_1, x_2, ..., x_n)_n =$ Match3' $(x_1, x_2, ..., x_n)$  for sequences of any length, assuming  $O(\log n)$  precision. We train small Transformers to verify that both decoders and encoders can perform Match3', and find that both models can perform Match3' with high accuracy.

Table 2	: Match3'	performance.
I GOIC D	1110000110	periornance

Model	Min Loss	Individual Token Accuracy	Full Sequence Accuracy
Medium Decoder (6 layer)	0.0001	<b>99.99</b> %	99.92%
Medium Encoder (6 layer)	0.0016	99.97%	99.50%

#### ADDITIONAL EXPERIMENTS

#### 6.1 GENERALIZATION ON ADDITION

We test the sample complexity and length generalization capabilities of decoders and encoders using addition tasks. We use the reversed addition format (\$123+456=975\$) from Lee et al. (2023). We find that encoders exhibit lower sample complexity compared to decoders, as seen in Figure 5, meaning they require fewer training examples to achieve similar performance. Additionally, encoders demonstrate superior ability to generalize to longer sequences, as shown in Figure 6. We provide more experimental details and results in Appendix C.3.





Decoder

ENTP

Figure 5: Addition Sample Complexity. The Figure 6: Addition Length Generalization. train and test datasets include numbers with up to 3 digits.

The train dataset includes numbers with up to 10 digits, and the test dataset has numbers with up to 15 digits.

#### 6.2 IN-CONTEXT LEARNING

We consider the problem of learning a function class  $\mathcal{F}$  using in-context example (Garg et al., 2022). In this problem, a function  $f \in \mathcal{F}$  is sampled from a distribution  $\mathcal{D}_{\mathcal{F}}$ , and a sequence of random in-puts is sampled i.i.d. from  $\mathcal{D}_{\mathcal{X}}$ , forming a prompt  $P: (x_1, f(x_1), x_2, f(x_2), \dots, x_N, f(x_N))$ . The objective is for the model to in-context learn the function f from the prompt P and predict  $f(x_{query})$ for a new input  $x_{\text{query}}$ . To this end, we train the model  $M_{\theta}$ , parameterized by  $\theta$ , by minimizing the expected loss over random sampled prompts as follows:

where  $\ell(\cdot, \cdot)$  is loss function and  $P^i$  refers to the partial prompt containing *i* in-context examples with  $(i+1)^{th}$  input (i.e.,  $(x_1, f(x_1), \ldots, x_i, f(x_i), x_{i+1})$ ). 

In this work, we examine four types of function classes: linear function, sparse linear function, two-layer neural network and decision tree. For all function classes, we sample  $x_i$  from a Gaussian

distribution  $N(0, \mathbf{I}_d)$  where d represents the dimension of  $x_i$ , and utilize the squared error as loss function. Detailed descriptions for each function class are provided in Appendix C.2. 

We provide the in-context-learning results according to the number of in-context examples for each function class in Figure 7. The encoder-only models demonstrate better performance compared to the decoder-only models in linear regression and sparse linear regression, while exhibiting compet-itive performance in two-layer NN regression and decision tree. 



Figure 7: **Results of in-context learning experiment.** The encoder-only models demonstrate superior or competitive performance across all function classes compared to the decoder-only models.

## 6.3 **OPENWEBTEXT**

We train Transformers on the OpenWebText dataset (Gokaslan et al., 2019), an open-source replication of the WebText dataset used to train GPT-2 (Radford et al., 2019), with the next-token prediction objective. We use medium models, described in Table 5, and hyperparameters from Table 6. As shown in Table 3, encoder-only Transformer slightly outperforms decoder-only Transformer.

Table 3: Minimum values of training and validation loss, as well as perplexity, for decoder-only and encoder-only Transformers on the OpenWebText dataset.

Model	Train Loss	Validation Loss	Train Perplexity	Validation Perplexity
Decoder-only	4.694	4.705	109.3	110.5
Encoder-only	4.643	4.650	103.9	104.6

#### **CONCLUSION AND FUTURE WORK**

In this work, we present theoretical and novel experimental results suggesting that, assuming compute is unlimited, decoder-only Transformers are not the ideal model for sequence modeling. We show that ENTP is more expressive without compromising generalization. Using Theorem 1 and 2, we find that the classes of functions encoder-only and decoder-only Transformers can exactly learn are different. We introduce the Count3 task and demonstrate, both theoretically (assuming Conjec-ture 1) and experimentally, that while encoders can perform this task easily, decoders cannot. We also find that encoders outperform decoders on a variety of auto-regressive tasks, including length generalization and in-context learning. However, the expressivity gained by removing the causal mask comes at the expense of efficiency.

