Continuous Chain of Thought Enables Parallel Exploration and Reasoning

Anonymous Authors¹

Abstract

Current language models generate chain-ofthought traces by autoregressively sampling tokens from a finite vocabulary. While this discrete sampling has achieved remarkable success, 015 conducting chain-of-thought with continuouslyvalued tokens (CoT2) offers a richer and more 018 expressive alternative. Our work examines the benefits of CoT2 through logical reasoning tasks 020 that inherently require search capabilities and provide optimization and exploration methods for CoT2. Theoretically, we show that CoT2 allows the model to track multiple traces in parallel and quantify its benefits for inference efficiency. No-025 tably, one layer transformer equipped with CoT2 can provably solve the combinatorial "subset sum problem" given sufficient embedding dimension. 028 These insights lead to a novel and effective su-029 pervision strategy where we match the softmax 030 outputs to the empirical token distributions of a set of target traces. Complementing this, we introduce sampling strategies that unlock policy optimization and self-improvement for CoT2. Our first strategy samples and composes K discrete 034 035 tokens at each decoding step to control the level of parallelism, and reduces to standard CoT when K = 1. Our second strategy relies on continuous exploration over the probability simplex. Experi-039 ments confirm that policy optimization with CoT2 indeed improves the performance of the model beyond its initial discrete or continuous supervision. 041

1. Introduction

043

044

045

046

047

049

050

051

052

053

054

000 001

002 003

008 009 010

> Chain-of-thought (CoT) strategies (Wei et al., 2022), when paired with strong base models, have achieved immense success and facilitated progress in remarkably challenging

tasks, such as solving AIME or IOI problems (Guo et al., 2025; Jaech et al., 2024). In essence, CoT boosts the expressive capability of the base model through autoregressive generation, a principle that also underlies the recent efforts on test-time compute scaling (Snell et al., 2024). Despite these advances, modern language model architectures may not fully utilize their potential for a few reasons. First is their discrete sampling of tokens-selecting a single token at each decoding step from a vocabulary of v tokens. This limits the model to emitting at most $\log_2(v)$ bits per sample, or more specifically, the Shannon entropy of the softmax output. This contrasts with the O(d) bits each token embedding can store, where d is the embedding dimension. Secondly, discrete sampling can cause the model to *commit* to certain solutions and avoid exploring alternatives (Yao et al., 2023). A practical method to address this is sampling multiple CoT traces and aggregating them, either through consistency (Wang et al., 2022) or best-of-N decoding (Ouyang et al., 2022) through more test-time computation.

In this work, we propose and investigate the use of *CoT* with <u>Continuous Tokens</u> (CoT2) to address these challenges, building on COCONUT (Hao et al., 2024). The fundamental idea in our CoT2 proposal is that rather than the model sampling a single token from the vocabulary, it samples or deterministically selects a continuous superposition of tokens according to the softmax output. Intuitively, this capability—effectively selecting multiple tokens simultaneously through a continuous superposition—would allow the model to pack more information within each token embedding and also enable it to track multiple reasoning paths in parallel—potentially emulating self-consistency or best-of-N decoding with a single trace. Toward this vision, we make the following technical contributions:

• Mechanistic and theoretical study of CoT2: We quantify the benefits of CoT2 along two directions. First, we examine the problem of *Minimum Non-Negative Sum* (MNNS) as a generalization of the classical Subset Sum problem. These problems, as well as related tasks like ProntoQA (Saparov & He, 2022), inherently benefit from parallel search capability. We show that a single layer transformer can solve MNNS using CoT2, showcasing the capability of transformers to track and expand multi-

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.





Figure 1: Illustration of CoT2 and discrete CoT for Minimum Non-Negative Sum (MNNS) task with m = 3. The input numbers are 2, 1, 4, and the correct path for this task (-2, -3, 1) is highlighted with yellow arrows and corresponds to the discrete CoT supervision. CoT2 supervision for the reasoning steps $t \in \{1, ..., m-1\}$ is the average of embeddings of reachable states, and for t = m is the embedding corresponding to the answer.

ple reasoning traces in latent space. Complementing this, under a certain trajectory decoupling assumption, we provide a theoretical study of CoT2 decoding methods

- Base CoT2: deterministic inference which creates and feeds continuous tokens using full softmax output at each step (Sec. 2&3);
- CoT2-MTS (multi-token sampling): our method which samples K discrete tokens from softmax and averages them to form a continuous token (Sec. 4);

and standard CoT which is a special case of CoT2-MTS with K = 1. We show that base CoT2 tracks and aggregates all reasoning paths whereas CoT2-MTS strictly generalizes CoT by tracking K paths; and establish the sample complexity benefits of the CoT2 methods.

Supervision and reinforcement for CoT2: We introduce the continuous supervision strategy CSFT for CoT2 models to explicitly track multiple teacher traces in parallel by fitting a target softmax map of the empirical distribution of tokens within the trace. Our method also reveals fundamental tradeoffs between the CoT2 accuracy and the embedding dimension. Complementing this, we introduce policy optimization methods for CoT2 (Section 4). We propose MTS as our primary strategy, which samples and composes K discrete tokens at each forward pass to control the level of parallelism. We also introduce a purely continuous sampling scheme over the probability simplex. Experiments on the MNNS, ProntoQA, and ProsQA tasks demonstrate that GRPO-based RL with CoT2 further improves the accuracy over SFT or CSFT (see Section 4.3). This demonstrates that the RL phase helps the model better prioritize relevant reasoning traces

and unlocks a promising strategy for training CoT2-based language models.

Ultimately, our results and methods underscore the strong potential of CoT2 and encourage further research. The rest of the paper is organized as follows: Section 2 introduces the technical setup, Section 3 describes our continuous supervision strategy as well as the MNNS, ProntoQA, and ProsQA tasks. Section 4 describes our sampling strategies and the resulting GRPO-based policy optimization methods. Section 5 provides theoretical guarantees and Section 6 concludes with a discussion.

1.1. Related Work

The efficacy of eliciting reasoning in LLMs through chainof-thought (CoT) prompting has been well-established (Nye et al., 2021; Wei et al., 2022; Kojima et al., 2022; Suzgun et al., 2023; Guo et al., 2025). CoT prompting provides a convenient way to increase inference-time compute and computational depth, both of which have been found to be independently useful (Pfau et al., 2024; Goyal et al., 2024; Feng et al., 2023; Merrill & Sabharwal, 2024). However, the discrete nature of CoT tokens forces sequential exploration of reasoning paths, resulting in longer reasoning paths and consequently increased inference-time compute. Furthermore, restricting reasoning to natural language can be inefficient, as groups of tokens can often be more effectively represented by a single continuous token. Thus, CoT2 offers an alternative strategy for compute-efficient reasoning and complements methods that aim to shorten/control the trace length of CoT (Aggarwal & Welleck, 2025; Zhang et al., 2025; Sui et al., 2025).

One way to address these challenges is by leveraging the 111 implicit reasoning capabilities of transformers (Yang et al., 112 2024; Shalev et al., 2024). Works such as (Deng et al., 2023; 113 2024; Yu et al., 2024) use various techniques to obtain mod-114 els that can perform reasoning internally without emitting 115 CoT tokens. Another line of work has found looped trans-116 formers to be effective on reasoning problems (Giannou 117 et al., 2023; Geiping et al., 2025), notably being able to 118 mimic CoT (Saunshi et al., 2025) with a sufficient number 119 of iterations. Our work is similar to this line of work in that

120 continuous representations are used to perform reasoning.

121 Our work is most related to a recent body of work intro-122 ducing LLMs capable of reasoning with explicit continuous 123 tokens decoded autoregressively. In particular, recently pro-124 posed COCONUT (Hao et al., 2024) autoregressively feeds 125 the last token's final-layer representation as input to the next 126 step. Given labeled CoT data, COCONUT is trained to pro-127 gressively replace discrete tokens with continuous tokens 128 (from left to right). Shen et al. (2025) propose CODI, where 129 an LLM with continuous CoT is supervised to produce the 130 correct answer, while also aligning its hidden representation 131 on the last reasoning token to that of a discrete CoT model 132 that shares the same backbone. Cheng & Van Durme (2024) 133 propose CCOT, where an auxiliary module is first trained 134 to decode autoregressively a compressed representation of 135 a discrete CoT trace, and later the main LLM is fine-tuned 136 to produce correct answers by additionally conditioning on 137 the generated continuous tokens. While COCONUT, CODI, 138 CCOT, and our CoT2 all aim to reason in continuous space, 139 we propose distinct algorithmic approaches that also address 140 the exploration challenge. Key differences include: (1) Our 141 continuous tokens are simplex-weighted compositions of 142 vocabulary tokens. (2) Our supervision method is novel and 143 explicitly targets implicit parallelism. (3) CoT2 does not 144 initialize from, nor attempt to mimic, discrete CoT. (4) By 145 introducing sampling strategies and associated GRPO variations, we realize the "Supervised Training \rightarrow Reinforcement 147 Learning" paradigm in the context of CoT2. We provide 148 further discussion of literature on multi-token prediction 149 and reinforcement learning in Appendix A. 150

2. Problem Setup

153 **Notation.** For an integer $n \ge 1$, we use the shorthand 154 $[n] = \{1, ..., n\}$. Throughout, we denote vectors by bold 155 lowercase letters (e.g. x) and matrices by bold uppercase 156 letters (e.g. X). For a vector $x \in \mathbb{R}^n$, the component x_i refers 157 to its *i*-th entry. The zero vector in \mathbb{R}^n is written as $\mathbf{0}_n$, and 158 the zero matrix in $\mathbb{R}^{m \times n}$ is $\mathbf{0}_{m \times n}$. Finally, we let $\Delta^{\nu-1}$ denote 159 the standard $\nu - 1$ simplex in \mathbb{R}^{ν} .

