SOFTMAX TRANSFORMERS ARE TURING-COMPLETE

Anonymous authors

Paper under double-blind review

ABSTRACT

Hard attention Chain-of-Thought (CoT) transformers are known to be Turing-complete. However, it is an open problem whether softmax attention Chain-of-Thought (CoT) transformers are Turing-complete. In this paper, we prove a stronger result that length-generalizable softmax CoT transformers are Turing-complete.

More precisely, our Turing-completeness proof goes via the CoT extension of the Counting RASP (C-RASP), which correspond to softmax CoT transformers that admit length generalization. We prove Turing-completeness for CoT C-RASP with causal masking over a unary alphabet (more generally, for the letter-bounded languages). While we show that this is actually not Turing-complete for arbitrary languages, we prove that its extension with relative positional encoding is Turing-complete for arbitrary languages. We empirically validate our theoretical results by training transformers for various languages that require complex (non-linear) arithmetic reasoning.

1 Introduction

Transformers (Vaswani et al., 2017) have enabled powerful Large Language Models (LLMs) with Chain-of-Thought (CoT) steps, which are capable of complex reasoning (cf. (Wei et al., 2022; OpenAI et al., 2024)). But what task can (and cannot) be done by CoT transformers? This fundamental question lies at the heart of the recent effort in understanding the ability of transformers through the lens of formal language theory (see the survey Strobl et al. (2024)). In particular, the question whether CoT transformers is *Turing-complete* — that is, capable of solving any problems solvable by Turing machines — is especially pertinent; see the work (cf. (Pérez et al., 2021; Bhattamishra et al., 2020; Merrill & Sabharwal, 2024; Qiu et al., 2025; Li & Wang, 2025)).

Are CoT transformers Turing-complete? All existing proofs of Turing-completeness of CoT transformers (cf. (Pérez et al., 2021; Bhattamishra et al., 2020; Merrill & Sabharwal, 2024; Qiu et al., 2025; Li & Wang, 2025)) employ *hardmax attention*, which is a rather unrealistic assumption. In particular, its use comes at the cost of a lack of a trainability guarantee. It is still an open question to date whether CoT transformers that use softmax attention are Turing-complete, and whether one can guarantee some sort of trainability.

Contributions. The main contributions of this paper are (i) to prove for the first time that softmax CoT transformers are Turing-complete, and (ii) to provide a kind of guarantee of trainability.

More precisely, we use the framework from Huang et al. (2025) of length-generalizable softmax transformers. In particular, the authors showed that a simple declarative language called C-RASP (with causal masking) Yang & Chiang (2024) can be converted into their framework, thereby also admitting length generalization. Our results use the extensions of these models with CoT steps.

We first show that CoT C-RASPs with causal masking are Turing-complete over a unary alphabet $\Sigma = \{a\}$. More generally, we show that Turing-completeness holds for *letter-bounded languages*, i.e., $L \subseteq a_1^* \cdots a_n^*$, where a_1, \ldots, a_n are distinct letters in the alphabet. Such languages are especially interesting because of their ability to model complex number-theoretic concepts (e.g., prime numbers, exponentiation, multiplication, etc.).

Interestingly, we can show that CoT C-RASPs with causal masking are *not* Turing-complete over arbitrary languages. In fact, simple languages (e.g. palindromes) cannot be solved by CoT C-RASPs.

To address this limitation, we extend CoT C-RASPs with *Relative Positional Encodings (RPEs)* (cf. Shaw et al. (2018); Liutkus et al. (2021); Dufter et al. (2022)), which assigns a positional information to any token *relative* to another token. We extend the framework of Huang et al. (2025) by adding RPEs, and show that length-generalizability still holds. Finally, we show Turing-completeness for CoT C-RASP[RPEs] over arbitrary languages.

We provide an experimental validation of our results for CoT C-RASP and CoT C-RASP[RPEs] by showing length generalization of transformers for complex number-theoretic concepts with *unary* representation (to be captured by CoT C-RASP) and with *binary* representation (to be captured by CoT C-RASP[RPEs]). For example, the concept of prime numbers will be represented as the language $L = \{a^p : p \text{ is prime}\}$ with unary representation, and as the language $L' = \{\text{bin}(p) : p \text{ is prime}\}$ with binary representation (where bin(p) denotes the binary representation of the number p, e.g., 5 is written as 101).

Novelty in our proofs. Instead of simulating Turing machines (as in most of the previous proofs), our Turing-completeness proofs provide a simulation of Minsky's Turing-powerful counter machines Minsky (1961) by CoT softmax transformers.

Organization. We start with the CoT models in Section 2. We then prove Turing-completeness results for the unary and letter-bounded cases in Section 3. Turing-completeness for the general case is proven in Section 4. We report our experiments in Section 5. Finally, we conclude in Section 6.

2 Models for Trainable Cot Transformers

2.1 Transformers and C-RASP

Softmax Transformers. We assume transformer decoders with softmax attention and causal masking (Softmax Attention Transformers, SMAT). Our formal definition of softmax transformers follows that of Huang et al. (2025). Attention weights are defined as

$$\bar{w} = \operatorname{softmax}(\log n \cdot \{\mathbf{v}_j^T \mathbf{K}^T \mathbf{Q} \mathbf{v}_i\}_{j=1}^i)$$
(1)

where \mathbf{v}_i denotes activations at position i, and \mathbf{K} , \mathbf{Q} transform these to keys and queries, respectively. Here, scaling with $\log n$ is included, as it is needed to theoretically represent sparse functions across unboundedly input strings and circumvent theoretical limitations of soft attention (Chiang & Cholak, 2022; Edelman et al., 2022). For the feedforward networks, we assume one-layer networks, where each hidden unit has either ReLU or Heaviside activation. Here, as in Huang et al. (2025), Heaviside is needed to theoretically represent functions with sharp thresholds; at any finite input length, it can be arbitrarily closely approximated using ReLU MLPs. As is standard, we encode an input $x \in \Sigma^*$ by applying a token embedding function $\mathrm{em}: \Sigma \to \mathbb{R}^k$ for some dimension k.

To define the computation of CoT via SMAT, we need the transformer to be able to output a token. We further define an output function $o: \mathbb{R}^d \to \Sigma$, parameterized by applying a linear function $\mathbb{R}^d \to \mathbb{R}^{|\Sigma|}$ followed by an argmax selecting the symbol receiving the highest score. Overall, we view an SMAT as a length-preserving map $T: \Sigma^* \to \Sigma^*$, where $T(x)_i$ indicates the symbol predicted after reading the prefix $x_1 \dots x_i$.

We refer to Appendix A for a formal definition and further discussion of design choices.

C-RASPs. C-RASP is equivalent to the fragment $K_t[\#]$ Yang & Chiang (2024); Yang et al. (2024) of LTL[Count] Barceló et al. (2024) with only past operator:

$$\varphi ::= Q_a (a \in \Sigma) | \varphi \wedge \varphi | \neg \varphi | \varphi \vee \varphi | t \sim t (t \in \{<, =, >\})$$

$$t ::= c (c \in \mathbb{N}) | \#[\varphi] | t + t$$

We will now define the semantics of C-RASP by structural induction on the C-RASP expressions. Suppose $w=w_1\cdots w_n\in \Sigma^+$. [As a side remark, it is possible to also allow the empty string ϵ as input, and for this we can use the "start-of-string" symbol \vdash . We do not do this to avoid clutter.] For syntactic category φ , we will define $[\![\varphi]\!]_w$ as a bitstring $h_1\cdots h_n\in\{0,1\}^n$. On the other hand, for syntactic category t, we will define $[\![t]\!]_w$ as a sequence $m_1\cdots m_n\in\mathbb{Z}^n$ of integers. For each sequence σ , we will write $\sigma(i)$ to denote the ith element in the sequence. We start with the two base cases:

• $\varphi = Q_a$. In this case, $h_i \in \{0, 1\}$ is 1 iff $w_i = a$.

• t = c. In this case, $m_i = c$ for each i.

We now proceed to the inductive cases:

- $\varphi = \psi_1 \wedge \psi_2$. Then, $h_i = \min\{ [\![\psi_1]\!]_w(i), [\![\psi_2]\!]_w(i) \}$.
- $\varphi = \psi_1 \vee \psi_2$. Then, $h_i = \max\{ \llbracket \psi_1 \rrbracket_w(i), \llbracket \psi_2 \rrbracket_w(i) \}$.
- $\varphi = \neg \psi$. Then, $h_i = 1 |\![\psi]\!]_w(i)$.
- $\varphi = t \sim t'$. Then, $h_i = 1$ iff $\llbracket t \rrbracket_w(i) \sim \llbracket t' \rrbracket_w(i)$.
- $t = \begin{aligned} &\neq \\ \hline \#[\varphi]$. Let <math>m_0 = 0$. Then, for each i > 0, $m_i = m_{i-1} + 1$ if $[\![\varphi]\!]_w(i) = 1$; else $m_i = m_{i-1}$.