**Future Work.** Several open avenues for further research remain in ENTP theory, such as proving Lemma 1 without relying on Conjecture 1 (or providing a counterexample) and generalizing Count3 to a broader class of tasks where encoders perform well but decoders do not. Further investigation into length generalization across tasks beyond addition is valuable to confirm the findings of this paper. Another natural extension of our work would be to explore whether we can further enhance expressivity past ENTP. In a more practical direction, as both training and inference with ENTP are compute-intensive, future work may focus on reducing computational costs. For instance, could we first train a decoder and then efficiently convert it to ENTP? Is it possible to combine decoders with ENTP? More broadly, are there better methods for efficiently training ENTP models? 

# 540 REFERENCES

565

566

567

571

572

573

574

581

582

583

Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical
report. *arXiv preprint arXiv:2303.08774*, 2023.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (eds.), *Advances in Neural Information Processing Systems*, volume 33, pp. 1877–1901. Curran Associates, Inc., 2020.

- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. Palm: Scaling language modeling with pathways. *Journal of Machine Learning Research*, 24(240): 1–113, 2023.
- Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li,
  Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. Scaling instruction-finetuned language models. *Journal of Machine Learning Research*, 25(70):1–53, 2024.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding, 2019. URL https://arxiv.org/abs/1810.04805.
  - Nan Ding, Tomer Levinboim, Jialin Wu, Sebastian Goodman, and Radu Soricut. CausalLM is not optimal for in-context learning. In *The Twelfth International Conference on Learning Representations*, 2024.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models.
   *arXiv preprint arXiv:2407.21783*, 2024.
  - Shivam Garg, Dimitris Tsipras, Percy Liang, and Gregory Valiant. What can transformers learn in-context? a case study of simple function classes. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho (eds.), *Advances in Neural Information Processing Systems*, 2022.
- Angeliki Giannou, Shashank Rajput, Jy-Yong Sohn, Kangwook Lee, Jason D. Lee, and Dimitris Papailiopoulos. Looped transformers as programmable computers. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett (eds.), Proceedings of the 40th International Conference on Machine Learning, volume 202 of Proceedings of Machine Learning Research, pp. 11398–11442. PMLR, 23–29 Jul 2023. URL https://proceedings.mlr.press/v202/giannou23a.html.
  - Aaron Gokaslan, Vanya Cohen, Ellie Pavlick, and Stefanie Tellex. Openwebtext corpus. http: //Skylion007.github.io/OpenWebTextCorpus, 2019.
- Edward J Hu, yelong shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. LoRA: Low-rank adaptation of large language models. In *International Conference on Learning Representations*, 2022.
- 587
   588
   588
   589
   590
   590
   591
   592
   593
   594
   594
   595
   595
   596
   596
   596
   597
   598
   598
   599
   598
   599
   599
   590
   590
   590
   590
   590
   590
   590
   590
   590
   590
   590
   590
   590
   590
   590
   590
   590
   590
   590
   590
   590
   590
   590
   590
   590
   590
   590
   590
   590
   590
   590
   590
   590
   590
   590
   590
   590
   590
   590
   590
   590
   590
   590
   590
   590
   590
   590
   590
   590
   590
   590
   590
   590
   590
   590
   590
   590
   590
   590
   590
   590
   590
   590
   590
   590
   590
   590
   590
   590
   590
   590
   590
   590
   590
   590
   590
   590
   590
   590
   590
   590
   590
   590
   590
   590
   590
   590
   590
   590
   590
   590
   590
   590
   590
   590
   590
   590
- Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdel rahman Mohamed, Omer
   Levy, Veselin Stoyanov, and Luke Zettlemoyer. Bart: Denoising sequence-to-sequence pre training for natural language generation, translation, and comprehension. In *Annual Meeting* of the Association for Computational Linguistics, 2019.