Assume that we are given an input context $X \in \mathbb{R}^{n \times d}$, where each of the *n* rows is a *d*-dimensional embedding vector. Our objective is to output *m* tokens given the context *X* with

164

151

152

*m*th output token being the final answer that is evaluated under some performance metric (e.g. accuracy or reward). For the first m-1 steps, the model outputs *continuous* tokens $\{z_t\}_{t \in [m-1]}$, which are *thought tokens* that enable a reasoning process. At the final step t = m, the model outputs a *discrete* token z_m from a vocabulary of size v. In the remainder of this paper, we investigate strategies for training this system in a way that improves final performance over standard discrete next-token prediction.

Formally, let $E = [e_1, \ldots, e_v]^\top \in \mathbb{R}^{v \times d}$ be the embedding matrix corresponding to the vocabulary of v tokens, where $e_i \in \mathbb{R}^d$ represents the embedding of the *i*th token. We define the next-token prediction model LM_θ parameterized by θ that assigns, at each step t, a probability distribution over possible next tokens given the prefix $z_{\leq t}$ and context X. Concretely, for $1 \leq t \leq m - 1$, the model outputs the following probability distribution over the v vocabulary entries via a softmax operation:

$$LM_{\theta}(\cdot \mid z_{< t}, X) := \alpha_t \text{ where}$$

$$\alpha_t = \left[\alpha_{t,1}, \ldots, \alpha_{t,v}\right] \in \Delta^{v-1}, \text{ i.e. } \alpha_{t,i} \ge 0 \text{ and } \sum_{i=1}^{v} \alpha_{t,i} = 1$$

We then form the continuous token as the convex combination of all tokens in the vocabulary:

$$z_t = E^{\top} \alpha_t \in \mathbb{R}^d, \quad \forall 1 \le t \le m-1$$

Hence each continuous token z_t is a linear combination of the vocabulary embeddings. At the final step t = m, the model samples a discrete token $z_m \in \{e_1, \ldots, e_v\}$ from its policy distribution $LM_\theta (\cdot | z_{\leq m}, X) = \alpha_m$. Finally, we note that we assume that the answer depends only on the final discrete token z_m merely for simplicity; the same framework naturally extends to decoding multiple final discrete tokens after continuous ones. We refer to this decoding strategy as **base CoT2** and observe that it results in a deterministic reasoning chain because the continuous tokens are precisely determined by the softmax map. In Section 4, we will introduce stochastic alternatives, such as **CoT2-MTS**, to facilitate generative reasoning.

3. CSFT: A Supervised Training Method for CoT2

In this section, we present our method of continuous supervised training to learn intermediate thought tokens as "soft" targets rather than "hard" target tokens, as described in Section 2. Specifically, we provide the model with convex combinations of vocabulary embeddings, which allows the model flexibility in those reasoning steps. Such an approach is particularly suitable when the task accuracy depends only on the final token or token distribution. Formally, at each reasoning step t = 1, ..., m - 1, the supervision specifies a



(a) Continuous vs. discrete SFT (Pass@k) (b) Training performance vs. embedding di- (c) Training performance vs. embedding diaccuracies for different temperatures on mension for continuous and discrete SFT on mension for continuous and discrete models MNNS task. on ProsQA.

Figure 2: The teacher distribution for the continuous model is derived from a search algorithm. (a): The figure illustrates that the discrete model requires multiple samplings (Pass@k) to match the single-sample performance of the continuous model on MNNS (10-run average). **Setting:** 4 input digits in 1 - 9; 1-layer, 1-head GPT2 with d = 24. (b-c): The figures reveal that above a certain embedding dimension threshold, the continuous model is superior to discrete in tasks involving search, like MNNS and ProsQA. **Setting (b):** 4 input digits in 1 - 9; 2-layer, 2-head GPT2 with $d \in \{16, 24, 32\}$. (c): 4-layer, 4-head GPT2 with $d \in \{24, 32, 40\}$.

target probability distribution

174

175

176

178

179

180

181 182

183

184

189 190

204

206

208

209

210

211

212

$$\boldsymbol{\alpha}_t^* = \left[\alpha_{t,1}^*, \ldots, \alpha_{t,v}^*\right] \in \Delta^{v-1},$$

185 where $\alpha_{t,i}^* \ge 0$ and $\sum_{i=1}^{\nu} \alpha_{t,i}^* = 1$. We train the model to align 186 its predicted distribution α_t to the supervision distribution 187 α_t^* rather than one-hot labels, with the help of a divergence-188 based loss:

$$\mathcal{L}_{\text{cont}}(\theta; X, t) = D\left(\boldsymbol{\alpha}_{t}^{*} \| \boldsymbol{\alpha}_{t}\right)$$

191 where $D(\cdot \| \cdot)$ is the cross-entropy (or equivalently KL divergence) between two distributions. This approach can 193 also be viewed as token-level knowledge-distillation, where 194 the teacher distribution α_t^* may be obtained through a 195 logic/search algorithm. At the final step t = m, we typ-196 ically have a discrete target $z_m^* \in \{e_1, \ldots, e_v\}$, so that α_m^* 197 is one-hot distribution placing probability 1 on that target 198 token and 0 elsewhere. This is equivalent to employing a 199 standard cross-entropy loss $-\log LM_{\theta}(z_m^* \mid z_{\leq m}, X)$ at the 200 final step. Hence, for each training example, the total loss 201 for the proposed continuous supervised training is the sum 202 of the continuous-token divergence losses: 203

$$\mathcal{L}_{\text{CSFT}}(\theta; X) = \sum_{t=1}^{m} \mathcal{L}_{\text{cont}}(\theta; X, t)$$
(1)

By minimizing $\mathcal{L}_{\text{CSFT}}(\theta)$, we teach the model to learn the soft targets α_t^* at each step and to predict the correct final discrete token. In the above training procedure, inspired by the discussions in Bachmann & Nagarajan (2024); Bengio et al. (2015), we consider two ways of providing prefixes to the language model:

1. **Teacher forcing:** Each step *t* is conditioned on the ground-truth prefix $z_{<t}^*$, meaning that the model has access to all ground-truth previous tokens during prediction. Concretely, for each step t' < t, the corresponding input $z_{t'}^* = E^{\top} \alpha_{t'}^*$ is a convex combination of all vocabulary tokens. 2. **Self-feeding:** Each step *t* autoregressively uses the model's previously generated outputs, $z_{<t}$, during training. In particular, as described in Section 2, the continuous output token $z_t = E^{\top} \alpha_t$, is a convex combination of vocabulary embeddings, which is then fed back to the model as part of the prefix.

It is also worth noting that one may apply *temperature scaling* or *thresholding* to α_t before forming z_t in order to filter the model's predictions. In our experiments, we find that teacher forcing leads to superior performance for CSFT, even though at inference time, the model runs in an autoregressive manner, as discussed below. See Appendix C for further discussion.

Inference. At inference time, the model does not rely on the ground-truth distributions α_t^* . Instead, at each continuous step t < m, the model autoregressively uses its own output distribution α_t by converting that distribution to a continuous token $z_t = E^{\top}\alpha_t$ and adds it to the prefix for the prediction in the next step. At the final step, the model generates a discrete sample from $\alpha_m = \text{LM}_{\theta}(\cdot \mid z_{< m}, X)$.

Discrete baseline. In this case, we use teacher-forced training where the next token prediction is performed conditioned on the previous ground-truth tokens with standard cross-entropy loss. Discrete baseline enforces z_t^* to be a token in vocabulary $\{e_1, \ldots, e_{\nu}\}$, which means that it is a special case of CSFT where the α_t^* are one-hot vectors rather than an arbitrary element of $\Delta^{\nu-1}$. The model minimizes the following objective, which is obtained by summing over all steps of teacher-forced next-token prediction:

$$\mathcal{L}_{\text{SFT}}(\theta; X) = \sum_{t=1}^{m} -\log \text{LM}_{\theta}\left(z_t^* \mid z_{< t}^*, X\right).$$
(2)

3.1. Tasks Requiring Exploration over States

In the next two subsections, we illustrate the CSFT training described in (1) on tasks that require *exploration* over multiple states. Consider that the vocabulary is sufficiently large that each state g of the task can be assigned a unique embedding. Then, we let Γ_t be the set of all states that could result from building upon step (t - 1), where $\Gamma_0 = \{g_0\}$ for the initial state g_0 . For each element $g \in \Gamma_t$, we assign a probability $\alpha_{t,g}^*$ that reflects how many times that state occurs under a search process. Then, α_t^* is formed by normalizing these probabilities into a distribution on Γ_t :

$$\alpha_{t,g}^* = \frac{\operatorname{count}_t(g)}{\sum_{h \in \Gamma} \operatorname{count}_t(h)},$$
(3)

where count_t(g) is the number of times state g appears among all expansions at step t. At the final step t = m, we select exactly one correct state from Γ_m , so that α_m^* is a one-hot vector:

$$\alpha_{m,g}^* = \begin{cases} 1, & \text{if } g \text{ is the correct final state,} \\ 0, & \text{otherwise.} \end{cases}$$

Remark. While we focus on search-based tasks MNNS and ProntoQA in the next two sections, one can also extend to training with continuous tokens in the language-model context. The distributions $\{\alpha_t^*\}$ at each step can be collected by (1) running a beam or best-first search to generate multiple partial trajectories; (2) scoring these trajectories with a reward function; and (3) curating them into a distribution that assigns higher mass to states that leads to higher rewards. This construction of α_t^* replaces one-hot supervision with soft supervision for intermediate reasoning steps.