Relative Positional Encodings. We also define an extension C-RASP[RPEs] (resp. SMAT[RPEs]) of C-RASP (resp. SMAT) with Relative Positional Encodings (RPEs), which are simply subsets $\mathfrak{R}\subseteq\mathbb{N}\times\mathbb{N}$. We start with C-RASP[RPEs]. In the sequel, the notation $[\![\mathfrak{R}]\!]$ refers to the function mapping each $(i,j)\in\mathbb{N}\times\mathbb{N}$ to $\{0,1\}$ such that $[\![\mathfrak{R}]\!](i,j)=1$ iff $(i,j)\in\mathfrak{R}$. For the syntactic category t, we allow the counting term:

$$\overleftarrow{\#}_{\mathfrak{R}}[\varphi]$$

which is to be interpreted at position j as the cardinality of:

$$\{i \in [1,j] : (i,j) \in \Re, i \models \varphi\}.$$

That is, we include i depending on the positional encoding of each i relative to j. [Alternatively, \Re can be construed as allowing positions at certain distances from each j.] This is a generalization of the class C-RASP[periodic, local] defined by Huang et al. (2025), where \Re is either periodic or local

As for SMAT[RPEs], the definition is a simple modification of SMAT in that the formula in Equation 1 becomes

$$\bar{w} = \operatorname{softmax}(\log n \cdot \{\langle A\mathbf{v}_i, B\mathbf{v}_j \rangle + \lambda \| \mathfrak{R} \| (i,j) \}_{i=1}^i). \tag{2}$$

Here, we interpret λ as a bias term and $[\mathfrak{R}](i,j)$ as 1 if $(i,j) \in [\mathfrak{R}]$; otherwise, it is 0.

Discussion of Relative Positional Encodings Relative positional encodings, which modify attention scores with positional information, are a popular approach for providing positional information to transformers. Our formalization of RPEs is a simple formal abstraction of *additive relative positional encodings*, which add a position-dependent term to the attention logits (Shaw et al., 2018; Dai et al., 2019; Xue et al., 2021; Press et al., 2022; He et al., 2021). Schemes in the literature differ in whether they are parameter-free (e.g., Press et al. (2022)) or involve learnable parameters. We consider the especially simple case where R is determined a-priori, parameter-free, and independent of the task at hand.

2.2 Extensions with Chain-of-Thought

Suppose Γ is the (finite) set of possible CoT tokens. CoT tokens in some $\Gamma_F \subseteq \Gamma$ are reserved to indicate that the computation is to terminate and that the input string is to be "accepted". Let $\Gamma_{\neg F} = \Gamma \backslash \Gamma_F$.

We define a CoT to be a map $F: \Sigma^* \to \Gamma^* \cup \Gamma^\omega$, where Γ is a finite set of CoT tokens, where all non-final symbols are in $\Gamma_{\neg F} \subseteq \Gamma$. Here, note that we include both finite (terminating) CoTs in Γ^* and infinite (non-terminating) CoTs in Γ^ω . Consideration of non-terminating CoTs is needed for Turing completeness.

The language L(F) recognized by F is the set of all $w \in \Sigma^*$ such that F(w) is finite and ends in an element of $\Gamma_F \subseteq \Gamma$.

CoT C-RASPs. To extend C-RASP (resp. C-RASP[RPEs]) with CoTs, a simple solution is as follows. A CoT C-RASP expression (over Γ) is a non-empty sequence $S=d_1,\ldots,d_l$ of definitions d_i of the form:

$$O_{a_i} \leftarrow \varphi_{a_i}$$
,

where $a_i \in \Gamma_{\neg F}$ and φ_{a_i} a normal C-RASP (resp. C-RASP[RPEs]) expression. The intuition of S is a *switch* condition, which will tell the program which token to output. S outputs a token on an input string $w \in (\Sigma \cup S)^+$ if $[\![\varphi_{a_i}]\!]_w(|w|) = 1$ for some i. The output of S on a string $w \in (\Sigma \cup \Gamma_{\neg F})^+$ is defined to be a_i , where i is the smallest index such that $[\![\varphi_{a_i}]\!]_w(|w|) = 1$ and that $[\![\varphi_{a_j}]\!]_w(|w|) = 0$ for each j < i. In this case, we write $S(w) = a_i$. Note that a CoT transformer might terminate without outputting a token if $[\![\varphi_{a_j}]\!]_w(|w|) = 0$ for each j; in this case, the input string w will be immediately rejected. Here, we write $S(w) = \bot$ (i.e. undefined).

A CoT C-RASP S generates the string $U=U_1\cdots U_m\in \Gamma^*$ on the input $w\in \Sigma^*$ if $S(wU_1\cdots U_{k-1})=U_k$ for each $k=1,\ldots,m$. Intuitively, this means that S autoregressively outputs the symbols in U. The language L(T) accepted by a CoT C-RASP S is defined to be the set of all $w\in \Sigma^*$ such that there exists a finite string $U\in \Gamma^*$ ending in an element of Γ_F such that T generates U on w, and non-last symbols in U are in $\Gamma_{\neg F}$.

We remark that, in many cases, the order of the sequence S is not so important, especially if we can ensure that at most O_{a_i} is going to be satisfied. We will use this in the sequel.

CoT SMATs. Recall that we view an SMAT T as a length-preserving map $T: \Sigma^* \to \Sigma^*$, where $T(x)_i$ indicates the symbol predicted after reading the prefix $x_1 \dots x_i$. An SMAT $T: (\Sigma \cup \Gamma)^* \to (\Sigma \cup \Gamma)^*$ generates the string $U = U_1 \cdots U_m \in \Gamma^*$ on the input w if T autoregressively predicts the string U – that is, if $T(wU_1 \cdots U_{k-1}) = U_k$ for each $k = 1, \dots, m$. The language L(T) accepted by a CoT SMAT T is defined to be the set of all $w \in \Sigma^*$ such that there exists a finite string $U \in \Gamma^*$ ending in Γ_F such that T generates T0 on T0, and non-last symbols in T1 are in T2.

Proposition 2.1. If a language is accepted by a CoT C-RASP (resp. C-RASP[RPEs]), then it is also accepted by a CoT SMAT (resp. SMAT[RPEs]).

The proof is a simple extension of (Huang et al., 2025, Thm. 9); see Appendix A.2.

2.3 Learnability with CoT

We now show that CoT C-RASP is learnable in the framework of Huang et al. (2025). Intuitively, this framework considers the case where transformers are trained on data from some bounded length and then deployed on data of larger lengths. We now make this formal. As before, we view SMATs as defining length-preserving maps $T: \Sigma^* \to \Sigma^*$. Define the hypothesis class Θ as the set of SMATs T where each parameter vector and matrix of T is represented at p bits of precision, for some p depending on T.

Definition 2.2. A language L is length-generalizably learnable with CoT if there is a CoT F with L(F) = L such that the following holds:

For each $i=1,2,3,\ldots$, use the idealized learning procedure from Definition 6 in Huang et al. (2025) to choose a sequence of SMATs $T_i\in\Theta$ $(i=1,2,3,\ldots)$ such that each T_i generates $F(w)_{1...i-|w|}$ on all inputs w, $|w|\leqslant i$. Then, there is some N_0 depending on L such that for all $i>N_0$, T_i will exactly recognize the language L with Cot.

We analogously define the same notions in the presence of RPEs. Given a set $\mathfrak{R} \subseteq \mathbb{N} \times \mathbb{N}$, define the hypothesis class $\Theta[\mathfrak{R}]$ as the set of SMAT[RPEs] T with the RPE \mathfrak{R} , where each parameter vector and matrix of T is represented at p bits of precision, for some p depending on T, and where each λ in (2) is fixed to 1. We then define *length-generalizably learnable with CoT with RPE* \mathfrak{R} by replacing Θ with $\Theta_{\mathfrak{R}}$ in Definition 2.2.

The intuition of this definition is that we can learn a single SMAT that works for all input lengths, even when training only on data from some bounded length, as long as the training length is sufficiently large. We note that the definition of the learning setup is substantially simpler than in Huang et al. (2025) since our transformers use no absolute positional encodings. Whereas Huang

¹Such a sequence always exists, as there is just a finite number of inputs at each length i.

et al. (2025) used separate hypothesis classes Θ_n at each context window size n, our learning setup requires a single hypothesis class Θ that works for all input lengths.

We then obtain the following guarantee:

Proposition 2.3. Consider a language expressible in C-RASP[RPEs] CoT, using RPE \Re . Then it is length-generalizably learnable with RPE \Re .