- 594 Zhiyuan Li, Hong Liu, Denny Zhou, and Tengyu Ma. Chain of thought empowers transformers to 595 solve inherently serial problems. arXiv preprint arXiv:2402.12875, 2024. 596
- William Merrill and Ashish Sabharwal. A logic for expressing log-precision transformers. Advances 597 in Neural Information Processing Systems, 36, 2024. 598
- William Merrill, Ashish Sabharwal, and Noah A Smith. Saturated transformers are constant-depth 600 threshold circuits. Transactions of the Association for Computational Linguistics, 10:843–856, 2022.
- Jorge Pérez, Pablo Barceló, and Javier Marinkovic. Attention is turing-complete. Journal of Ma-603 *chine Learning Research*, 22(75):1–35, 2021. 604
- 605 Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 606 Language models are unsupervised multitask learners. 2019. URL https://api. 607 semanticscholar.org/CorpusID:160025533.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi 609 Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text 610 transformer. Journal of machine learning research, 21(140):1–67, 2020. 611
- 612 Clayton Sanford, Daniel Hsu, and Matus Telgarsky. One-layer transformers fail to solve the induc-613 tion heads task. arXiv preprint arXiv:2408.14332, 2024a.
- 614 Clayton Sanford, Daniel J Hsu, and Matus Telgarsky. Representational strengths and limitations of 615 transformers. Advances in Neural Information Processing Systems, 36, 2024b. 616
- 617 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, 618 Ł ukasz Kaiser, and Illia Polosukhin. Attention is all you need. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (eds.), Advances in Neural 619 Information Processing Systems, volume 30. Curran Associates, Inc., 2017. 620
- 621 Thomas Wang, Adam Roberts, Daniel Hesslow, Teven Le Scao, Hyung Won Chung, Iz Beltagy, 622 Julien Launay, and Colin Raffel. What language model architecture and pretraining objective 623 works best for zero-shot generalization? In International Conference on Machine Learning, pp. 624 22964-22984. PMLR, 2022. 625
  - Gail Weiss, Yoav Goldberg, and Eran Yahav. Thinking like transformers. In International Conference on Machine Learning, pp. 11080–11090. PMLR, 2021.
- 628 Shaohua Wu, Xudong Zhao, Tong Yu, Rongguo Zhang, Chong Shen, Hongli Liu, Feng Li, Hong 629 Zhu, Jiangang Luo, Liang Xu, et al. Yuan 1.0: Large-scale pre-trained language model in zero-630 shot and few-shot learning. arXiv preprint arXiv:2110.04725, 2021.
- Chulhee Yun, Srinadh Bhojanapalli, Ankit Singh Rawat, Sashank Reddi, and Sanjiv Kumar. Are 632 transformers universal approximators of sequence-to-sequence functions? In International Con-633 ference on Learning Representations, 2020. 634
- 635 Hattie Zhou, Arwen Bradley, Etai Littwin, Noam Razin, Omid Saremi, Joshua M. Susskind, Samy Bengio, and Preetum Nakkiran. What algorithms can transformers learn? a study in length 636 generalization. In The Twelfth International Conference on Learning Representations, 2024. 637
- 638 639

627

631

601

602

608

- 640
- 641 642
- 643
- 644
- 645

646

#### 648 NOTATION AND DEFERRED PROOFS А 649

#### A.1 NOTATION

650 651 •  $\mathcal{E}, \mathcal{D}$ : Encoder/Decoder with positional embeddings. 652 •  $\widetilde{\mathcal{E}}, \widetilde{\mathcal{D}}$ : Encoder/Decoder without positional embeddings (with same model weights as  $\mathcal{E}, \mathcal{D}$ ) 653 654 (which we call position free). 655 • L: number of Transformer blocks. 656 • D: embedding dimension. 657 •  $p_i$ : The *i*'th positional embedding where  $p_i \in \mathbb{R}^D$ . 658 •  $\mathbf{W}_{Q/K/V}^{i}$ : The weight matrix of the Query, Key and Value respectively of the *i*'th attention 659 block. 661 • I: The identity matrix. 662 • 0: The zero matrix. 663 •  $(x_1, x_2, x_3, \dots)$ : A sequence of values/tokens. •  $\mathcal{T}_{\mathcal{E}}(x_1, x_2, x_3, \dots) \coloneqq (\mathcal{E}(x_1)_1, \mathcal{E}(x_1, x_2)_2, \mathcal{E}(x_1, x_2, x_3)_3, \dots).$ 665 666 For convention, we will use bold capital letters  $\mathbf{X}$  to denote matrices, and unbold lowercase letters 667 x to denote vectors. 668 669 For theorem 1 and theorem 2, our model of Transformers will use single-head attention and assume 670 the dimension of the Query and Key are equal to the embedding dimension D. We also omit layer 671 normalization and scaling of attention scores by  $1/\sqrt{D}$ . 672 673 A.2  $T_{\mathcal{D}} = \mathcal{D}$ 674 *Proof.* The main observation is that the MLP layers of a decoder are element-wise, and the 675 attention layers of a decoder are causal (i.e. the contextual embedding of the *i*'th token is computed 676 only using the tokens  $j \leq i$ ). We thus have that 677

$$\mathcal{D}(x_1,\ldots,x_i,\ldots,x_n)_i = \mathcal{D}(x_1,\ldots,x_i)_i$$

for any  $i \in [n]$ . As a result,

$$\mathcal{D}(x_1, \dots, x_i, \dots, x_n)_i = \mathcal{D}(x_1, \dots, x_i)_i$$
  
=  $f_{\mathcal{D}}(x_1, \dots, x_i)$   
=  $\mathcal{T}_{\mathcal{D}}(x_1, \dots, x_i, \dots, x_n)_i$ ,

for all  $i \in [n]$ , i.e.  $\mathcal{D}(x_1, \ldots, x_n) = \mathcal{T}_{\mathcal{D}}(x_1, \ldots, x_n)$ .