3.1.1. MINIMUM NON-NEGATIVE SUM TASK

We now introduce the Minimum Non-Negative Sum (MNNS) task, where the goal is to assign signs to a list of numbers so that their sum is as small as possible while being nonnegative. The MNNS task can be viewed as partitioning a set of numbers into two subsets with a minimal difference, which makes it closely related to the subset sum problems explored in Dziri et al. (2023); Thomm et al. (2024). Formally, we are given *m* integers d_1, \ldots, d_m , and the task is to assign signs $\sigma_i \in \{+1, -1\}$ such that $s = \sigma_1 d_1 + \dots + \sigma_m d_m \ge 0$ and s is minimized. Let $\sigma^{\text{opt}} = (\sigma_1^{\text{opt}}, \dots, \sigma_m^{\text{opt}})$ denote the optimal assignment that achieves the minimal nonnegative sum s^{opt} out of 2^m possible sign assignments. Here, every possible partial sum $\sigma_1 d_1 + \cdots + \sigma_t d_t \in \Gamma_t$ is assigned a unique embedding $e_{\phi(\sigma_1 d_1 + \dots + \sigma_t d_t)}$, where $\phi(\cdot)$ maps each sum to a distinct id in [v]. We now describe the two modes of supervision, where the input digits are processed one by one by accumulating partial sums, as illustrated in Figure 1:

• Supervision for CoT2 model: At step t, there are $|\Gamma_t| = 2^t$

partial sums of length *t*, and accordingly, we provide the following target distribution α_t^* :

$$\alpha_{t,i}^* = \begin{cases} \frac{\operatorname{count}_t(i)}{2^t}, & \text{if token } i \text{ appears } \operatorname{count}_t(i) \text{ times} \\ & \text{as a partial sum of length } t; \\ 0, & \text{otherwise.} \end{cases}$$

At the final step t = m, the distribution α_m^* assigns probability 1 to the correct sum $e_{\phi(\sigma_1^{\text{opt}}d_1+\dots+\sigma_m^{\text{opt}}d_m)}$ and 0 to all others.

• Supervision for discrete model: We supervise the discrete model along the correct chain of partial sums by providing $e_{\phi(\sigma_1^{\text{opt}}d_1+\dots+\sigma_t^{\text{opt}}d_t)}$ for $1 \le t \le m$ as target tokens, and train following the standard cross-entropy objective described in (2).

While constructing the dataset, we split the training and validation sets by ensuring that any permutation of numbers appears in exactly one split. The aim behind this is to prevent memorization and make a fair evaluation. We also encode input and output numbers with separate tokens in our vocabulary. As an example, an input appears as $\langle BOS \rangle d_1 d_2 \dots \rightarrow$, and the corresponding output as $s_1 s_2 \dots s^{opt} \langle EOS \rangle$, where s^{opt} is the minimal nonnegative sum for $\{d_1, \dots, d_m\}$. For the model, we use the GPT2 architecture (Radford et al., 2019) with different head, layer, and embedding dimension configurations, and train it from scratch. Please refer to Appendix B for more experimental details.

3.1.2. PRONTOQA AND PROSQA DATASETS

Other datasets we explore in our investigation of the CSFT approach are the ProntoQA (Saparov & He, 2022) and ProsQA (Hao et al., 2024) datasets, which are logical reasoning tasks that require exploration over multiple possible paths. Specifically, each question in ProntoQA asks whether a certain target word (node) B is reachable from a root word (node) A within a fixed number of hops, while for ProsQA it asks which of the target words B or C is reachable. We use 5-hop questions and present the graph in a structured format. In particular, for each problem, we represent nodes and edges using embeddings, which we use as the model input rather than text input. The detailed structured format and examples are provided in Appendix B.2.

The graph structure of the ProntoQA and ProsQA datasets naturally obeys the supervision in (3), so that we determine the words that can be reached using t edges from A and supervise intermediate tokens on the resulting distribution. At the final reasoning step m, the supervision assigns probability 1 to the correct label: yes or no for ProntoQA, and B or C for ProsQA. For the standard discrete model, we provide an explicit chain of nodes from A to the target node

274

275 (B or C) as the target path at each step. We direct readers to 276 Appendix **B**.2 for additional details on the supervision.

278 3.2. Results and Discussion of CoT2 Supervision

277

311

312

279 In our experiments on the MNNS, ProsQA, and ProntoQA 280 tasks, we observe that CSFT significantly outperforms dis-281 crete baseline once the embedding dimension exceeds a 282 moderate threshold, as shown in Figures 2b and 2c. We 283 also note that continuous tokens allow for a faster conver-284 gence above this threshold, as illustrated by Figures 2b 285 and 2c. Through continuous tokens, the model gains the 286 capacity to represent multiple partial expansions in parallel, 287 enabling a 'search-like' ability that results in higher accu-288 racy. We also demonstrate in Figure 2a that the discrete 289 model requires multiple sampling (pass@k) to approach 290 the performance achieved by a single attempt of the CoT2 291 model. This indicates that continuous tokens effectively 292 hedge against early mistakes, as there's no accumulating 293 error. This argument aligns with the previous "snowballing 294 errors" discussions for discrete autoregressive generation 295 in (Bachmann & Nagarajan, 2024). Moreover, although 296 the continuous approach needs more embedding capacity 297 to allow its distributional representations at each step, it can then achieve strong performance with fewer layers and 299 heads compared to the discrete model, as further demon-300 strated by the results in Appendix C.1. We also provide 301 experiments on ProntoQA in Appendix C.1, which confirm 302 similar findings to results on ProsQA. 303

304 We also evaluated the performance of the discrete model under more sparse supervision scenarios, such as providing 306 only a subset of the correct partial sums or, more extremely, 307 only the final answer. In these experiments, we observed 308 that denser supervision improved the discrete CoT model's 309 performance; see Appendix C.1 for further details. 310

4. Reinforcement Learning Methods for CoT2

313 In this section, we describe how to apply RL with continuous 314 output tokens. Specifically, we explore GRPO training on 315 top of continuous or discrete models that are supervised 316 trained based on Section 3 for the MNNS, ProntoQA, and 317 ProsQA tasks. By illustrating two sampling methods for 318 GRPO, we demonstrate how a model trained with discrete 319 SFT can be adapted to produce continuous outputs. We 320 assume a sparse reward setting where the reward is 1 for a 321 correct final answer and 0 otherwise.

In our setup, a language model LM_{θ} acts as a policy 323 over tokens. Let $\{\mathbf{Z}^{(i)}\}_{i=1}^{G}$ be a group of G trajectories 324 sampled from old policy $LM_{\theta_{old}}$ such that each trajectory $\mathbf{Z}^{(i)} = (\mathbf{z}_1^{(i)}, \dots, \mathbf{z}_m^{(i)})$ contains *m* output tokens given a fixed 325 326 327 input X. We denote by $\hat{A}_{i,t}$ the advantage estimate at step t in trajectory i and note that $\hat{A}_{i,t} = \hat{A}_i$ is identical across 328 329

Algorithm 1 Multi-Token Sampling GRPO for Continuous **Token Generation**

Input: Initial policy $LM_{\theta_{init}}$; hyperparameters K, G, m, ϵ, β . 1: $LM_{\theta}, LM_{\theta_{ref}} \leftarrow LM_{\theta_{init}}$

- 2: for iteration = 1, 2, ..., I and for step = 1, 2, ..., S do
- Sample a batch of inputs $\{X^{(b)}\}_{h=1}^{B}$ 3:
- 4: Update $LM_{\theta_{old}} \leftarrow LM_{\theta}$
- for each input X in the batch and for each trajectory i =5: 1, \ldots , *G* from that **X** do
- 6: for each token step $t = 1, \ldots, m$ do
- if t < m then 7: Sample *K* tokens $\{e_{i_1}, \ldots, e_{i_K}\}$ from $\alpha_t^{(i), \text{old}}$ to create continuous token $z_t \leftarrow \frac{1}{K} \sum_{r=1}^{K} e_{i_r}$. 8:
 - Policy ratio $r_t(\theta) \leftarrow \left(\prod_{r=1}^K \alpha_{t,i_r}^{(i)} / \prod_{r=1}^K \alpha_{t,i_r}^{(i), \text{old}}\right)^{\frac{1}{K}}$.

10: else 11:

9:

- se Sample $z_m = e_j$ from $\alpha_m^{(i), \text{ old}}$.
- Policy ratio for discrete token $r_m(\theta) \leftarrow \alpha_{m,i}^{(i)} / \alpha_{m,j}^{(i), \text{old}}$. 12: end if
- 13: 14: end for
- Obtain advantage estimates $\hat{A}_{i,t}$ for each token t in each 15: trajectory $\mathbf{Z}^{(i)}$ and calculate objective.
- 16: end for
- 17: Update θ to minimize $\mathcal{L}_{\text{GRPO}}(\theta)$.