The proof is a straightforward adaptation of (Huang et al., 2025, Thm. 7); see Appendix A.2.

3 UNARY CASE

Theorem 3.1. Each recursively enumerable language over a unary alphabet $\Sigma = \{a\}$ can be recognized by SMAT in the CoT setting.

This theorem follows from the following proposition.

Proposition 3.2. Each recursively enumerable language over a unary alphabet $\Sigma = \{a\}$ can be recognized by C-RASP in the CoT setting.

In turn, this follows directly from the following proposition; recall that a language $L \subseteq \Sigma^+$ is letter-bounded if it is a subset of $a_1^*a_2^*\cdots a_n^*$ for some distinct letters $a_1,\ldots,a_n\in\Sigma$.

Proposition 3.3. Each recursively enumerable letter-bounded language over any alphabet Σ can be recognized by C-RASP in the CoT setting.

We will deduce Proposition 3.3 from the following proposition, which will be most convenient for our construction. Given an alphabet Σ with $\Sigma = \{a_1, \ldots, a_n\}$, the corresponding $Parikh\ map$ is the map $\Psi \colon \Sigma^* \to \mathbb{N}^n$, where $w \in \Sigma^*$ is mapped to $(|w|_{a_1}, \ldots, |w|_{a_n})$, where $|w|_{a_i}$ is the number of occurrences of a_i in w. In other words, $\Psi(w)$ is the vector that contains all letter counts in w. Notice that for $u, v \in \Sigma^*$, we have $\Psi(u) = \Psi(v)$ if and only if v can be obtained from v by re-arranging the letters, or by v we have v we say that a language v is v is v in v in

Proposition 3.4. Each recursively enumerable permutation-invariant language over any alphabet Σ can be recognized by C-RASP in the CoT setting.

To prove this proposition, we give a reduction from counter machines. To define these, we define Φ_k to the set of expressions φ of the following form: a conjunction of counter tests of the form $x_i \sim 0$, where x_i indicates the *i*th counter and $\infty \in \{>, =\}$. A *k-counter machine (k-CM)* is a tuple

$$(P, \Delta, q_0, F)$$

where P is a set of states,

$$\Delta \subseteq P \times \Phi_k \times P \times \{-1, 0, 1\}^k$$

is a set of *transitions*, $q_0 \in P$ is the *initial state*, and $F \subseteq P$ is the set of final states. We also assume that the machine is *deterministic*, i.e., whenever there are two transitions $(p, \varphi, q, \boldsymbol{u})$ and $(p, \varphi', q', \boldsymbol{u}')$ starting in the same state p, but with $(q, \boldsymbol{u}) \neq (q', \boldsymbol{u}')$, then φ and φ' cannot be true at the same time (i.e. $\varphi \wedge \varphi'$ is unsatisfiable). For a transition $\tau = (p, \varphi, q, \boldsymbol{u})$, we will employ the notation $\mathrm{src}(\tau) := p$, $\mathrm{tgt}(\tau) := q$, $\varphi_\tau := \varphi$, and $\boldsymbol{u}_\tau := \boldsymbol{u}$.

A configuration of such a k-CM is a tuple $(q, x) \in P \times \mathbb{Z}^k$, where $q \in P$ and $x \in \mathbb{Z}^k$. For configurations $(p, x), (q, y) \in P \times \mathbb{Z}^k$, we write $(p, x) \to (q, y)$ if there is a transition $\tau \in \Delta$ with $\operatorname{src}(\tau) = p$, $\operatorname{tgt}(\tau) = q$, $\varphi_{\tau}(x)$ is true, and $y = x + u_{\tau}$. By $\stackrel{*}{\to}$, we denote the reflexive transitive closure of the relation \to on the configurations. A configuration (q, x) is initial if $q = q_0$. We say that an initial configuration (q_0, x) is accepted if $(q_0, x) \stackrel{*}{\to} (p, y)$ with $p \in F$. In other words, if there exists a run of the k-CM that eventually arrives in a final state.

We will employ the following variant of the fact that counter machines are Turing-complete. Note that if one uses CM as language acceptors, with input-reading transitions, then just two counters are sufficient for Turing-completeness. In our construction, it will be most convenient to provide the input of the CM at its counters. In this setting, it is known that three additional counters (aside from the input counters) are sufficient for Turing-completeness:

271

272 273

274

275

276

277

278

279

280

281

282

283

284 285

286

287

288

289

290

291 292

293

294

295 296

297

298 299

300

301 302

303

304

305

306 307

308

310

311

312

313 314

315 316

317

318 319

320 321

322

323

Lemma 3.5. For every recursively enumerable set $S \subseteq \mathbb{N}^n$, there is a (n+3)-CM so that for every $x \in \mathbb{N}^n$, the configuration $(q_0, x, 0, 0, 0)$ is accepted if and only if $x \in S$.

This is a direct consequence of CM, as language acceptors are able to recognize all recursively enumerable languages (this is implicit in (Minsky, 1961, Theorem Ia), and explicit in (Fischer et al., 1968, Theorem 3.1)) and that k-CM accept the same languages as 3-CM (Greibach, 1976, Theorem 2.4). Moreover, if $S \subseteq \mathbb{N}^n$ is recursively enumerable, then the language $L := \{a_1^{x_1} \cdots a_n^{x_n}\}$ $(x_1, \dots, x_n) \in S$ is a recursively enumerable language, and so there exists a three-counter machine M that recognizes L. This three-counter machine can easily be turned into a (n+3)-CM as we need it: whenever M reads a letter a_i , our CM will decrement the i-th counter; and when M uses counter $j \in \{1, 2, 3\}$, then our CM will use counter n + j.

Corollary 3.6. For every recursively enumerable permutation-invariant language $L \subseteq \Sigma^+$, there is a(n+3)-CM so that for every $w \in \Sigma^+$, we have $w \in L$ if and only if $(q_0, \Psi(w), 0, 0, 0)$ is accepted. *Proof.* Follows from Lemma 3.5: For a recursively enumerable $L \subseteq \Sigma^+$, the Parikh image $\Psi(L)$ is recursively enumerable; since L is permutation-invariant, we have $w \in L$ iff $\Psi(w) \in \Psi(L)$.

Proof of Proposition 3.4. Let $\Sigma = \{a_1, \ldots, a_n\}$ and take a permutation-invariant recursively enumerable language $L \subseteq \Sigma^*$. From Corollary 3.6, we get a (n+3)-CM such that from the configuration $C_0 := (q_0, x_1, \dots, x_n, 0, 0, 0)$, the CM will reach F if and only if $a_1^{x_1} \cdots a_n^{x_n} \in L$.

We define the set Γ of CoT tokens to be Σ unioned with the transition relation Δ . Note that the C-RASP is going to be evaluated at the last position on input wv where $v \in \Gamma^*$. The construction of the C-RASP CoT transformer considers the following cases.

Initial step. At the beginning, the last symbol in the input to the C-RASP is in Σ . This indicates that the CM is in the initial state q_0 . We add the following rules to our CoT C-RASP expression S

$$O_{\tau} \leftarrow \varphi(\overline{\#}[Q_{a_1}], \dots, \overline{\#}[Q_{a_n}]) \wedge Q_a,$$

for each $a \in \Sigma$ and each transition $\tau \in \Delta$ with $\operatorname{src}(\tau) = q_0$. The order in which the rules are added is not important since the counter machine is deterministic.

Non-initial step. After an initial step, the last symbol in the input is always a transition of the CM, which indicates which state the CM is in. We add the following rules to our CoT C-RASP expression S (in no particular order):

$$O_{\tau'} \leftarrow \varphi_{\tau'}(t_1, \dots, t_{n+3}) \wedge Q_{\tau},$$

for any $\tau, \tau' \in \Delta$ with $tgt(\tau) = src(\tau')$. Here, t_1, \ldots, t_{n+3} are the count-valued C-RASP terms

$$t_{i} = \overleftarrow{\#}[Q_{a_{i}}] + \sum_{\rho \in \Delta} \mathbf{u}_{\rho}(i) \cdot \overleftarrow{\#}[Q_{\rho}] \qquad \text{for } i = 1, \dots, n$$

$$t_{i} = \sum_{\rho \in \Delta} \mathbf{u}_{\rho}(i) \cdot \overleftarrow{\#}[Q_{\rho}] \qquad \text{for } i = n + 1, n + 2, n + 3.$$

$$(3)$$

$$t_i = \sum_{\rho \in \Delta} \mathbf{u}_{\rho}(i) \cdot \overleftarrow{\#}[Q_{\rho}]$$
 for $i = n + 1, n + 2, n + 3.$ (4)

Intuitively, each $[t_i]_w$ will tell us the value of the *i*th counter. For $i=1,\ldots,n$, we have the additional summand $\#[Q_{a_i}]$ because this is the initial value of the *i*th counter, according to Lemma 3.5.