#### A.3 ONE-LAYER ENCODER AND DECODER ARE EQUIVALENT

*Proof.* Let  $\mathcal{D}$  and  $\mathcal{E}$  have the same parameters. Then the query key, and value vectors for the attention layer (denoted  $q_i, k_i, v_i$  for each  $x_i$  respectively) will be the same for both models.

For one-layer decoder: 
$$\mathcal{D}(x_1, ..., x_n)_i = \text{MLP}\left(\sum_{j=1}^i \text{Softmax}(q_i, k_j)v_j\right)$$
.  
For one-layer encoder:  $\mathcal{E}(x_1, ..., x_n)_i = \text{MLP}\left(\sum_{j=1}^n \text{Softmax}(q_i, k_j)v_j\right)$ .  
For one-layer encoder:  $\mathcal{E}(x_1, ..., x_n)_i = \text{MLP}\left(\sum_{j=1}^n \text{Softmax}(q_i, k_j)v_j\right)$ .  
If  $i = n$ , then  $\mathcal{D}(x_1, ..., x_n)_i = \mathcal{E}(x_1, ..., x_n)_i$ . Thus,  $\mathcal{T}_{\mathcal{D}}(x_1, x_2, ..., x_n) = \mathcal{T}_{\mathcal{E}}(x_1, x_2, ..., x_n)$ .

697

678

679

684

685 686

687 688

689

690

A.4 PROOF OF THEOREM 1

699 *Proof.* We first provide a construction for  $\mathcal{D}$ . For the attention-block of the first two layers, we use 700 the same weight matrices. Namely, we set  $\mathbf{W}_{K}^{1} = \mathbf{W}_{K}^{2} = \mathbf{W}_{Q}^{1} = \mathbf{W}_{Q}^{2} = \mathbf{0}$  and  $\mathbf{W}_{V}^{1} = \mathbf{W}_{V}^{2} = \mathbf{I}$ . 701 For every other attention block, we set them to the constant zero function by setting  $\mathbf{W}_{V}^{i} = \mathbf{0}$  for

 $i \geq 3$ . We similarly set the weights and biases of every MLP block to zero. Thus, in essence  $\widetilde{\mathcal{D}}$  is just two duplicate attention blocks stacked on top of each other with a skip connection after each attention block.

Now consider three arbitrary vectors  $x_1, x_2, x_3 \in \mathbb{R}^D$  and its corresponding sequence  $(x_1, x_2, x_3)$ . Let us first compute the output of  $\tilde{\mathcal{D}}$  on  $(x_1, x_2, x_3)$ . The first attention block and skip connection will map the input sequence to the sequence

$$\left(2x_1, x_2 + \frac{x_1 + x_2}{2}, x_3 + \frac{x_1 + x_2 + x_3}{3}\right)$$

712 The second attention block and skip connection will then map to the following:

$$\left(4x_1, \frac{7x_1+9x_2}{4}, \frac{23x_1+17x_2+32x_3}{18}\right)$$

716 Our first observation from this mapping is that there clearly exists  $\tilde{x}_1, \tilde{x}_2, \tilde{x}_3 \in \mathbb{R}^D$  such that 717  $\widetilde{\mathcal{D}}(\tilde{x}_1, \tilde{x}_2, \tilde{x}_3)_3 \neq \widetilde{\mathcal{D}}(\tilde{x}_2, \tilde{x}_1, \tilde{x}_3)_3$ . Our second observation is that

$$4x_1 = \frac{7x_1 + 9x_2}{4} \quad \Longleftrightarrow \quad x_1 = x_2$$

which follows from simplifying the left hand equation.

Now, let us for sake of contradiction assume the existence of some encoder  $\mathcal{E}$  such that  $\mathcal{T}_{\mathcal{E}}$  exactly replicates  $\widetilde{\mathcal{D}}$  on every input sequence. We first claim that the first two positional embeddings  $p_1, p_2$ of  $\mathcal{E}$  must differ. This follows by our first observation and thus the requirement that  $\mathcal{E}(\tilde{x}_1, \tilde{x}_2, \tilde{x}_3)_3 \neq$  $\mathcal{E}(\tilde{x}_2, \tilde{x}_1, \tilde{x}_3)_3$  — which can only happen if  $p_1 \neq p_2$  due to the permutation invariance of encoders when there are no positional embeddings. Now as  $p_1 \neq p_2$ , there exists vectors  $y_1, y_2, c \in \mathbb{R}^D$  such that  $y_1 \neq y_2$  and  $y_1 + p_1 = y_2 + p_2 = c$ . It follows immediately that