18: end for

```
Output: LM_{\theta}
```

all steps of a trajectory under sparse reward setting. To quantify how the new policy LM_{θ} differs from the old one on token $z_t^{(i)}$ from *i*th trajectory, we define the policy ratio $r_t^{(i)}(\theta) = \frac{\sum_{l \in \mathcal{M}_{\theta}(z_t^{(i)} | z_{< t}^{(i)}, X)}}{\sum_{l \in \mathcal{M}_{\theta_{old}}(z_t^{(i)} | z_{< t}^{(i)}, X)}}.$ We update the model parameters θ by minimizing the clipped surrogate objective (Shao et al., 2024: Yu et al., 2025):

$$\begin{split} \mathcal{L}_{\text{GRPO}}(\theta) &= -\frac{1}{\sum_{i=1}^{G} |\mathbf{Z}^{(i)}|} \sum_{i=1}^{G} \sum_{t=1}^{|\mathbf{Z}^{(i)}|} \left[\min\left(r_t^{(i)}(\theta) \, \hat{A}_{i,t}, \right. \\ & \left. \operatorname{clip}\left(r_t^{(i)}(\theta), 1\!-\!\epsilon, 1\!+\!\epsilon\right) \, \hat{A}_{i,t}\right) - \beta \, \mathbb{D}_{\text{KL}}\left[\mathrm{LM}_{\theta} || \mathrm{LM}_{\theta_{\text{ref}}} \right] \right] \end{split}$$

As the output length is fixed in our setting, we have $|\mathbf{Z}^{(i)}| =$ *m* for each trajectory. Here, ϵ is a clipping parameter that bounds the ratio $r_t(\theta)$, and β controls the strength of KLdivergence from a reference policy $LM_{\theta_{ref}}$ which is the SFTinitialized policy. We set the number of GRPO iterations $\mu = 1$ and estimate the KL divergence with the Schulman Approximator as in Shao et al. (2024).

4.1. Multi-Token Sampling

We emulate the rollout of a continuous token by sampling a fixed number of $K \le v$ discrete tokens and averaging them at steps $t = 1, \ldots, m - 1$. We refer to this hybrid method as CoT2-MTS (multi-token sampling). For the GRPO objective, we propose the following method to calculate the policy ratio for continuous tokens. Specifically, assume at step t we sample discrete tokens e_{i_1}, \ldots, e_{i_k} with probabilities $\alpha_{t,i_1}, \ldots, \alpha_{t,i_K}$ under the current policy and probabilities $\alpha_{t,i_1}^{\text{old}}, \dots, \alpha_{t,i_K}^{\text{old}}$ under the old policy. We define the policy ra-

	K	Val. Accuracy (%)		Val. Entropy (SFT \rightarrow SFT+GRPO)					
п		SFT	SFT+GRPO	token1	token ₂	token3	token ₄		
24	1 3 6	39.76	49.01 52.60 49.69	$\begin{array}{c} 0.3218 \rightarrow 0.0314 \\ 0.3717 \rightarrow 0.0647 \\ 0.4524 \rightarrow 0.1242 \end{array}$	$\begin{array}{c} 0.5858 \rightarrow 0.0712 \\ 0.7461 \rightarrow 0.2120 \\ 0.7738 \rightarrow 0.3361 \end{array}$	$0.5499 \rightarrow 0.1576$ $0.8006 \rightarrow 0.3312$ $0.8364 \rightarrow 0.6615$	$\begin{array}{c} 0.4786 \rightarrow 0.1718 \\ 0.5338 \rightarrow 0.1529 \\ 0.5134 \rightarrow 0.2159 \end{array}$		
32	1 3 6	43.50	51.61 55.66 50.38	$\begin{array}{c} 0.3618 \rightarrow 0.0143 \\ 0.3886 \rightarrow 0.0412 \\ 0.4224 \rightarrow 0.0632 \end{array}$	$\begin{array}{c} 0.6331 \rightarrow 0.0395 \\ 0.7101 \rightarrow 0.0904 \\ 0.7915 \rightarrow 0.2161 \end{array}$	$\begin{array}{c} 0.3518 \rightarrow 0.2631 \\ 0.5447 \rightarrow 0.5757 \\ 0.6077 \rightarrow 0.8481 \end{array}$	$\begin{array}{c} 0.1962 \rightarrow 0.1182 \\ 0.2863 \rightarrow 0.1734 \\ 0.2780 \rightarrow 0.1514 \end{array}$		

tio for these continuous steps by dividing geometric means:

333 334 335

345

346

347

349

350

351

352

353

354

355

356

357

358

359

360

361

362

363 364

365

366

367

368

369

370

371

373

374

375

376

$$r_t(\theta) = \frac{\mathrm{LM}_{\theta}\left(z_t \mid z_{< t}, X\right)}{\mathrm{LM}_{\theta_{\mathrm{old}}}\left(z_t \mid z_{< t}, X\right)} = \left(\frac{\alpha_{t, i_1} \cdots \alpha_{t, i_K}}{\alpha_{t, i_1}^{\mathrm{old}} \cdots \alpha_{t, i_K}^{\mathrm{old}}}\right)^{1/K}, \quad (4)$$

for t = 1, ..., m - 1. The geometric mean ensures that the ratio for each continuous step remains on the same scale as the final discrete token's ratio and, thus, helps avoid overly large or small updates in the GRPO objective and provides more stable training compared to the direct multiplication of probabilities. Once this ratio is computed, we then average the *K* sampled tokens to form z_t , which is fed to the model as the query for the next prediction step. At the final step t = m, where the token $z_m = e_j$ is discrete with $j \in [v]$ denoting its index, the policy ratio is simply the probability ratio of selecting that token:

$$r_m(\theta) = \frac{\mathrm{LM}_{\theta}\left(z_m \mid z_{< m}, X\right)}{\mathrm{LM}_{\theta_{\mathrm{old}}}\left(z_m \mid z_{< m}, X\right)} = \frac{\alpha_{m,j}}{\alpha_{m,j}^{\mathrm{old}}}.$$
(5)

Inference. During inference after GRPO training, we apply the same multi-token sampling procedure at each of the first m - 1 steps to form the continuous token via the average of *K* sampled embeddings.

Remark. An alternative to the normalization of ratios given by (4) is to directly scale down the logits by 1/K before applying softmax. However, using this approach during inference leads to a distribution shift relative to the SFTtrained model and ultimately degrades performance.

4.2. Dirichlet Sampling

In this section, we present another method for generating continuous tokens at each step by interpreting the model's output distribution $\alpha_t \in \Delta^{\nu-1}$ as concentration parameters of a Dirichlet distribution over the $\nu - 1$ simplex. We introduce a scaling hyperparameter $\gamma > 0$ and define the Dirichlet distribution with the parameters $\gamma \alpha_t = (\gamma \alpha_{t,1}, \dots, \gamma \alpha_{t,\nu})$. Without this scaling, directly using α_t as parameters often causes training instability, particularly when many $\alpha_{t,i}$ values are small. We then sample a point $\hat{\alpha}_t \in \Delta^{\nu-1}$ from the resulting distribution Dir (α_t) . After sampling, we form the continuous token by mapping $z_t = E^{\top}\hat{\alpha}_t \in \mathbb{R}^d$, which becomes the query for the next step. We denote the Dirichlet densities induced by current and old policies as $f_{\theta}(z_t)$ and $f_{\theta_{\text{old}}}(z_t)$, respectively. Accordingly, we define the policy ratio at a continuous step t < m as:

$$r_t(\theta) = \frac{\mathrm{LM}_{\theta}(z_t \mid z_{< t}, X)}{\mathrm{LM}_{\theta_{\mathrm{old}}}(z_t \mid z_{< t}, X)} = \frac{f_{\theta}(z_t)}{f_{\theta_{\mathrm{old}}}(z_t)}$$

The above definition parallels how we compute probability ratios for discrete actions but replace the categorical pmf with continuous Dirichlet pdf. At the final step t = m, we sample a discrete token $z_m \in \{e_1, \ldots, e_v\}$ from α_m , and use the standard policy ratio given by (5). At inference, we follow the autoregressive procedure in Section 3 by creating a convex combination of vocabulary tokens.

4.3. Results and Discussion of Policy Optimization for CoT2

MNNS evaluation: Table 1 provides our results for the MNNS task and demonstrates that, for each $K \in \{1, 3, 6\}$, CoT2-MTS significantly improves validation accuracy relative to the discrete SFT baseline (39.76%), with moderate K yielding the best final performance. We also observe that smaller K-values correspond to larger reductions in tokenlevel entropies, suggesting that the model becomes more confident at each intermediate step by learning to commit to fewer tokens. This suggests a curriculum on K-starting small and gradually increasing-could potentially further improve the training on the MNNS task. Interestingly, the third token's entropy remains relatively high, which might indicate that the model continues to hedge among several partial expansions at that step, which may help preserve useful diversity of reasoning. Therefore, CoT2-MTS enables a model trained with discrete SFT to produce continuous outputs and helps it achieve better performance. Finally, Ap385Table 2: Validation accuracies on ProsQA and ProntoQA386for CoT2 and Discrete CoT, evaluated at K = 6, 8, 12 us-387ing CoT2-MTS sampling GRPO. All models use a 4-layer,3884-head GPT2 with embedding dimension 32. The SFT val-389ues are constant as they represent initial accuracies before390GRPO. Remarkably, GRPO with our multi-token sampling391scheme results in consistent improvements.

392

395

396

399

400 401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

]	ProsQA	ProntoQA		
		SFT	SFT+GRPO	SFT	SFT+GRPO	
K=6	CoT2	93.37	93.83	75.36	76.15	
	Discrete CoT	68.50	68.24	59.58	62.28	
K=8	CoT2	93.37	94.09	75.36	76.66	
	Discrete CoT	68.50	71.58	59.58	71.53	
K=12	CoT2	93.37	94.21	75.36	77.64	
	Discrete CoT	68.50	72.76	59.58	74.03	

pendix C.2 shows that the CoT2 model with CSFT achieves a strong performance once the embedding dimension is sufficiently large (compared to the results in Table 1), however, it can be further improved with GRPO with Dirichlet Sampling.