Output symbols. The desired output symbols for acceptance are any $\tau \in \Delta$ for which $\operatorname{tgt}(\tau) \in F$.

Correctness. The C-RASP directly simulates the CM, so correctness is immediate.

Finally, Proposition 3.3 follows easily from Proposition 3.4: We can modify our C-RASP to check (e.g. in each step) that the (initial) input word belongs to $a_1^* \cdots a_n^*$. See Appendix B for the proof.

GENERAL CASE

Given that Propositions 3.3 and 3.4 show that for letter-bounded or permutation-invariant languages, CoT C-RASP are Turing-complete, this raises the question of whether they are even Turing-complete for a general language $L \subseteq \Sigma^*$. The following shows that they are not:

Proposition 4.1. *C-RASP in the CoT setting is not Turing-complete over* $\Sigma = \{a, b\}$.

This follows from the following lemma (e.g. take PALINDROME).

Lemma 4.2. If a language L is recognized by CoT C-RASP, then for each n the restriction $L_n \subseteq L$ to all inputs of length $\leq n$ is recognized by an automaton of size polynomial in n.

This is an immediate corollary of the logarithmic communication complexity of Limit Transformers and hence C-RASP (Theorem 12 in Huang et al. (2025)). See Appendix C for details.

However, we will show that in the presence of relative positional encodings, CoT C-RASP are in fact fully Turing-complete:

Theorem 4.3. Every recursively enumerable language over an arbitrary alphabet Σ can be recognized by C-RASP[RPEs] in the CoT setting.

The CoT C-RASP[RPEs] constructed in Theorem 4.3 is based on the following idea. Given an input $w \in \Sigma^*$ with say $|\Sigma| = n$, our CoT C-RASP[RPEs] first computes an encoding of $w \in \Sigma^*$ as a vector in \mathbb{N}^n . After this, it uses a construction similar to above to simulate a CM on this encoding.

To avoid confusion between multiplication of 0 and 1 on the one hand and concatenation of words, we will use different symbols for the numbers $0,1\in\mathbb{N}$ and the letters 0 and 1. Then for a=0, b=1, c=0, and d=1, we can distinguish between ab=0 and cd=01. To convert between these objects, we use the notation $\overline{0}:=0$, $\overline{0}=0$, $\overline{1}=1$, and $\overline{1}=1$.

Encoding words over two letters We first describe how to encode two-letter words. Formally, we have a partial function $\beta \colon \mathbb{N} \to \{0,1\}^*$, where \to means that β is partial, i.e. not every number represents a word. However, if a number represents a word, then it is unique. A number $x \in \mathbb{N}$ will represent a word if and only if $x \neq 0$. Hence, suppose $x \neq 0$. Then we can write $x = \sum_{i=0}^m b_i 2^i$, where $b_m, \ldots, b_0 \in \{0,1\}$, and $b_m = 1$. Let $j = \max\{i \mid b_i = 0\}$ be the left-most position of a zero when writing the most significant bit first. Then we set

$$\beta(x) := \overline{b_{j-1}} \, \overline{b_{j-2}} \, \cdots \, \overline{b_0}.$$

In other words, $\beta(x)$ is the word consisting of all digits of x's binary representation, when reading from most significant bit first, and starting after the left-most zero. For example, we have

$$\beta(2^5 + 2^3 + 2^1) = 1010,$$
 $\beta(2^6) = 00000,$ $\beta(2^4 + 2^3 + 2^1) = 10.$

Encoding words over arbitrary alphabets Now suppose Σ is an arbitrary alphabet with $\Sigma = \{a_1, \ldots, a_n\}$. Then we encode words in Σ^* by vectors in \mathbb{N}^n . Similar to above, we define a partial function $\sigma \colon \mathbb{N}^n \to \Sigma^*$. Let us first describe the domain of σ . We say that an n-tuple (w_1, \ldots, w_n) of words $w_1, \ldots, w_n \in \{0, 1\}^*$ is consistent if (i) the words w_1, \ldots, w_n have the same length, say $m \in \mathbb{N}$ and (ii) for every position $i \in [1, m]$, there is exactly one $j \in [1, n]$ such that w_j has the letter 1 at position i. Intuitively, the consistent n-tuples correspond exactly to the words in Σ^* : A word $w \in \Sigma^*$ of length m corresponds to the n-tuple (w_1, \ldots, w_n) where each w_i has length n, and the 1's in w_i are exactly at those positions that carry a_i in w. This leads to an intermediate partial function $\mu \colon (\{0,1\}^*)^n \to \Sigma^*$, where $\mu(w_1, \ldots, w_n)$ is defined if and only if (w_1, \ldots, w_n) is consistent, and in that case, $\mu(w_1, \ldots, w_n) \in \Sigma^*$ is the word corresponding to w_1, \ldots, w_n

With this, we are ready to define σ . The domain of σ consists of those $\boldsymbol{x}=(x_1,\ldots,x_n)\in\mathbb{N}^n$ where (i) all entries are non-zero and (ii) the tuple $(\beta(x_1),\ldots,\beta(x_n))$ is consistent. Moreover, for $\boldsymbol{x}=(x_1,\ldots,x_n)\in\mathrm{dom}\,\sigma$, we set

$$\sigma(\boldsymbol{x}) := \mu(\beta(x_1), \dots, \beta(x_n)).$$

For example, for n = 2, we have

$$\sigma(2^4 + 2^0, 2^4 + 2^2 + 2^1) = \mu(\beta(2^4 + 2^0, 2^4 + 2^2 + 2^1)) = \mu(001, 110) = a_2 a_2 a_1.$$

An important property of σ is that if we change $\mathbf{x} = (x_1, \dots, x_n)$ by introducing further 1's on the left of some binary representation of x_i , then $\sigma(\mathbf{x})$ remains the same. For example, we also have

$$\sigma(2^5 + 2^4 + 2^0, 2^4 + 2^2 + 2^1) = \mu(\beta(2^5 + 2^4 + 2^0, 2^4 + 2^2 + 2^1)) = \mu(001, 110) = a_2 a_2 a_1.$$

even though we modified the left-most entry by introducing the term 2^5 . This means, for every word $w \in \Sigma^*$ and every bound $k \in \mathbb{N}$, there is an $x \in \mathbb{N}^n$ such that (i) all entries in x are $\ge k$ and (ii) $\sigma(x) = w$. This will be important in the proof.

The relative positional encoding A key ingredient in our proof is the relative positional encoding (recall that we have shown that without RPE, Theorem 4.3 does not hold). Perhaps surprisingly, the RPE we use in the proof does not depend on the language we are accepting: It is the same relation for every Turing machine we want to simulate. Its definition is based on the partial function $\beta \colon \mathbb{N} \to \{0,1\}^*$ above. We define the relation $\mathfrak{R} \subseteq \mathbb{N} \times \mathbb{N}$ as

 $(i,j)\in\mathfrak{R}\iff i\leqslant j,\,i\in[1,|\beta(j)|],$ and the word $\beta(j)\in\{0,1\}^*$ has 1 at position i for every $(i,j)\in\mathbb{N}\times\mathbb{N}.$ For example, if $j=2^6+2^5+2^3+2^1+2^0$, then we have $\beta(j)=$ 1011 and hence $(1,j),(3,j),(4,j)\in\mathfrak{R},$ but $(2,j)\notin\mathfrak{R}.$

Overview Our C-RASP with CoT will work in *two phases*. During the *first phase*, it prolongs the input so that subsequently, a σ -encoding of the original input word can be computed using Count-Valued Operations. For this, it relies on the RPE \Re . In the *second phase*, our C-RASP simulates a counter machine, similar to the permutation-invariant case.

Phase I: Constructing encoding of the input word In order to compute the σ -encoding $x \in \mathbb{N}^n$ of the input word $w \in \Sigma^*$, our CoT C-RASP proceeds as follows. It compute the entries $x(1), \ldots, x(n)$ of x in this order. Suppose (w_1, \ldots, w_n) is the consistent tuple representing w, i.e. $\mu(w_1, \ldots, w_n) = w$. To compute x(1), our CoT C-RASP appends a dummy letter \square_1 until the current word length ℓ satisfies $\beta(\ell) = w_1$. Note that this is possible since there are infinitely many ℓ with $\beta(\ell) = w_1$. Once this holds, we place a special letter \boxplus_1 . Then, the CoT C-RASP appends a dummy letter \sqsubseteq_2 until the current word length satisfies $\beta(\ell) = w_2$, and then places \boxplus_2 , etc.