$$\mathcal{E}(y_1)_1 = \mathcal{E}(y_1 + p_1)_1$$

$$=\mathcal{E}(c)_1$$

$$=\widetilde{\mathcal{E}}(c,c)_2$$

$$=\widetilde{\mathcal{E}}(y_1+p_1,y_2+p_2)_2$$
$$=\mathcal{E}(y_1-y_2)_2$$

$$= \mathcal{E}(y_1, y_2)_2.$$

But since  $y_1 \neq y_2$ , by the second observation we made, it must be that  $\widetilde{\mathcal{D}}(y_1)_1 \neq \widetilde{\mathcal{D}}(y_1, y_2)_2$ . Since we assumed that  $\mathcal{T}_{\mathcal{E}}$  exactly replicates  $\widetilde{\mathcal{D}}$  on every input sequence, it thus follows that  $\mathcal{E}(y_1)_1 \neq \mathcal{E}(y_1, y_2)_2$ — a contradiction. Hence, no such encoder  $\mathcal{E}$  exists, which directly implies that we can always find some sequence  $(x_1, x_2, ...)$  where  $\widetilde{\mathcal{D}}(x_1, x_2, ...) \neq \mathcal{T}_{\mathcal{E}}(x_1, x_2, ...)$ .

742 A.5 PROOF OF THEOREM 2

743 744 745 746 746 746 747 747 747 748 Proof. We first provide a construction for  $\tilde{\mathcal{E}}$ . For the attention-block of the first two layers, we use the same weight matrices. Namely, we set  $\mathbf{W}_{K}^{1} = \mathbf{W}_{Q}^{2} = \mathbf{W}_{Q}^{1} = \mathbf{W}_{Q}^{2} = \mathbf{W}_{V}^{1} = \mathbf{W}_{V}^{2} = \mathbf{I}$ . For every other attention block, we set them to the constant zero function by setting  $\mathbf{W}_{V}^{i} = \mathbf{0}$  for  $i \ge 3$ . We similarly set the weights and biases of every MLP block to zero. Thus, in essence  $\tilde{\mathcal{E}}$  is just two duplicate attention blocks stacked on top of each other with a skip connection after each attention block.

Now consider two arbitrary vectors  $x_1, x_2 \in \mathbb{R}^D$  and its corresponding sequence  $(x_1, x_2)$ . A brief inspection will reveal that

$$\mathcal{T}_{\widetilde{\mathcal{E}}}(x_1, x_2) = (4x_1, \alpha x_1 + \beta x_2), \tag{7}$$

753 where  $\alpha, \beta > 0$ . 754

Next, we assume the existence of some  $\mathcal{D}$  where  $p_1 \neq p_2$  and exactly replicates  $\mathcal{\tilde{E}}$ . We fix  $x_1 = 0$ and let  $x_2 = p_1 - p_2$ . Observe that  $x_2 \neq 0$  as  $p_1 \neq p_2$ . It follows that there is some constant vector

752

709 710 711

713 714 715

719

720

734

735 736

737

738

739

 $c \in \mathbb{R}^D$  where  $c \in \mathbb{R}^D$ 

758

759

760 761 762

763

764

765

Now from equation 7, we have that  $\mathcal{T}_{\tilde{\mathcal{E}}}(x_1, x_2) = (0, \beta x_2)$  for some  $\beta > 0$ . As  $x_2 \neq 0$ , it follows that  $\beta x_2 \neq 0$  and hence  $(0, \beta x_2)$  is not a constant sequence — contradicting the output of the decoder. Hence, no  $\mathcal{D}$  where  $p_1 \neq p_2$  can exactly replicates  $\tilde{\mathcal{E}}$ .

= (c, c).

 $\mathcal{D}(x_1, x_2) = \widetilde{\mathcal{D}}(x_1 + p_1, x_2 + p_2)$ 

 $= \widetilde{\mathcal{D}}(p_1, p_1)$ 

## A.6 PROOF OF LEMMA 2

*Proof.* Lemma 2 follows from Algorithm 4. From Weiss et al. (2021), a RASP program can be compiled to a Transformer, with a fixed number of Transformer-layers and attention heads. However, it assumes that the MLPs can perform any element-wise operation. Thus, it suffices to show that each MLP needs O(1) layers with respect to sequence length n.

We first observe that the largest internal value possible within Algorithm 4 is  $n^2$ , so there are  $n^2 + 1$ distinct internal values including 0. Using  $\lfloor 2 \log_2 n \rfloor + 1$  bits, we can represent all of these values as unsigned integers. Using floating-point numbers we would also need  $\Theta(\log n)$  bits.

All linear element-wise operations in Algorithm 4 can be implemented trivially, since an MLP can perform arbitrary linear transformations. Therefore, we focus on the only nonlinear element-wise function  $g: [2n-1] \times [n] \rightarrow [n-1]$ :<sup>7</sup>

778

779 780 781

782 783 Using a constant number of linear operations and ReLU functions, g can be constructed as follows:

 $g(a,b) = \begin{cases} a, & a < b\\ a-b, & a \ge b. \end{cases}$ 

$$g(a,b) = \operatorname{ReLU}(a - M \operatorname{ReLU}(a - b + \epsilon)) + \operatorname{ReLU}(a - b),$$

where  $0 < \epsilon < 1$  and  $M \ge \frac{2n}{\epsilon}$ .