ProsQA and ProntoQA evaluation: Table 2 provides an evaluation of the benefits of GRPO with CoT2-MTS on models trained with discrete or continuous SFT. Remarkably, we observe that both CoT2 and discrete CoT models consistently improve across all rollout sizes (K = 6, 8, 12), with larger K values yielding better final accuracies by promoting more exploration. Notably, the discrete CoT model benefits relatively more from reinforcement learning compared to CoT2, likely because the CoT2 model already incorporates exploration implicitly through continuous supervision (CSFT). Aligning with this, we observe that for the ProntoQA dataset, the final performance of discrete CoT approaches that of the CoT2 model. Finally, we expect that, the performance gain of the MNNS task could be more limited compared to ProsOA and ProntoOA due to the highly structured nature of the MNNS task which makes CSFT supervision a natural choice and difficult to improve over.

5. Theoretical Analysis

In this section, we first present the construction of a singlelayer transformer that solves the MNNS task using an attention layer followed by a mixture-of-experts MLP layer. We then provide a theoretical comparison between base CoT2, CoT2-MTS, and discrete CoT models.

5.1. Solving the Minimum Non-Negative Sum Task

Proposition 1 (Solving MNNS). There exists a 1-layer transformer architecture with a mixture-of-experts MLP layer that solves the MNNS task using CoT2 by storing (sine,

cosine) embeddings of all 2^k states at the k-th iteration in a non-overlapping manner.

The above construction utilizes trigonometric embeddings, inspired by the mechanistic insights given by Nanda et al. (2023). Our approach leverages these trigonometric embeddings to provide a theoretical guarantee that **the transformer can track and add/subtract multiple numbers in parallel** by benefiting from the embedding capacity and reading off the minimum non-negative number at the final step. An important observation regarding our construction is that the trajectories at each intermediate reasoning step are truly decoupled as it stores each state using non-overlapping (sine, cosine) representations. This is similar to the left side of Figure 1, but we also utilize rotations/shifts to ensure distinct states are orthogonal and are easy to read out.

5.2. Understanding and Formalizing the Benefits of CoT2 and Comparison to CoT

To proceed, we argue that CoT2 equips the model with the ability to track multiple paths in parallel which can be formalized through Assumption 1 below. Building on this condition, we will provide a formal comparison of CoT2 and discrete CoT models in the remainder of this section.

Assumption 1. Recall the model LM_{θ} in Section 2. For any step t and prefix tokens $z_{\leq t}$, we assume (i) the next token probabilities depend only on the last token z_t and the query X and (ii) if the last token is $z_t = \sum_{j=1}^{\nu} \alpha_{t,j} e_j$ so that $\sum_{j=1}^{\nu} \alpha_{t,j} = 1$, the output distribution α_{t+1} decouples as follows:

$$\mathrm{LM}_{\theta}^{(t)}(\cdot \mid \boldsymbol{z}_{t}, \boldsymbol{X}) \stackrel{(i)}{=} \mathrm{LM}_{\theta}^{(t)}(\cdot \mid \boldsymbol{z}_{\leq t}, \boldsymbol{X}) \stackrel{(ii)}{=} \sum_{j=1}^{\nu} \alpha_{t,j} \, \mathrm{LM}_{\theta}^{(t)}(\cdot \mid \boldsymbol{e}_{j}, \boldsymbol{X}).$$

Under Assumption 1, the token distribution $\alpha_{t+1} = LM_{\theta}^{(t)}(\cdot | z_t, X)$ evolves with the equation $\alpha_{t+1} = \alpha_t M_t(z_t; X)$, starting from $\alpha_1 = LM_{\theta}^{(0)}(\cdot | X)$ until α_m . Here, $M_t(z_t; X) \in \mathbb{R}^{\nu \times \nu}$ is a Markov transition matrix that is allowed to depend on the input X and the last token z_t (see (IIdiz et al., 2024) for related discussion). To keep exposition cleaner, we omit z_t and X in the notation of $M_t(z_t; X)$, and use M_t instead. We first observe the following regarding base CoT2, standard discrete CoT, and CoT2-MTS inference strategies.

• *Base CoT2*: At each step t = 1, ..., m, the model outputs the continuous token $z_t = E^{\top} \alpha_t$ and uses it as the query for the next step.

Interpretation: Base CoT2 keeps track of all possible traces simultaneously. Over *m* steps, it tracks and aggregates all v^m traces where the trace $(i_t)_{t=1}^m$ has a weight of $\prod_{i=1}^m \alpha_{t,i_i}$.

Discrete CoT: At each step 1 ≤ t ≤ m, the model samples exactly one token z_t = e_i, from α_t, and uses it as the query

440 for the next step.

- 441 **Interpretation:** Discrete CoT samples a single trace out 442 of v^m traces with a likelihood of $\prod_{i=1}^m \alpha_{t,i_t}$ for trace $(i_t)_{t=1}^m$.
- **CoT2-MTS (multi-token sampling)**: At each step $1 \le t \le m$, i.i.d. sample *K* tokens e_{i_1}, \ldots, e_{i_K} from to α_t , average these tokens to form $z_t = \frac{1}{K} \sum_{r=1}^{K} e_{i_r}$, which it uses as query for the next step.
- **Interpretation:** CoT2-MTS tracks *K* traces in parallel according to their discrete CoT likelihoods. However, these traces are not statistically independent.
- 450

With these three inference methods defined, we present the
following result on the statistical consistency of the final
outputs of these methods.

455 **Proposition 2** (Consistency of CoT and CoT2 inference). 456 Under Assumption 1 and given X, the output of base CoT2 457 is $z_m = \sum_{j=1}^{v} \alpha_{m,j} e_j$ where $\alpha_m = \alpha_1 \prod_{t=1}^{m-1} M_t$. Discrete 458 CoT and CoT2-MTS have the same output once we take the 459 expectation over their stochastic sampling. 460

461 **Remark:** Please note that α_m represents a distribution over 462 the vocabulary. The proposition above shows that as the 463 number of samples approaches infinity, the empirical dis-464 tribution $\hat{\alpha}_m$ obtained from CoT2-MTS (or discrete CoT) 465 traces converge in probability to the deterministic distri-466 bution α_m of the base CoT2 model. α_m is computed in a 467 sampling-free fashion and thus, is not a random variable.

468 Overall Proposition 2 establishes the statistical consistency 469 of all three inference methods as they all estimate the same 470 distribution α_m over the vocabulary. However, they differ in 471 how many samples are needed to approximate that distribu-472 tion with repeated samplings. In particular, the base CoT2 473 model outputs the entire probability distribution over tokens 474 at every intermediate step, thereby, implicitly tracking all 475 possible trajectories in parallel as continuous embeddings. 476 Consequently, it directly computes the exact final token dis-477 tribution in one forward pass without repeated sampling. In 478 contrast, due to stochasticity, discrete CoT or CoT2-MTS 479 needs multiple i.i.d. samples to approximate this distribu-480 tion. This observation motivates us to study and contrast the 481 sample complexities of discrete CoT and CoT2-MTS. Our 482 next proposition establishes our distribution approximation 483 guarantee with respect to the total variation distance and 484 shows that CoT2-MTS reduces the sample complexity of 485 estimation compared to discrete CoT by a factor of K. 486

Proposition 3. Let α_m be the final expected output distribution after m steps of CoT according to Proposition 2. Let $\hat{\alpha}_m$ be the distribution resulting from averaging the outputs of i.i.d. CoT2-MTS traces with parallelism K. Then, to guarantee $\|\hat{\alpha}_m - \alpha_m\|_1 \le \epsilon$ with high probability, the total number of samples (traces) required scales as $\Theta\left(\frac{\nu}{K\epsilon^2}\right)$.

493 494 Recall that CoT2-MTS generalizes discrete CoT (K = 1).

For K = 1, which corresponds to the discrete CoT model, the above proposition reduces to the known $\Theta(\frac{v}{\epsilon^2})$ sample complexity of approximating a *v*-category distribution in ℓ^1 distance (Kamath et al., 2015). Note that as $K \to \infty$, CoT2-MTS converges to the base CoT2, and the above proposition recovers the one-shot performance of base CoT2. Thus, although the three models yield the same final distribution, the discrete model requires substantially more rollouts for accurate approximation due to inherent noise from singletoken sampling. In contrast, the base CoT2 model carries the entire mixture of partial expansions at each step and computes the distribution in one shot, while the CoT2-MTS model captures multiple partial expansions in each step and proportionally reduces the sample complexity. This theoretical intuition aligns with our empirical findings in the Pass@k experiments in Section 3, where we observe that the base CoT2 approach achieves comparable performance to discrete CoT while requiring substantially fewer samples.

6. Limitations

One limitation of CoT2 is that it may require larger embedding dimensions to represent multiple parallel reasoning traces in more complex tasks. As a broader impact, while CoT2 can increase performance by searching more reasoning paths in parallel, this representation shadows the model's intermediate decision process, and might potentially reduce the interpretability of the model.

References

- Aggarwal, P. and Welleck, S. L1: Controlling how long a reasoning model thinks with reinforcement learning. *arXiv preprint arXiv:2503.04697*, 2025.
- Bachmann, G. and Nagarajan, V. The pitfalls of next-token prediction. In *Forty-first International Conference on Machine Learning*, 2024. URL https://openreview. net/forum?id=76zq8Wk16Z.
- Bengio, S., Vinyals, O., Jaitly, N., and Shazeer, N. Scheduled sampling for sequence prediction with recurrent neural networks, 2015. URL https://arxiv.org/abs/ 1506.03099.
- Blei, D. M., Ng, A. Y., and Jordan, M. I. Latent dirichlet allocation. J. Mach. Learn. Res., 3(null):993–1022, March 2003. ISSN 1532-4435.
- Cheng, J. and Van Durme, B. Compressed chain of thought: Efficient reasoning through dense representations. *arXiv* preprint arXiv:2412.13171, 2024.
- Deng, Y., Prasad, K., Fernandez, R., Smolensky, P., Chaudhary, V., and Shieber, S. Implicit chain of thought

reasoning via knowledge distillation. *arXiv preprintarXiv:2311.01460*, 2023.