Initially, the last letter will be some $a \in \Sigma$. Then, our CoT C-RASP simply outputs \square_1 : We have

$$O_{\square_1} \leftarrow Q_a$$

for every $a \in \Sigma$. When we have a letter \square_i at the end, our CoT C-RASP checks whether the current length ℓ already satisfies $\beta(\ell) = w_i$:

$$O_{\boxplus_i} \leftarrow Q_{\square_i} \land \overleftarrow{\#}_{\mathfrak{R}}[Q_{a_i}] = \overleftarrow{\#}[Q_{a_i}] \land \overleftarrow{\#}_{\mathfrak{R}}[\top] = \overleftarrow{\#}[Q_{a_i}]$$
 (5)

$$O_{\square_i} \leftarrow Q_{\square_i} \land (\overleftarrow{\#}_{\mathfrak{R}}[Q_{a_i}] \neq \overleftarrow{\#}[Q_{a_i}] \lor \overleftarrow{\#}_{\mathfrak{R}}[\top] \neq \overleftarrow{\#}[Q_{a_i}]) \tag{6}$$

for each $i=1,\ldots,n$. If we evaluate rule 5 on a word of length ℓ , we check that (i) the last letter is \square_i , (ii) the number of positions j with $(j,\ell) \in \mathfrak{R}$ that carry a_i equals the total number of positions that carry a_i , and (iii) the number of positions j with $(j,\ell) \in \mathfrak{R}$ equals the number of positions that carry a_i . Thus, conditions (ii) and (iii) say that the positions j with $(j,\ell) \in \mathfrak{R}$ are precisely those that carry an a_i . In other words, $\beta(\ell) = w_i$. If these conditions are met, then the output letter is \boxplus_i .

Moreover, if we evaluate rule 6, we check that $\beta(\ell)$ does not equal w_i yet. In this case, the output letter is again \square_i , and the whole check will be repeated with the next word length.

If the last letter is \coprod_i with $i \le n-1$, then we start computing x(i+1): We output \coprod_{i+1} in 7:

$$O_{\square_{i+1}} \leftarrow Q_{\boxplus_i}$$
 for each $i = 1, \dots, n-1$ (7)

$$O_{\tau} \leftarrow Q_{\mathbb{H}_n}$$
 for each transition $\tau \in \Delta$ with $\operatorname{src}(\tau) = q_0$ (8)

If the last letter is \coprod_n , we initiate the CM run by outputting some initial transition τ . This is rule 8.

After the above process, we have placed $\boxplus_1, \dots, \boxplus_n$. Thus, the current input word is then of the form $w' = w \square_1^{f_1} \boxplus_1 \square_2^{f_2} \boxplus_2 \dots \square_n^{f_n} \boxplus_n$, where for the tuple $\mathbf{x} = (x_1, \dots, x_n)$ with $x_i = |w| + f_1 + \dots + f_i$, we have $\sigma(\mathbf{x}) = w$. A count-valued operation can then access the encoding of w using the terms

$$X_i = \# [\# [\boxplus_i] = 0] \qquad \text{for } i = 1, \dots, n \tag{9}$$

Thus, X_i is the number of positions that have no occurrence of \coprod_i to their left (and do not carry \coprod_i themselves). Since there is exactly one occurrence of \coprod_i , this means X_i is exactly the position of \coprod_i , minus one. Therefore, the term X_i evaluates to x(i), meaning we have $\sigma(X_1, \ldots, X_n) = w$.

Phase II: Simulating the counter machine During the first phase, our CoT C-RASP appended letters to make an encoding $x \in \mathbb{N}^n$ of the input word available through C-RASP terms Eq. (9). We now use a CM that starts with this encoding in its counters and then decides whether $w \in L$. Such a counter machine exists because of Lemma 3.5 and the fact that $S = \{x \in \mathbb{N}^n \mid \sigma(x) \in L\}$ is recursively enumerable (since σ is computable). The simulation of the CM on x works exactly like in Section 3, except that in the terms defined in equation 3, instead of using $\#[Q_{a_i}]$ for $i = 1, \ldots, n$, we use the C-RASP term X_i defined in equation 9. See Appendix C for details.

5 EMPIRICAL EXPERIMENTS

We empirically validate our Turing-completeness results on some complex arithmetical concepts. Our theory predicts that CoT C-RASP with NoPE suffices for unary representation (of numbers), while RPEs are needed for binary representation. Accordingly, we conduct three experiments: 1) unary without \Re , 2) binary without \Re , and 3) binary with \Re , here \Re refers to Section 4.

The arithmetic tasks include Prime, Exponential, Division, Greatest Common Divisor, and Multiplication; see Appendix D for details. For each task, we construct two counter machines (CMs), one for the unary representation and one for the binary representation. We then generate a training set of 40k samples and a validation set val_0 of 10k samples over the range [1–100]. In addition, we generate two out-of-distribution validation sets, val_1 and val_2 , each containing 10k samples over the ranges [101–200] and [201–300], respectively. The SMATs are trained using AdamW (weight decay 0.01), a batch size of 64, and up to 30k (60k for harder tasks) steps, with early stopping once out-of-distribution accuracy reaches 100%. For unary, we use an embedding size of 128, 2 layers, 4 heads, and dropout 0.0. For binary, we use an embedding size of 384, 6 layers, 4 heads, an initial scaling factor of $\lambda=1$, and dropout 0.0.

As shown in Table 1, SMAT achieves strong in-distribution performance over unary representations, with val_0 accuracies exceeding 99.7%. Although generalization on val_2 degrades slightly, it still learns these tasks quite well. In contrast, the binary representation with $\mathfrak R$ yields near-perfect accuracy on val_0 and val_1 , and outstanding performance on val_2 . Without $\mathfrak R$, however, SMAT models achieve comparable accuracy only on val_0 , performance on val_1 and val_2 drops to nearly zero. These results demonstrate that unary benefits from NoPE, whereas binary requires $\mathfrak R$ for length generalization.

Language	Unary representation			Binary representation		
Language	$val_{0(\%)}$	$val_{1(\%)}$	$val_{2(\%)}$	$val_{0(\%)}$	$val_{1(\%)}$	$val_{2(\%)}$
Prime	100.00	100.00	100.00	100.00	100.00	99.69
Exponential	99.98	99.92	99.88	100.00	100.00	99.92
Division	99.92	99.03	99.51	100.00	100.00	84.20
Greatest Common Divisor	99.94	99.92	99.30	99.93	99.33	93.23
Multiplication	99.79	99.96	99.93	99.93	99.71	99.22

Table 1: Generalization Performance Across Validation Splits for Unary and Binary Representations

6 CONCLUDING REMARKS

Related work. Similar to our work, Hou et al. (2025) aims to provide length-generalizing constructions for Turing completeness. However, there are two key differences. First, we demonstrate the existence of softmax transformer constructions, whereas Hou et al. (2025) only demonstrated constructions in RASP (Weiss et al., 2021). Second, the approach of Hou et al. (2025) ensures length generalization only if no n-grams are repeated, for some fixed n, which is likely to be unrealistic in the limit of long inputs. In contrast, our approach theoretically ensures full-length generalizability.

Future work. Recent results have refined Turing-completeness for transformers (albeit with hard attention) by relating the number of CoT steps and complexity classes. For example, in (Merrill & Sabharwal, 2024), CoT transformers with polynomially number of CoT steps correspond to classes solvable by polynomial-time algorithms. Similar results were also recently derived in (Li & Wang, 2025), which relate to space complexity. We leave it for future work to refine our Turing-completeness results with computational complexity.

REFERENCES

Pablo Barceló, Alexander Kozachinskiy, Anthony Widjaja Lin, and Vladimir V. Podolskii. Logical languages accepted by transformer encoders with hard attention. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net, 2024. URL https://openreview.net/forum?id=gbrHZq07mq.