<sup>785</sup> Because g can be implemented using a constant number of linear operations and ReLU functions, it can be implemented by an MLP with ReLU activation functions and O(1) depth.

Because of skip-connections, we can concatenate MLPs by zeroing out attention. Then we can create any MLP of O(1) depth.

Thus, it is possible to construct an encoder  $\mathcal{E}$  with L = O(1) and D = O(1) such that  $\mathcal{E}(x_1, \ldots, x_m)_m = f_{\text{TC}}(x_1, \ldots, x_m)$  for all sequences of length  $m \le n$ .

792 793 794

796

797 798

799 800

801

802

803

804 805

806 807

808

809

# **B** ATTENTION PATTERNS OF DIFFERENT TRANSFORMER ARCHITECTURES

In Figure 8, we provide the visualization of attention patterns of encoder-only, decoder-only, prefix decoder-only, and encoder-only models.

# C EXPERIMENT DETAILS

In all experiments, we train encoders in the same manner as decoders, processing entire sequences in each batch with a single gradient optimization step. Although this approach does not offer the same efficiency benefits for encoders as it does for decoders, we adopt it to maintain consistency between the training processes of both models.

C.1 Count3 WITH LLM

In the main paper, we fine-tuned two LLMs for Count3: Llama3-8B (Dubey et al., 2024) and GPT-40 (Achiam et al., 2023). Here, we provide the details about the fine-tuning of each model. For

<sup>&</sup>lt;sup>7</sup>This function implements modular division for a bounded range.  $g(a, b) = a \mod b$ , when  $0 \le a < 2b$ .



Figure 8: Attention patterns of different Transformer architectures in next token prediction. Encoder-decoder and prefix decoder-only models first perform full attention on the input (prefix), and then use causal attention to predict subsequent tokens. In contrast, decoder-only models apply causal attention to all tokens, without distinguishing between input and output. Encoder-only models also do not separate input and output, but they recalculate attention from scratch for each token prediction, performing full attention on all tokens.

GPT-40, we used the official API, setting the batch size to 4 and the learning rate multiplier to 10. For Llama3-8B, we employed LoRA fine-tuning Hu et al. (2022) with a batch size of 16 and a learning rate of  $1.4 \times 10^{-4}$ . Regarding prompt design, we included the algorithm code in the prompt so that the LLMs could leverage its knowledge of natural language (see Table 4). We note that the loss was applied only to the answer part of the prompt.

Additionally, we verify the validity of the prompt design used for LLM fine-tuning. To this end, we modify the task from Count3 to Count2<sup>8</sup> and fine-tune the Llama3-8B using prompts that include algorithmic code, as in the main experiment. As shown in Figure 9, the model successfully learns Count2 with the proposed prompt design, achieving high sequence accuracy. This demonstrates that the model's difficulty in learning Count3 is due to the characteristics of causal decoder-only architecture, rather than an issue with the training strategy, such as the prompt design.



Figure 9: **Results of LLM fine-tuning on** Count2. The Llama3-8B successfully learns Count2 when using the same prompt format as Count3. This demonstrates that the reason LLMs struggle with Count3 is not due to the complexity of the prompt, but rather the characteristics of the decoder-only model.

#### C.2 IN-CONTEXT LEARNING

As introduced in the main paper, four types of function classes are considered in the in-context learning experiment (Garg et al., 2022): linear function, sparse linear function, two-layer neural network, and decision tree. For all function classes, the input  $x_i$  is drawn from a Gaussian distribution

<sup>8</sup>Count2(
$$x_1, x_2, \dots, x_n$$
) :=  $\left| \left\{ i \in n : x_i + x_n \equiv 0 \pmod{n} \right\} \right| \pmod{n}$ 

Table 4: **Example of prompt used for LLM experiments.** We include the code for the algorithm in the prompt to leverage the knowledge of the LLMs.

 $N(0, \mathbf{I}_d)$  where d represents the dimension of  $x_i$ . We provide detailed descriptions of each function class below.

**Linear function.** We consider the class of linear functions  $\mathcal{F} = \{f \mid f(x) = w^{\top}x, w \in \mathbb{R}^d\}$ where w is drawn from a Gaussian distribution  $N(0, \mathbf{I}_d)$ . Following previous work (Garg et al., 2022), we set d = 20 and the number of data points N to 40.

**Sparse linear function.** We also consider a sparse linear function, which is similar to a linear function setup. The difference is that after drawing w from  $N(0, \mathbf{I}_d)$ , only k randomly selected coordinates are kept, while the remaining ones are set to zero. Following previous work (Garg et al., 2022), we set k = 3.