- 498 Deng, Y., Choi, Y., and Shieber, S. From explicit cot to
 499 implicit cot: Learning to internalize cot step by step.
 500 arXiv preprint arXiv:2405.14838, 2024.
- Dziri, N., Lu, X., Sclar, M., Li, X. L., Jiang, L., Lin, B. Y., West, P., Bhagavatula, C., Bras, R. L., Hwang, J. D., Sanyal, S., Welleck, S., Ren, X., Ettinger, A., Harchaoui, Z., and Choi, Y. Faith and fate: Limits of transformers on compositionality, 2023. URL https://arxiv.org/ abs/2305.18654.
- Feng, G., Zhang, B., Gu, Y., Ye, H., He, D., and Wang, L.
 Towards revealing the mystery behind chain of thought:
 A theoretical perspective. In *Thirty-seventh Conference* on Neural Information Processing Systems, 2023. URL https://openreview.net/forum?id=qHrADgAdYu.
- 514 Geiping, J., McLeish, S., Jain, N., Kirchenbauer, J., Singh,
 515 S., Bartoldson, B. R., Kailkhura, B., Bhatele, A., and
 516 Goldstein, T. Scaling up test-time compute with latent
 517 reasoning: A recurrent depth approach. *arXiv preprint*518 *arXiv:2502.05171*, 2025.
- 519 Giannou, A., Rajput, S., Sohn, J.-Y., Lee, K., Lee, J. D., 520 and Papailiopoulos, D. Looped transformers as pro-521 grammable computers. In Krause, A., Brunskill, E., Cho, 522 K., Engelhardt, B., Sabato, S., and Scarlett, J. (eds.), 523 Proceedings of the 40th International Conference on Ma-524 chine Learning, volume 202 of Proceedings of Machine 525 Learning Research, pp. 11398-11442. PMLR, 23-29 526 Jul 2023. URL https://proceedings.mlr.press/ 527 v202/giannou23a.html. 528
- Gloeckle, F., Idrissi, B. Y., Rozière, B., Lopez-Paz, D., and
 Synnaeve, G. Better & faster large language models via
 multi-token prediction. In *Proceedings of the 41st International Conference on Machine Learning*, ICML'24.
 JMLR.org, 2024.
- Goyal, S., Ji, Z., Rawat, A. S., Menon, A. K., Kumar, S., and
 Nagarajan, V. Think before you speak: Training language
 models with pause tokens. In *The Twelfth International Conference on Learning Representations*, 2024. URL
 https://openreview.net/forum?id=ph04CRkPdC.
- Guo, D., Yang, D., Zhang, H., Song, J., Zhang, R., Xu, R.,
 Zhu, Q., Ma, S., Wang, P., Bi, X., et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.
- Hao, S., Sukhbaatar, S., Su, D., Li, X., Hu, Z., Weston,
 J., and Tian, Y. Training large language models to reason in a continuous latent space. *arXiv preprint arXiv:2412.06769*, 2024.

- Ildiz, M. E., Huang, Y., Li, Y., Rawat, A. S., and Oymak, S. From self-attention to markov models: Unveiling the dynamics of generative transformers, 2024. URL https://arxiv.org/abs/2402.13512.
- Jaech, A., Kalai, A., Lerer, A., Richardson, A., El-Kishky, A., Low, A., Helyar, A., Madry, A., Beutel, A., Carney, A., et al. Openai o1 system card. arXiv preprint arXiv:2412.16720, 2024.
- Kamath, S., Orlitsky, A., Pichapati, D., and Suresh, A. T. On learning distributions from their samples. In Grünwald, P., Hazan, E., and Kale, S. (eds.), *Proceedings of The 28th Conference on Learning Theory*, volume 40 of *Proceedings of Machine Learning Research*, pp. 1066– 1100, Paris, France, 03–06 Jul 2015. PMLR. URL https: //proceedings.mlr.press/v40/Kamath15.html.
- Kojima, T., Gu, S. S., Reid, M., Matsuo, Y., and Iwasawa, Y. Large language models are zero-shot reasoners. *Advances in neural information processing systems*, 35: 22199–22213, 2022.
- Liu, A., Feng, B., Xue, B., Wang, B., Wu, B., Lu, C., Zhao, C., Deng, C., Zhang, C., Ruan, C., et al. Deepseek-v3 technical report. arXiv preprint arXiv:2412.19437, 2024.
- Merrill, W. and Sabharwal, A. The expressive power of transformers with chain of thought. In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id= NjNGlPh8Wh.
- Nanda, N., Chan, L., Lieberum, T., Smith, J., and Steinhardt, J. Progress measures for grokking via mechanistic interpretability, 2023. URL https://arxiv.org/abs/ 2301.05217.
- Nye, M., Andreassen, A. J., Gur-Ari, G., Michalewski, H., Austin, J., Bieber, D., Dohan, D., Lewkowycz, A., Bosma, M., Luan, D., et al. Show your work: Scratchpads for intermediate computation with language models. *arXiv:2112.00114*, 2021.
- Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright, C., Mishkin, P., Zhang, C., Agarwal, S., Slama, K., Ray, A., et al. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35:27730–27744, 2022.
- Pfau, J., Merrill, W., and Bowman, S. R. Let's think dot by dot: Hidden computation in transformer language models. In *First Conference on Language Modeling*, 2024. URL https://openreview.net/forum?id=NikbrdtYvG.
- Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., and Sutskever, I. Language models are unsupervised multitask learners. *OpenAI Technical*

- 550 551 552 553 554 555 556 557 558 559 560 561 562 563 564 565 566 567 568 569 570 571 572 573 574 575 576 577 578 579 580
 - Report, 2019. URL https://cdn.openai.com/ better-language-models/language_models_ are_unsupervised_multitask_learners.pdf.
 - Saparov, A. and He, H. Language models are greedy reasoners: A systematic formal analysis of chain-of-thought. *arXiv preprint arXiv:2210.01240*, 2022.
 - Saunshi, N., Dikkala, N., Li, Z., Kumar, S., and Reddi,
 S. J. Reasoning with latent thoughts: On the power of
 looped transformers. In *The Thirteenth International Conference on Learning Representations*, 2025. URL
 https://openreview.net/forum?id=din0lGfZFd.
 - Shalev, Y., Feder, A., and Goldstein, A. Distributional reasoning in llms: Parallel reasoning processes in multihop reasoning. *arXiv preprint arXiv:2406.13858*, 2024.
 - Shao, Z., Wang, P., Zhu, Q., Xu, R., Song, J., Bi, X., Zhang, H., Zhang, M., Li, Y. K., Wu, Y., and Guo, D. Deepseekmath: Pushing the limits of mathematical reasoning in open language models, 2024. URL https://arxiv.org/abs/2402.03300.
 - Shazeer, N. Glu variants improve transformer. *arXiv* preprint arXiv:2002.05202, 2020.
 - Shen, Z., Yan, H., Zhang, L., Hu, Z., Du, Y., and He, Y. Codi: Compressing chain-of-thought into continuous space via self-distillation. *arXiv preprint arXiv:2502.21074*, 2025.
- Silver, D., Schrittwieser, J., Simonyan, K., Antonoglou,
 I., Huang, A., Guez, A., Hubert, T., Baker, L., Lai, M.,
 Bolton, A., and et al. Mastering the game of go without human knowledge. *Nature*, 550(7676):354–359, Oct
 2017. doi: 10.1038/nature24270.
- 583
 584
 584
 585
 586
 586
 586
 587
 587
 588
 589
 580
 580
 581
 581
 582
 583
 584
 585
 586
 587
 587
 587
 588
 588
 588
 589
 589
 589
 580
 580
 580
 581
 581
 582
 582
 583
 584
 584
 585
 586
 586
 587
 587
 587
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
 588
- Stiennon, N., Ouyang, L., Wu, J., Ziegler, D. M., Lowe,
 R., Voss, C., Radford, A., Amodei, D., and Christiano,
 P. Learning to summarize from human feedback, 2022.
 URL https://arxiv.org/abs/2009.01325.
- Sui, Y., Chuang, Y.-N., Wang, G., Zhang, J., Zhang, T.,
 Yuan, J., Liu, H., Wen, A., Zhong, S., Chen, H., et al.
 Stop overthinking: A survey on efficient reasoning for
 large language models. *arXiv preprint arXiv:2503.16419*,
 2025.
- Suzgun, M., Scales, N., Schärli, N., Gehrmann, S., Tay, Y., Chung, H. W., Chowdhery, A., Le, Q., Chi, E., Zhou, D., and Wei, J. Challenging BIG-bench tasks and whether chain-of-thought can solve them. In Rogers, A., Boyd-Graber, J., and Okazaki, N. (eds.), *Findings of the Association for Computational Linguistics: ACL 2023*, pp.

13003–13051, Toronto, Canada, July 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023. findings-acl.824. URL https://aclanthology.org/ 2023.findings-acl.824/.