- Satwik Bhattamishra, Arkil Patel, and Navin Goyal. On the computational power of transformers and its implications in sequence modeling. In *CoNLL*, pp. 455–475. Association for Computational Linguistics, 2020.
 - David Chiang and Peter Cholak. Overcoming a theoretical limitation of self-attention. In Smaranda Muresan, Preslav Nakov, and Aline Villavicencio (eds.), *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2022, Dublin, Ireland, May* 22-27, 2022, pp. 7654–7664. Association for Computational Linguistics, 2022. doi: 10.18653/V1/2022.ACL-LONG.527. URL https://doi.org/10.18653/v1/2022.acl-long.527.
 - Zihang Dai, Zhilin Yang, Yiming Yang, Jaime G Carbonell, Quoc Le, and Ruslan Salakhutdinov. Transformer-xl: Attentive language models beyond a fixed-length context. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 2978–2988, 2019.
 - Philipp Dufter, Martin Schmitt, and Hinrich Schütze. Position information in transformers: An overview. *Computational Linguistics*, 48(3):733–763, 09 2022. ISSN 0891-2017. doi: 10.1162/coli_a_00445. URL https://doi.org/10.1162/coli_a_00445.
 - Benjamin L Edelman, Surbhi Goel, Sham Kakade, and Cyril Zhang. Inductive biases and variable creation in self-attention mechanisms. In *International Conference on Machine Learning*, pp. 5793–5831. PMLR, 2022.
 - Patrick C Fischer, Albert R Meyer, and Arnold L Rosenberg. Counter machines and counter languages. *Mathematical systems theory*, 2(3):265–283, 1968. doi: 10.1007/BF01694011.
 - Sheila A. Greibach. Remarks on the complexity of nondeterministic counter languages. *Theor. Comput. Sci.*, 1(4):269–288, 1976. doi: 10.1016/0304-3975(76)90072-4.
 - Michael Hahn. Theoretical limitations of self-attention in neural sequence models. *Transactions of the Association for Computational Linguistics*, 8:156–171, 2020.
 - Pengcheng He, Xiaodong Liu, Jianfeng Gao, and Weizhu Chen. Deberta: Decoding-enhanced bert with disentangled attention. In *International Conference on Learning Representations*, 2021.
 - Kaiying Hou, Eran Malach, Samy Jelassi, David Brandfonbrener, and Sham M. Kakade. Universal length generalization with turing programs. In *Forty-second International Conference on Machine Learning*, 2025. URL https://openreview.net/forum?id=RNSd6G3lcD.
 - Xinting Huang, Andy Yang, Satwik Bhattamishra, Yash Raj Sarrof, Andreas Krebs, Hattie Zhou, Preetum Nakkiran, and Michael Hahn. A formal framework for understanding length generalization in transformers. In *The Thirteenth International Conference on Learning Representations, ICLR 2025, Singapore, April 24-28, 2025.* OpenReview.net, 2025. URL https://openreview.net/forum?id=U49N5V51rU.
 - Qian Li and Yuyi Wang. Constant bit-size transformers are turing complete. *CoRR*, abs/2506.12027, 2025. doi: 10.48550/ARXIV.2506.12027. URL https://doi.org/10.48550/arXiv.2506.12027.
 - Antoine Liutkus, Ondrej Cífka, Shih-Lun Wu, Umut Simsekli, Yi-Hsuan Yang, and Gaël Richard. Relative positional encoding for transformers with linear complexity. In Marina Meila and Tong Zhang (eds.), *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event*, volume 139 of *Proceedings of Machine Learning Research*, pp. 7067–7079. PMLR, 2021. URL http://proceedings.mlr.press/v139/liutkus21a.html.
 - William Merrill and Ashish Sabharwal. The expressive power of transformers with chain of thought. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024.* OpenReview.net, 2024. URL https://openreview.net/forum?id=NjNGlPh8Wh.
 - Marvin L Minsky. Recursive unsolvability of post's problem of 'tag' and other topics in theory of turing machines. *Annals of Mathematics*, 74(3):437–455, 1961. doi: 10.2307/1970290.

541

543

544

546

547

548

549

550

551

552

553

554

558

559

561

564

565

566

567

568

569

570

571

572

573

574

575

576

577

578

579

580

581

582

583

584

585

586

588

590

592

OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mohammad Bavarian, Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher Berner, Lenny Bogdonoff, Oleg Boiko, Madelaine Boyd, Anna-Luisa Brakman, Greg Brockman, Tim Brooks, Miles Brundage, Kevin Button, Trevor Cai, Rosie Campbell, Andrew Cann, Brittany Carey, Chelsea Carlson, Rory Carmichael, Brooke Chan, Che Chang, Fotis Chantzis, Derek Chen, Sully Chen, Ruby Chen, Jason Chen, Mark Chen, Ben Chess, Chester Cho, Casey Chu, Hyung Won Chung, Dave Cummings, Jeremiah Currier, Yunxing Dai, Cory Decareaux, Thomas Degry, Noah Deutsch, Damien Deville, Arka Dhar, David Dohan, Steve Dowling, Sheila Dunning, Adrien Ecoffet, Atty Eleti, Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix, Simón Posada Fishman, Juston Forte, Isabella Fulford, Leo Gao, Elie Georges, Christian Gibson, Vik Goel, Tarun Gogineni, Gabriel Goh, Rapha Gontijo-Lopes, Jonathan Gordon, Morgan Grafstein, Scott Gray, Ryan Greene, Joshua Gross, Shixiang Shane Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris, Yuchen He, Mike Heaton, Johannes Heidecke, Chris Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele, Brandon Houghton, Kenny Hsu, Shengli Hu, Xin Hu, Joost Huizinga, Shantanu Jain, Shawn Jain, Joanne Jang, Angela Jiang, Roger Jiang, Haozhun Jin, Denny Jin, Shino Jomoto, Billie Jonn, Heewoo Jun, Tomer Kaftan, Łukasz Kaiser, Ali Kamali, Ingmar Kanitscheider, Nitish Shirish Keskar, Tabarak Khan, Logan Kilpatrick, Jong Wook Kim, Christina Kim, Yongjik Kim, Jan Hendrik Kirchner, Jamie Kiros, Matt Knight, Daniel Kokotajlo, Łukasz Kondraciuk, Andrew Kondrich, Aris Konstantinidis, Kyle Kosic, Gretchen Krueger, Vishal Kuo, Michael Lampe, Ikai Lan, Teddy Lee, Jan Leike, Jade Leung, Daniel Levy, Chak Ming Li, Rachel Lim, Molly Lin, Stephanie Lin, Mateusz Litwin, Theresa Lopez, Ryan Lowe, Patricia Lue, Anna Makanju, Kim Malfacini, Sam Manning, Todor Markov, Yaniv Markovski, Bianca Martin, Katie Mayer, Andrew Mayne, Bob McGrew, Scott Mayer McKinney, Christine McLeavey, Paul McMillan, Jake McNeil, David Medina, Aalok Mehta, Jacob Menick, Luke Metz, Andrey Mishchenko, Pamela Mishkin, Vinnie Monaco, Evan Morikawa, Daniel Mossing, Tong Mu, Mira Murati, Oleg Murk, David Mély, Ashvin Nair, Reiichiro Nakano, Rajeev Nayak, Arvind Neelakantan, Richard Ngo, Hyeonwoo Noh, Long Ouyang, Cullen O'Keefe, Jakub Pachocki, Alex Paino, Joe Palermo, Ashley Pantuliano, Giambattista Parascandolo, Joel Parish, Emy Parparita, Alex Passos, Mikhail Pavlov, Andrew Peng, Adam Perelman, Filipe de Avila Belbute Peres, Michael Petrov, Henrique Ponde de Oliveira Pinto, Michael, Pokorny, Michelle Pokrass, Vitchyr H. Pong, Tolly Powell, Alethea Power, Boris Power, Elizabeth Proehl, Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh, Cameron Raymond, Francis Real, Kendra Rimbach, Carl Ross, Bob Rotsted, Henri Roussez, Nick Ryder, Mario Saltarelli, Ted Sanders, Shibani Santurkar, Girish Sastry, Heather Schmidt, David Schnurr, John Schulman, Daniel Selsam, Kyla Sheppard, Toki Sherbakov, Jessica Shieh, Sarah Shoker, Pranav Shyam, Szymon Sidor, Eric Sigler, Maddie Simens, Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin Sokolowsky, Yang Song, Natalie Staudacher, Felipe Petroski Such, Natalie Summers, Ilya Sutskever, Jie Tang, Nikolas Tezak, Madeleine B. Thompson, Phil Tillet, Amin Tootoonchian, Elizabeth Tseng, Preston Tuggle, Nick Turley, Jerry Tworek, Juan Felipe Cerón Uribe, Andrea Vallone, Arun Vijayvergiya, Chelsea Voss, Carroll Wainwright, Justin Jay Wang, Alvin Wang, Ben Wang, Jonathan Ward, Jason Wei, CJ Weinmann, Akila Welihinda, Peter Welinder, Jiayi Weng, Lilian Weng, Matt Wiethoff, Dave Willner, Clemens Winter, Samuel Wolrich, Hannah Wong, Lauren Workman, Sherwin Wu, Jeff Wu, Michael Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qiming Yuan, Wojciech Zaremba, Rowan Zellers, Chong Zhang, Marvin Zhang, Shengjia Zhao, Tianhao Zheng, Juntang Zhuang, William Zhuk, and Barret Zoph. Gpt-4 technical report, 2024. URL https://arxiv.org/abs/2303.08774.