**Two-layer neural network** We examine the class of two-layer ReLU neural networks  $\mathcal{F} = \{f \mid f(x) = \mathbf{W}_{(2)}\sigma(\mathbf{W}_{(1)}x), \mathbf{W}_{(2)} \in \mathbb{R}^{1 \times h}, \mathbf{W}_{(1)} \in \mathbb{R}^{h \times d}\}$  where  $\sigma(\cdot) = \max(0, \cdot)$  (i.e., ReLU function). We set h = 100, d = 20, and the number of data point N to 100.

**Decision tree** We consider the class of decision trees represented by full-binary trees of fixed height. In these trees, the leaf node values are drawn from N(0, 1), and the non-leaf nodes are sampled from random integers between 0 and d, indicating an index of the input x. At each non-leaf node, if the value of the input at the specified index is positive, we move to the right; otherwise, we move to the left. Given an input, we start at the root node and repeat this process until reaching a leaf node. The value of the leaf node becomes the output of the function.

911 C.3 ADDITION

We test the sample complexity of encoder-only and decoder-only Transformers using addition tasks,
with up to 3-digit numbers. We sample the dataset of all possible 3-digit addition examples using
a method similar to the method described in Lee et al. (2023). We start with all 1,000,000 3-digit,
2-digit, and 1-digit addition examples. Then we randomly remove 90% of the 3-digit addition
examples, adjusting the ratio of 3-digit to 2-digit examples from around 100:1 to around 10:1. Next
we split the data into training, testing, and validation splits, stratified by the number of digits and

carries in each addition example. All 1-digit addition examples were put into the training split. Since
our models tend to over fit the training dataset, we save the model with the lowest loss on a validation
dataset. We test decoder-only and encoder-only Transformers on plain and reversed addition tasks,
using between 1.25k to 20k training examples. All sample complexity tests are run with at least 5
different seeds. We test small models, described in Table 5.



Figure 10: Addition Sample Complexity. The train and test datasets include numbers with up to 3 digits. The dataset uses the plain addition format (\$123+456=579\$), unlike the results in Figure 5.

We train Transformers to add larger numbers and evaluate their ability to perform length generalization. Training is performed on numbers with up to 10 digits, while testing extends to numbers
with up to 15 digits. Each model is trained on 100,000 examples using the reversed addition format
(Lee et al., 2023). The numbers are sampled to ensure equal probability for each length, without
duplicates. Consequently, in larger datasets, there are fewer 1-digit addition examples compared to
10-digit ones, as the total number of possible 1-digit examples is smaller. All length generalization
tests are run with 3 different seeds. We test Small-Deep models described in Table 5.

## C.4 MODEL SIZES

In Table 5, we provide the configurations of the Transformer architectures used in the experiments from the main paper.

Name	Number of Layers	Number of Heads	Embedding Dimension
Small	3	3	192
Medium	6	6	384
Large	12	12	768
Small-Deep	8	2	128

Table 6: OpenWebText Hyperparameters

Table 5: Model specifications.

## C.5 OPENWEBTEXT

Parameter	Value
warmup_iters	2000
lr_decay_iters	600,000
min_lr	0.00006
max_lr	0.0006
beta1	0.9
beta2	0.95
weight_decay	0.1
block_size	128
batch_size	32

#### 972 D ADDITIONAL EXPERIMENTAL RESULTS ON Count3

In this section, we conducted Count3 experiments with additional Transformer variants. First, we
evaluated the performance of Prefix decoder-only models (Raffel et al., 2020; Wu et al., 2021),
which perform non-causal attention for the prefix portion. Using the same experimental setup, we
initiated sequences with 16 random seeds, which were used as the prefix. As shown in Figure 11,
while the Prefix decoder-only model slightly outperforms the decoder-only model, it also fails to
learn the triplet counting task. This demonstrates that performing full attention to certain parts of a
sequence in a decoder-only model is insufficient to solve tasks requiring Count3-level complexity.

Next, we conducted experiments on the Count3 task using BERT (Devlin et al., 2019), a representative encoder-only architecture. Specifically, we fine-tuned BERT using the ENTP approach under the same experimental conditions. As shown in Figure 11, BERT combined with ENTP successfully learned triplet counting. Notably, as BERT is pre-trained and larger compared to the medium transformer used in the paper, it converged more quickly. This result indicates that ENTP is effective not only for the model configurations specified in the main paper but also for larger pre-trained models.



Figure 11: The performance of the Prefix decoder-only model and BERT combined with ENTP on the Count3 task.