- Thomm, J., Camposampiero, G., Terzic, A., Hersche, M., Schölkopf, B., and Rahimi, A. Limits of transformer language models on learning to compose algorithms, 2024. URL https://arxiv.org/abs/2402.05785.
- Wang, X., Wei, J., Schuurmans, D., Le, Q., Chi, E., Narang, S., Chowdhery, A., and Zhou, D. Self-consistency improves chain of thought reasoning in language models. *arXiv preprint arXiv:2203.11171*, 2022.
- Wei, J., Wang, X., Schuurmans, D., Bosma, M., Xia, F., Chi, E., Le, Q. V., Zhou, D., et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837, 2022.
- Xiong, Z., Cai, Z., Cooper, J., Ge, A., Papageorgiou, V., Sifakis, Z., Giannou, A., Lin, Z., Yang, L., Agarwal, S., Chrysos, G. G., Oymak, S., Lee, K., and Papailiopoulos, D. Everything everywhere all at once: Llms can incontext learn multiple tasks in superposition, 2024. URL https://arxiv.org/abs/2410.05603.
- Yang, S., Gribovskaya, E., Kassner, N., Geva, M., and Riedel, S. Do large language models latently perform multi-hop reasoning? In Ku, L.-W., Martins, A., and Srikumar, V. (eds.), *Proceedings of the* 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pp. 10210– 10229, Bangkok, Thailand, August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024. acl-long.550. URL https://aclanthology.org/ 2024.acl-long.550/.
- Yao, S., Yu, D., Zhao, J., Shafran, I., Griffiths, T., Cao, Y., and Narasimhan, K. Tree of thoughts: Deliberate problem solving with large language models. *Advances in neural information processing systems*, 36:11809–11822, 2023.
- Yu, P., Xu, J., Weston, J., and Kulikov, I. Distilling system 2 into system 1. *arXiv preprint arXiv:2407.06023*, 2024.
- Yu, Q., Zhang, Z., Zhu, R., Yuan, Y., Zuo, X., Yue, Y., Fan, T., Liu, G., Liu, L., Liu, X., Lin, H., Lin, Z., Ma, B., Sheng, G., Tong, Y., Zhang, C., Zhang, M., Zhang, W., Zhu, H., Zhu, J., Chen, J., Chen, J., Wang, C., Yu, H., Dai, W., Song, Y., Wei, X., Zhou, H., Liu, J., Ma, W.-Y., Zhang, Y.-Q., Yan, L., Qiao, M., Wu, Y., and Wang, M. Dapo: An open-source llm reinforcement learning system at scale, 2025. URL https://arxiv.org/abs/2503.14476.

605 606 607	Zhang, X., Huang, Z., Ni, C., Xiong, Z., Chen, J., and Oymak, S. Making small language models efficient rea- soners: Intervention supervision reinforcement <i>arXiv</i>
608	nrenrint arYiv:2505.07061, 2025
600	preprint arXiv.2505.07901, 2025.
610	
611	
612	
612	
013	
615	
010	
010	
01/	
618	
619	
620	
621	
622	
623	
624	
625	
626	
627	
620	
620	
621	
632	
633	
634	
635	
636	
637	
638	
639	
640	
641	
642	
643	
644	
645	
646	
647	
648	
649	
650	
651	
652	
653	
654	
655	
656	
657	
650	

APPENDIX

We discuss additional related work in Appendix A. We provide further implementation details in Appendix B, including those
for the MNNS task (Appendix B.1), the ProntoQA/ProsQA datasets (Appendix B.2), and GRPO training (Appendix B.3).
We present additional experimental results in Appendix C, and we offer details on continuous supervised training and GRPO
in Appendix C.1 and Appendix C.2, respectively. Finally, we include the proofs of Propositions 1, 2, and 3 in Appendix D.

A. Further Related Work

The proposed CoT2 approach simultaneously tracks all possible trajectories and superposes them within continuous tokens. This approach is similar to that of Xiong et al. (2024), who superpose multiple candidate outputs into a single final token. Our approach also shares similarities with decoding algorithms like self-consistency (Wang et al., 2022) and Best-of-N-Sampling (Stiennon et al., 2022), which generate multiple trajectories by running inference multiple times and then select a final answer based on the aggregate statistics. In contrast, our algorithm performs a single inference, superposing different trajectories all at once and determining the final answer in one pass. Furthermore, our Dirichlet sampling approach for generating multiple rollouts in GRPO training draws connections to previous works such as Latent Dirichlet Allocation (LDA) (Blei et al., 2003), which introduces Dirichlet priors within a hierarchical Bayesian framework, and AlphaGo (Silver et al., 2017), which injects Dirichlet noise to encourage exploration.

Our work also tangentially relates to research on multi-token prediction (Bachmann & Nagarajan, 2024; Liu et al., 2024; Gloeckle et al., 2024), which aims to improve the efficiency and quality of generation by predicting multiple tokens at once. It is hypothesized that effective future prediction necessitates the exploration of many possible continuations, which is similar to our CoT2 approach.

B. Implementation Details

Computational Resources: All experiments were run on a Slurm-managed cluster using L40S GPUs with 48GB of memory. Each experiment fits on a single GPU. In the case of 4 input digits, the SFT or CSFT training takes approximately 3 hours on a single GPU. For 5-digit inputs, the dataset size increases by roughly a factor of 10, and the training time increases proportionally. The entire codebase was implemented in PyTorch.

B.1. Implementation Details of Experiments on MNNS Task

Dataset Details: For the MNNS task, the vocabulary consists of a range of numbers from [-S, S] for some positive integer *S*, together with $\langle BOS \rangle$, $\langle EOS \rangle$, and \rightarrow special tokens. The integer *S* is chosen so that all possible partial sums of the selected input digits lie within [-S, S]. For example, when the input digits lie in the range 1–10, we set S = 36, whereas for digits in 5–14, we set S = 40. We performed our experiments on the 4 and 5 input digit scenarios. A sample input line with *m* numbers is:

$$\langle BOS \rangle D_1 D_2 \dots D_m \rightarrow$$

Accordingly, the output will be *m* sum tokens, where the final token corresponds to the answer, followed by *<EOS>* token:

$$S_1 S_2 \ldots S_m < EOS >$$

As a concrete example, consider the input 2, 1, 4 (m = 3), following Figure 1. In this case, the solution for the MNNS task is -2 - 1 + 4 = 1. Therefore, for the discrete model, the input is $\langle BOS \rangle D_2 D_1 D_4 \rightarrow$ and we supervise it along the trajectory of correct output tokens $S_{-2} S_{-3}$, $S_1 \langle EOS \rangle$, as illustrated in Figure 1. On the other hand, the continuous supervision at the first step holds S_2 and S_{-2} as possibilities. Then, for the next step, we add 1 or -1 to these numbers, and the resulting possibilities are S_3 , S_1 , S_{-1} , S_{-3} . Finally, at the last step, the model is supervised to pick the correct answer S_1 as the token.

We split the datasets by ensuring that each permutation of a set of numbers is exactly in one of the train and validation datasets, as the answer to the question is permutation-invariant. This way, we prevent the models from memorizing the answer and make a fair comparison. We also use 0.8-0.2 split for train-val datasets.

Model and Hyperparameters: We use the GPT2 model, with 1 layer 1 head, 2 layer 2 head, and 4 layer 4 head as the configurations. For each configuration, we experiment with embedding dimensions of 16, 24, or 32. We train with a learning rate of $1r = 10^{-4}$ and use AdamW (no weight decay). The batch size is 16 for 4-digit inputs and 64 for 5-digit inputs.

Final answer of the models: To make a proper comparison, we only check the final answer of the models, as checking the correctness of the full path of the discrete model would be unfair.

Pass@k Experiment: We perform our experiments for temperatures 0, 0.4, 0.8, and 1 by repeating the evaluation 10 times for each *k* value where k changes from 1 to 14.

B.2. Implementation Details of Experiments on ProntoQA/ProsQA Datasets

Dataset Details: Different from the original ProntoQA/ProsQA datasets which described the structured ontology in natural language as a set of known conditions, we use a more structured format through a token-level representation. An example prompt is shown below.

Description of the structured ontology: Each component of the ontology and associated questions is represented through discrete tokens with their own learned embeddings, rather than as raw textual input. Specifically, we use the GPT-2 architecture and encode the ontology's structural components. Below are two examples demonstrating how natural-language assertions are mapped to our tokenized format:

Brimpuses are not luminous \rightarrow 'A' 'not in' 'B' '.'.

Shumpuses are amenable; Each yumpus is a lorpu; Every lorpus is floral \rightarrow 'C' 'in' 'D' '.'.

Below, we have the ProntoQA and ProsQA datasets' input-output format.

The structure of ProntoQA:

Input:'Description' '{' 'A' 'not in' 'B' '.' ... 'C' 'in' 'D' '.' '}' 'Question' '{' 'C'
'not in' 'F' '.' '}'

Output:'Steps' { 'C' 'in' 'D' '.' ... 'D' 'in' 'E' '.' '}' 'Answer' '{' 'False' '}'

The structure of ProsQA:

Input:'Description' '{' 'A' 'in' 'B' '.' ... 'C' 'in' 'D' '.' '}' 'Question' '{' 'C' 'in' 'F'
'or' 'E' '}'

Output:'Steps' { 'C' 'in' 'D' '.' ... 'D' in 'E' '.' '}' 'Answer' '{' 'F' '}'

Each distinct component or relation (e.g., 'A', 'in', 'not in') is treated as a unique token, and singular/plural variants (such as 'lempus' and 'lempuses') are collapsed into a single token to simplify the vocabulary. Alongside these concept tokens, special structural tokens ('Description', '', '.', 'or', etc.) are also included, which results in a vocabulary size of 31 tokens. To avoid biases, we balance the dataset. In ProntoQA, "yes" and "no" each appear with 50% probability, and in ProsQA, the correct answer is randomly permuted at the first or second position. For all the other experimental and training settings, we follow (Hao et al., 2024).

² Model and Hyperparameters: We use the GPT2 model, with 2 layer 2 head, and 4 layer 4 head as the configurations. We tested embedding dimensions 24, 32, 40 with these configurations. We set batch size 64. We train with a learning rate of 10^{-4} and use AdamW (no weight decay).