Jorge Pérez, Pablo Barceló, and Javier Marinkovic. Attention is turing-complete. *J. Mach. Learn. Res.*, 22:75:1–75:35, 2021. URL https://jmlr.org/papers/v22/20-302.html.

Ofir Press, Noah Smith, and Mike Lewis. Train short, test long: Attention with linear biases enables input length extrapolation. In *International Conference on Learning Representations*, 2022.

Ruizhong Qiu, Zhe Xu, Wenxuan Bao, and Hanghang Tong. Ask, and it shall be given: On the turing completeness of prompting. In *The Thirteenth International Conference on Learning Representations*, 2025. URL https://openreview.net/forum?id=AS8SPTyBgw.

Peter Shaw, Jakob Uszkoreit, and Ashish Vaswani. Self-attention with relative position representations. In Marilyn Walker, Heng Ji, and Amanda Stent (eds.), *Proceedings of the 2018 Con-*

ference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers), pp. 464–468, New Orleans, Louisiana, June 2018. Association for Computational Linguistics. doi: 10.18653/v1/N18-2074. URL https://aclanthology.org/N18-2074/.

Lena Strobl, William Merrill, Gail Weiss, David Chiang, and Dana Angluin. What formal languages can transformers express? A survey. *Trans. Assoc. Comput. Linguistics*, 12:543–561, 2024. doi: 10.1162/TACL_A_00663. URL https://doi.org/10.1162/tacl_a_00663.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In Isabelle Guyon, Ulrike von Luxburg, Samy Bengio, Hanna M. Wallach, Rob Fergus, S. V. N. Vishwanathan, and Roman Garnett (eds.), Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA, pp. 5998–6008, 2017. URL https://proceedings.neurips.cc/paper/2017/hash/3f5ee243547dee91fbd053c1c4a845aa-Abstract.html.

Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed H. Chi, Quoc V. Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language models. In *Proceedings of the 36th International Conference on Neural Information Processing Systems*, NIPS '22, Red Hook, NY, USA, 2022. Curran Associates Inc. ISBN 9781713871088.

Gail Weiss, Yoav Goldberg, and Eran Yahav. Thinking like transformers. In *International Conference on Machine Learning*, pp. 11080–11090. PMLR, 2021.

Linting Xue, Noah Constant, Adam Roberts, Mihir Kale, Rami Al-Rfou, Aditya Siddhant, Aditya Barua, and Colin Raffel. mt5: A massively multilingual pre-trained text-to-text transformer.
 In Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pp. 483–498, 2021.

Andy Yang and David Chiang. Counting like transformers: Compiling temporal counting logic into softmax transformers. *CoRR*, abs/2404.04393, 2024. doi: 10.48550/ARXIV.2404.04393. URL https://doi.org/10.48550/arXiv.2404.04393.

Andy Yang, Lena Strobl, David Chiang, and Dana Angluin. Simulating hard attention using soft attention. *CoRR*, abs/2412.09925, 2024. doi: 10.48550/ARXIV.2412.09925. URL https://doi.org/10.48550/arXiv.2412.09925.

A Additional material on Section 2

A.1 FORMAL DEFINITION OF SOFTMAX TRANSFORMERS.

Our definition of softmax transformers follows that of Huang et al. (2025), though we use a highly simplified notation here for exposition. In a SoftMax Averaging Transformers (SMAT), each layer consists of two affine transformations A, B and a feedforward network C. Given a sequence

$$\mathbf{v}_1, \dots, \mathbf{v}_n$$

the layer outputs

$$\mathbf{w}_1, \ldots, \mathbf{w}_n$$

where

$$\mathbf{w}_i := \mathbf{v}_i + C(\mathbf{v}_i')$$

where $\mathbf{v}_i' := \sum_{j=1}^i \bar{w}(j)\mathbf{v}_j$ and

$$\bar{w} = \operatorname{softmax}(\log n \cdot \{\mathbf{v}_{j}^{T} \mathbf{K}^{T} \mathbf{Q} \mathbf{v}_{i}\}_{j=1}^{i})$$
(10)

where \mathbf{v}_i denotes activations at position i, and \mathbf{K} , \mathbf{Q} transform these to keys and queries, respectively. Here, scaling with $\log n$ is included, as it is needed to theoretically represent sparse functions across unboundedly input strings and circumvent theoretical limitations of soft attention (Chiang & Cholak, 2022; Edelman et al., 2022). Here, we show the case of a single head, extension to multiple heads is straightforward.

We assume C is a one-layer feedforward layer, where each hidden unit has either ReLU or Heaviside activation. Here, as in Huang et al. (2025), Heaviside is needed to theoretically represent functions with sharp thresholds; at any finite input length, it can be arbitrarily closely approximated using ReLU MLPs.

Huang et al. (2025) also assume that attention logits are rounded to fixed precision; we do not require this for our results here. Also, whereas Huang et al. (2025) consider Absolute Positional Encodings (APE), which necessitated introducing fixed context windows and positional offsets, we do not consider APE here, and so do not need to introduce offsets. Thus, SMATs considered in the present paper are uniformly applicable to arbitrarily long inputs.

To interface SMAT with an input string $w \in \Sigma^+$, we apply a token embedding function $em: \Sigma \to \mathbb{R}^k$ for some dimension k; these are followed by some number of SMAT layers. To define a CoT SMAT, we need the transformer to be able to output a token. To this end, we define an output function $o: \mathbb{R}^d \to \Sigma$, parameterized by applying a linear function $\mathbb{R}^d \to \mathbb{R}^{|\Sigma|}$ followed by an argmax selecting the symbol receiving the highest score.

Overall, we view an SMAT as a length-preserving map $T: \Sigma^* \to \Sigma^*$, where $T(x)_i$ indicates the symbol predicted after reading the prefix $x_1 \dots x_i$.

Discussion Our formalization of SMAT follows the setting of Huang et al. (2025), which was designed to study the learnability of transformers. We note two aspects, which are needed to enable softmax transformers to represent functions across arbitrarily long inputs, and overcome well-known theoretical limitations of softmax attention (Hahn, 2020; Chiang & Cholak, 2022). First, scaling attention logits with $\log n$ is necessary to represent sparse attention to specific positions, which otherwise would be impossible to achieve using softmax attention (Hahn, 2020; Chiang & Cholak, 2022; Edelman et al., 2022). Importantly, this scaling does not involve any new learnable parameters. Second, using Heaviside activations is necessary to represent functions with sharp thresholds, as is needed to perform exact comparison of counts across unboundedly long lengths. At any finite input length, Heaviside can be arbitrarily closely approximated using ReLU MLPs. We view Heaviside (which is not differentiable) as a theoretical proxy for steep ReLU network as is standardly trainable.

A.2 PROOFS FOR COT EXPRESSIVENESS AND LEARNABILITY

Proof of Proposition 2.1. This is a simple extension of Theorem 9 in Huang et al. (2025), as we now explain.

We define a CoT as a map $\Sigma^* \to \Sigma^*$ from an input string $w \in \Sigma^*$ to the sequence $w_2 \dots, w_N$ generated by a CoT C-RASP or CoT SMAT on the input string w. Starting from a CoT generated by a CoT C-RASP program, we aim to translate it to a CoT generated by a CoT SMAT.

We first explain the case without RPEs. We need to show that, if a CoT is generated in C-RASP CoT, then there is an SMAT generating the same CoT. In the case of language acceptance by a single binary label computed at the final token, Theorem 9 in Huang et al. (2025) shows that C-RASP can be simulated by a *limit transformer* without positional information. Our first observation is that, in the model of Huang et al. (2025), a limit transformer without positional information is equivalent to a standard transformer without positional encodings and infinite context window, which in turn is equivalent to an SMAT as defined in our paper here. The proof of Theorem 9 in Huang et al. (2025) builds a transformer that computes the values of all boolean predicates computed in the C-RASP program at each position in the string, with one dimension in the model's activations fo each boolean predicate. This means that the truth values of the expressions φ_{a_i} appearing in the switch condition S can also be computed. In order to evaluate the switch condition, we add another layer (whose attention heads have zero value matrices, i.e., don't contribute), then linearly project the relevant entries onto a binary vector of length $|\Gamma|$, and apply a piecewise linear function to convert this into a one-hot vector selecting the lowest-index token a_i such that φ_{a_i} is true. We now have a limit transformer which at each position outputs a one-hot vector indicating which CoT token to output. This means, whenever a CoT is expressible in C-RASP CoT, it is also expressible by SMAT with CoT.