## **E** IMPLEMENTATION OF ATTENTION USING O(D) MEMORY

Algorithm 3 Implementation of Attention Using O(D) Memory **Require:**  $q \in \mathbb{R}^{n \times d}, k \in \mathbb{R}^{n \times d}, v \in \mathbb{R}^{n \times D}$ 1:  $y_n \leftarrow \mathbf{0}^D$ 2:  $a \leftarrow \mathbf{0}^D$ 3:  $b \leftarrow 0$ 4: for j = 1, ..., n do 5:  $c \leftarrow \exp(q_n^T k_j)$ 6:  $a \leftarrow a + cv_i$ 7:  $b \leftarrow b + c$ 8:  $y_n \leftarrow \frac{a}{b}$ 9: return  $y_n$ 

# 1026 F RASP ALGORITHMS

<sup>1028</sup> In Algorithm 4, 5, and 6, we provide the Python RASP implementation (Zhou et al., 2024) for Count3 and Match3'.

1030 1031 1032

1052 1053 1054

1055

```
Algorithm 4 Count3 RASP Encoder Implementation
```

```
def g(a, b):
1033
               return a if a < b else a - b
1034
1035
          def count_triplets(x):
               idxs = indices(x)
1036
1037
               # set n[i] = len(x) and last_x[i] = x[-1] for all i (only possible with encoder)
               n = sel_width(select(k=x, q=x, pred=true))
last_x = kqv(k=idxs, q=n - 1, v=x, pred=equals, reduction="mean")
1038
1039
               # g(a, b) is equivalent a % b if 0 <= a < 2 \star b
               y = seq_map(n - x, n, g) # y[i] = -x[i] % n
z = seq_map(x + last_x, n, g) # z[i] = (x[i] + x[-1]) % n
1040
1041
               # conut the number of (i, j) such that y[i] == z[j]
# c = sum(A), where A[i, j] = 1 if y[i] == z[j] else 0
1042
1043
               c = kqv(
                    k=full(x, 1),
1044
                    q=full(x, 1),
1045
                    v=sel_width(select(k=z, q=y, pred=equals)) * n, # sum(v) = mean(v * n)
                    pred=equals,
1046
                    reduction="mean",
1047
               )
1048
               # conpute count % n
1049
               c -= idxs * n
               # because count <= n^2, there exists i such that c[i] = count % n \text{ or } c[i] = n
# the case c[i] = n is handled by the default value (0) when no keys are selected
1050
               return kqv(k=c, q=n, v=c, pred=lambda a, b: 0 <= a and a < b, reduction="mean")
1051
```

#### Algorithm 5 Match3' RASP Decoder Implementation

```
def has_triplet(x):
1056
              idxs = indices(x)
              first_x = kqv(k=idxs, q=full(x, 0), v=x, pred=equals, reduction="mean", causal=True)
1057
              # use bitmask to compute mod
y = -x & 127  # y[i] = -x[i] % 128
z = (first_x + x) & 127  # z[i] = (x[0] + x[i]) % 128
1058
              У
1059
1060
              \# \max_{count}[-1] > 0 if there exists (i, j) such that y[i] == z[j]
1061
              max_count = kqv(
                   k=full(x, 1),
q=full(x, 1),
1062
1063
                   v=sel_width(select(k=y, q=z, pred=equals)),
                   pred=equals,
reduction="max",
1064
              )
1065
1066
              return tok_map(max_count, lambda a: min(a, 1)) # return 0 or 1
1067
1068
```

1072 1073 1074

1069 1070 1071

1075

1076

1077

1078

Algo	rithm 6 Count3 RASP Decoder COT Implementation
def (	count_triplets(x):
:	idxs = indices(x) n = kgy(k=x, g=full(x, EOS), y=idxs, pred=eguals, reduction="min", causal=True
1	n = tok_map(n, lambda a: a if a else -2)
	<pre>seq_len = kqv(k=x, q=x, v=idxs, pred=true, reduction="mean") seq_len = kqv(k=x, q=x, v=idxs, pred=true, reduction="max", causal=True)</pre>
	i = seq_len - n
	j = seq_len - 2 * n ki = kgv(k=idxs, g=i, v=x, pred=eguals, reduction="max", causal=True)
1	<pre>xj = kqv(k=idxs, q=j, v=x, pred=equals, reduction="max", causal=True)</pre>
	y = (n - xi) + 1
:	$z = (last_x + xj) % n + 1$
-	y_mask_write = (n <= idxs) & (idxs < 2 ★ n) z mask_write = (2 ★ n <= idxs) & (idxs < 3 ★ n)
-	$\sum_{n=1}^{n} (1 - 1) = (1$
:	z_mask_read = (z * n < 10xs) & (10xs <= 3 * n)
:	z_count = sel_width( select(
	k=x * y_mask_read,
	q=2, pred=lambda a, b: a == b and a != 0,
	causal=True,
(	count = kqv (
	q=z_mask_read,
	v=n * x * z_mask_read, pred=lambda a, b: a & b,
	reduction="mean",
ė	ans = count ≈ n
ė	ans_mask_write = idxs == 3 * n eos_mask_write = idxs > 3 * n
	return (
	y * y_mask_write
	+ z_count * z_mask_write + ans * ans_mask_write
	+ EOS * eos_mask_write