We now consider the case with RPEs. We again build on Theorem 9 in Huang et al. (2025). We first note that the definition of attention logits with RPE exactly matches the definition of attention logits

in Limit Transformers with functions ϕ in Huang et al. (2025), where $\phi(i,j)$ is simply $[\![R]\!](i,j)$. Hence, for the purpose of expressivity, any SMAT[RPEs] transformer is equivalent to a limit transformer. Then, when translating from C-RASP to SMAT, implementing an RPE into an attention head proceeds along exactly the same lines as the translation of the special case $\#[j\leqslant i:\psi(i,j)]P(j)$ in the proof of that theorem.

Proof of 2.3. We first consider the case without RPEs. We build on Theorem 7 in Huang et al. (2025) and its variant for transformers without positional encodings, Corollary 18 in Huang et al. (2025). First, from Proposition 2.1, we know that if a language is expressible in C-RASP CoT, then it is also expressible by SMAT with CoT. The proof of that proposition further notes that our model of SMAT is equivalent to a limit transformer without positional information. Then, by Corollary 18 in Huang et al. (2025), any input-output map expressible by a limit transformer without positional information is length-generalizably learnable. This proves the result for the case without RPEs.

We now consider the case with RPEs. The proof is similar to Proposition 2.3; however, we need to (i) show that C-RASP[RPEs] can be simulated by SMATs with RPE, (ii) length generalization for SMAT RPE transformers follows from expressibility by SMATs with RPE. First, regarding (i), we again build on Theorem 9 in Huang et al. (2025), extending our argument from the proof of Proposition 2.1. We first note that the definition of attention logits with RPE exactly matches the definition of attention logits in Limit Transformers with functions ϕ in Huang et al. (2025), where $\phi(i,j)$ is simply $[\![R]\!](i,j)$. Hence, for the purpose of expressivity, any SMAT[RPEs] transformer is equivalent to a limit transformer. Then, when translating from C-RASP to SMAT, implementing an RPE into an attention head proceeds along exactly the same lines as the translation of the special case $\#[j \le i : \psi(i,j)]P(j)$ in the proof of that theorem. Second, regarding (ii), we use Corollary 18 in Huang et al. (2025) and note that the addition of fixed (not learned) RPE to attention heads in both the learned transformers and limit transformers has no impact on the argument.

B ADDITIONAL MATERIAL ON SECTION 3

In this subsection, we prove Proposition 3.3 from Proposition 3.4.

Suppose $\Sigma = \{a_1, \dots, a_n\}$. If $L \subseteq a_1^* \cdots a_n^*$ is recursively enumerable, then so is the language $K = \{u \in \Sigma^* \mid \exists v \in L \colon \Psi(u) = \Psi(v)\}$ of all permutations of L. Moreover, K is permutation-invariant, and thus recognized by a CoT C-RASP according to Proposition 3.4. Since $L = K \cap a_1^* \cdots a_n^*$, to turn that CoT C-RASP into a CoT C-RASP for L, it remains to check that the input word belongs to the set $a_1^* \cdots a_n^*$. Therefore, for all rules $O_a \leftarrow P$, where P is a C-RASP expression, we use

$$O_a \leftarrow P \land \bigwedge_{1 \le i < j \le n} \overset{\leftarrow}{\#} [Q_{a_i} \land \overset{\leftarrow}{\#} [Q_{a_j}] > 0] = 0,$$

where the second conjunct says that there are no positions carrying an a_i that have at least one a_j with j > i to their left. Then, the modified C-RASP clearly recognizes $K \cap a_1^* \cdots a_n^* = L$.

C ADDITIONAL MATERIAL ON SECTION 4

Details of Phase II In this section, we present the details of Phase II of the construction in Section 4. For this, first observe that

$$S = \{ \boldsymbol{x} \in \mathbb{N}^n \mid \sigma(\boldsymbol{x}) \in L \}$$

is recursively enumerable (since σ is computable). is recursively enumerable, since the partial function σ is computable. Therefore, by Lemma 3.5, there is a (n+3)-counter machine (P, Δ, q_0, F) such that for any $\boldsymbol{x} \in \mathbb{N}^n$, we have $\boldsymbol{x} \in S$ if and only if from the configuration $(q_0, \boldsymbol{x}, 0, 0, 0)$, the counter machine eventually reaches a control state in F.

We simulate a step of the counter machine using the following rule. If the CoT C-RASP finds the letter τ as the last letter, then for each possible next transition τ' , it checks whether its guard $\varphi_{\tau'}$ is satisfied, and if so, executes τ' by outputting τ' . Thus, we have

$$O_{\tau'} \leftarrow \varphi_{\tau'}(t_1, \dots, t_{n+3}) \wedge Q_{\tau}$$

for any two transitions $\tau, \tau' \in \Delta$ for which $\operatorname{tgt}(\tau) = \operatorname{src}(\tau')$. Here, t_1, \ldots, t_{n+3} are the following terms:

$$t_i = X_i + \sum_{
ho \in \Delta} oldsymbol{u}_
ho(i) \cdot \overleftarrow{\#}[Q_
ho] \qquad \qquad ext{for } i = 1, \dots, n ext{, and}$$

$$t_i = \sum_{
ho \in \Delta} oldsymbol{u}_
ho(i) \cdot \overleftarrow{\#}[Q_
ho] \qquad \qquad ext{for } i = n+1, n+2, n+3,$$

where X_i is the count-valued C-RASP term from (9). For $i \in \{n+1, n+2, n+3\}$, t_i is just the sum of counter effects on counter i. Equivalently, t_i is the current value of counter i after executing all these transitions. For $i \in [1, n]$, t_i we also add X_i , which has the effect that the counters $1, \ldots, n$ are initialized with X_i .

Finally, our CoT C-RASP accepts if the output symbol is any $\tau \in \Delta$ with $\operatorname{tgt}(\tau) \in F$.

Other Proofs

Proof of Lemma 4.2. If L is recognized by a CoT C-RASP, then it is also recognized by an SMAT C-RASP by Lemma 2.1. In fact, our model of SMAT is equivalent to the NoPE special case of the Limit Transformers of Huang et al. (2025). Now Theorem 12 in Huang et al. (2025) shows the following: Take any k. For each string $w \in \Sigma^*$, let $F(w) \in \Gamma^* \cup \Gamma^\omega$ be the associated CoT by which the language is recognized via an SMAT. Assume Alice has access to the prefix of wF(w) of length k, and Bob has access to the remainder, then Alice needs to communicate just $\mathcal{O}(\log k)$ bits to allow Bob to compute the output of the SMAT at all positions $k+1,k+2,\ldots$. In fact, Theorem 12 in Huang et al. (2025) is stated for the special case where k is half the input length, but the argument is entirely general, as it only relies on the length of Alice's part.

Note that, if the CoT terminates before k - |w| steps, Alice can just communicate that. Now given the SMAT recognizes L via CoT, Bob can determine² from Alice's communication if a given string is in the language or not.

Now we construct a family of NFAs accepting the language as follows.

For $x, y \in \Sigma^*$, define $x \equiv_{AB} y$ if and only if, for all $z \in \Sigma^*$, Alice communicates the same to Bob on xz and yz. By definition, each equivalence class of this relation is a subclass of a Nerode equivalence class of L (\dagger).

Given any length bound $n \in \mathbb{N}$, let Q_n be the set of all \equiv_{AB} -classes represented by at least some words of length $\leqslant n$. By the result described above, $|Q_n|$ is bounded by $\leqslant \sum_{k=1}^n 2^{\mathcal{O}(\log k)} = \mathcal{O}(poly(n))$. Now, by definition of the congruence, Q_n is the state set of an automaton computing \equiv_{AB} -equivalence classes. By (\dagger), it recognizes L.

D ADDITIONAL MATERIAL ON SECTION 5

Language	Unary Representation	Binary Representation
Prime	$\{ a^p : p \in \mathbb{P} \}$	$\{ bin(p) : p \in \mathbb{P} \}$
Exponential	$\{a^{2^i}: i \geqslant 0\}$	$\{ \operatorname{bin}(i) \# \operatorname{bin}(j) : j = 2^i \}$
Division	$\{a^ib^j:j\mid i\}$	$\{ bin(i) \# bin(j) : j \mid i \}$
Greatest Common Divisor	$\{a^ib^jc^k: k=\gcd(i,j)\}$	$\{ bin(i) #bin(j) #bin(k) : k = gcd(i, j) \}$
Multiplication	$\{a^ib^jc^k:k=i\cdot j\}$	$\{ bin(i) \# bin(j) \# bin(k) : k = i \times j \}$

Table 2: unary and binary representation of arithmetic languages. Here \mathbb{P} is the set of prime numbers, $j \mid i$ denotes divisibility, $\gcd(i,j)$ is the greatest common divisor, and $i \times j$ is multiplication.

²This is not decidable, but Bob in this model is a computationally unconstrained agent, with communication between Alice and Bob as the only bottleneck.