Maj@k Experiment: We use majority voting for evaluation instead of Pass@k, because both ProntoQA and ProsQA are binary questions. We perform our experiments for temperatures 0, 0.4, 0.8, and 1 by repeating the evaluation 10 times for each k value where k changes from 1 to 21. If two or more answers end up with the same top vote, we pick one randomly.

B.3. Implementation Details of GRPO Training

In (Hao et al., 2024), the reference model is updated by $LM_{\theta_{ref}} \leftarrow LM_{\theta}$ in each iteration (epoch). This approach is reasonable for their setting with a large dataset and a small number of epochs over it. For our setting, however, we set the reference model to the initial model and never update it through iterations as we have a smaller dataset. Meanwhile, we update the old model before every batch $LM_{\theta_{old}} \leftarrow LM_{\theta}$.

In our experiments, we use G = 8 trajectories per input data point, use clipping parameter $\epsilon = 0.1$, and set the KL-divergence coefficient $\beta = 0$ in most cases (with $\beta = 0.1$ in a few). For the CoT2 model with MTS sampling, we change the number of tokens to sample K from 1 to 12. In the MNNS task, the 5-digit case has about ten times more data than the 4-digit

case, so we typically focus on 4-digit MNNS because of computational considerations and use a batch size of 16 in those
 experiments.

Learning rates differ by model and setting. We use $lr = 5 \times 10^{-5}$ for CoT2-MTS sampling (figures in the main text), 1r = 1×10^{-5} for discrete CoT with Dirichlet sampling, and $lr = 1 \times 10^{-6}$ for CoT2 with Dirichlet sampling. For ProntoQA and ProsQA experiments, we perform a grid search over learning rates ranging from 1×10^{-4} to 1×10^{-8} and select and report results using the best-performing configuration. For most settings, we find $lr = 1 \times 10^{-5}$ optimal; however, for CoT2 and discrete CoT models with K = 6, we set $lr = 1 \times 10^{-6}$. We also use AdamW with a weight-decay of 0.01. For Dirichlet experiments on the MNNS task, we try various scale parameters γ , but we find $\gamma = 20$ to work best in most settings. Unless stated otherwise, we report the best validation accuracy found during training for each setting.

780781 C. Experimental Results

782 783 C.1. Continuous Supervised Training Results

Teacher Forcing and Self-feeding Comparison: As described in Section 3, we tested two approaches of providing prefixes during training the CoT2 model with CSFT. Although the model autoregressively generates at inference time, teacher forcing yields better performance than self-feeding during CSFT training. Our results demonstrate that, We also tested curriculum settings, where we switch to self-feeding after a pre-determined number of epochs in the training. Still, the accuracies didn't improve beyond pure teacher-forcing training. The results are illustrated in Figure 4, where we refer to teacher-forcing as "hard-teacher" and refer to self-feeding as "soft-teacher".

790 Sparse Supervision for Discrete Baseline: We also tested providing a subset of the correct path to the discrete model. We 791 observed that a sparsely supervised discrete model can achieve better performance than the fully supervised discrete model 792 when the distribution is "easier" to handle by the model. As an example, we tested the case when we have 5 input digits from 793 the range of 11 to 19. In this case, in nearly all of the cases, the answer to our MNNS game is (sum of minimum 3 numbers) 794 - (sum of maximum 2 numbers) out of the 5 input numbers. In this case, when only 1 token from the correct path is provided 795 to the discrete model, it's better than 3 and 5 token cases. However, when we change the distribution to a range of numbers 796 from 5 to 13, which makes the question reasonably harder, the discrete model with 1 token supervision performs worse than 797 the other two, and the discrete model with full supervision performs best. The results are demonstrated in Figure 3. 798

799 Further Results on CoT2 vs Discrete CoT: The results in Figure 5 also indicate that above an embedding dimension 800 threshold, the CoT2 model has superior performance and trains significantly faster than the discrete CoT model. Moreover, 801 combining the results of Figure 5 with Figure 2c, we see that the CoT2 model with one layer and one head GPT2 model 802 performs better than discrete CoT model with two layers and two heads at embeddings 24 and 32. While the continuous 803 approach requires greater embedding capacity to support its distributional representations at each step, it can outperform 804 the discrete model using fewer layers and attention heads. Supporting the findings in Figure 2, Figure 6 illustrates that on 805 the ProntoQA task, CoT2 consistently outperforms the discrete CoT baseline when the embedding dimension is above a 806 threshold. Likewise, as depicted in Figure 7 and Figure 8, the discrete CoT model requires multiple samplings (Maj@k) 807 to match the single-shot performance of CoT2 on both ProntoQA and ProsQA, which indicates that CoT2 model is more 808 sample-efficient.

- 809
- 810
- 811
- 812
- 813 814
- 815
- 816
- 817 818
- 819
- 820
- 821
- 822
- 823
- 824



Figure 3: The figure illustrates that when the range of digits makes the question non-trivial on an MNNS task, the discrete CoT model trained with full token supervision outperforms sparse supervisions; in particular, single token supervision yields the worst performance. **Setting:** 5 input digits in 5 - 13; 2-layer, 2-head GPT2 with d = 32.



Figure 4: The comparison between the hard and soft teachers for different embedding dimensions. The figure illustrates that the hard teacher is superior to the soft teacher. Setting: 4 input digits in 1 - 9; 4-layer, 4-head GPT2 with $d \in \{16, 24, 32\}$.



Figure 5: Comparison between CoT2 and discrete CoT2 model for different embedding dimensions. The figure demonstrates that above a certain embedding dimension threshold, the CoT2 model outperforms the discrete CoT model in the MNNS task. **Setting:** 4 input digits in 1 - 9; 1-layer, 1-head GPT2 with $d \in \{16, 24, 32\}$.



Figure 6: Comparison between CoT2 and discrete CoT2 model for different embedding dimensions in ProntoQA task. The figure shows that above an embedding dimension threshold, the CoT2 model outperforms the discrete CoT model. Setting: 4 input digits in 1 - 9; 4-layer, 4-head GPT2 with $d \in \{24, 32, 30\}$.



Figure 7: The figure illustrates that the discrete CoT2 model requires multiple samplings (Maj@k) to match the single-shot performance of the CoT2 model on ProntoQA (10-run average). Setting: 4-layer, 4-head GPT2 with d = 32.



Figure 8: The figure illustrates that the discrete CoT2 model requires multiple samplings (Maj@k) to match the single-shot performance of the CoT2 model on ProsQA (10-run average). Setting: 4-layer, 4-head GPT2 with d = 32.

935 C.2. GRPO Results

Discussion on ProntoQA/ProsQA Datasets: Table 2 illustrates that GRPO training using CoT2-MTS sampling consistently improves discrete CoT and CoT2 models over their initial SFT accuracy. Moreover, we observe that the improvement in the discrete CoT model is greater, which might indicate that the CoT2 model already gains an RL-like exploration mechanism through CSFT training. We observe that while increasing K initially increases the accuracy by sampling more tokens at each step, beyond some K, the improvements diminish. This observation is consistent with Table 1, where we see that a moderate K value offers the best final performance. One possible explanation is that while higher K promotes better exploration, it also raises the chance of sampling unhelpful tokens that disrupt the averaged token representation. Indeed, for larger K, we observe that the RL objective saturates to near zero which suggests that most rollouts fail once the averaged token contains too many distracting tokens.

Discussion on Dirichlet Sampling: We also investigate the effects of Dirichlet sampling in GRPO training discrete CoT and CoT2 models. The results in Table 3 indicates that applying Dirichlet sampling ($\gamma = 20$) in GRPO training of discrete CoT model consistently improves over the initial SFT training accuracies. Similar to the CoT2 + MTS sampling results in Table 1, we observe that the entropy at the third token remains relatively high, which suggests a beneficial diversity in model's predictions for that token. Moreover, the Table 4 demonstrates that Dirichlet sampling also improves the CoT2 model's SFT accuracy, even though it has a high initial SFT accuracy. As illustrated in Table 4, we find there is an optimal value for the scale parameter γ , since larger γ typically yields more uniform sampling distributions, whereas smaller γ concentrates the distribution more sharply. Thus, adjusting γ provides a balance between exploration and stability in GRPO training.

Table 3: Discrete CoT models trained with GRPO after SFT using Dirichlet sampling ($\gamma = 20$) and a learning rate of 1 × 10⁻⁵. We show validation accuracy (%) and token-level entropy (SFT \rightarrow SFT+GRPO) for each (Layers, Heads) setting, with an embedding dimension of 24 for GPT2 model.

0	Lavers	Heads	Val. Accuracy (%)		Val. Entropy (SFT \rightarrow SFT+GRPO)					
62	, 110	1104405	SFT	SFT+GRPO	token ₁	token ₂	token ₃	token ₄		
3	1	1	39.76	46.25	$0.4851 \rightarrow 0.1701$	$0.5165 \rightarrow 0.6380$	$0.3243 \rightarrow 0.6590$	$0.1597 \rightarrow 0.4878$		
4 5	2	2	70.26	75.84	$0.4851 \rightarrow 0.4027$	$0.5165 \rightarrow 0.4413$	$0.3243 \rightarrow 0.2907$	$0.1597 \to 0.1386$		

Table 4: Validation accuracies GRPO training with CoT2 models using different Dirichlet sampling scales (γ) with learning rate of 1 × 10⁻⁶. We show the baseline SFT accuracy (87.84%) and final performance after GRPO.

Dirichlet Scale (γ)	SFT Val. Acc (%)	SFT + GRPO Val. Acc (%				
10		89.76				
20	87.84	90.75				
40		90.37				

D. Theoretical Details

Justification for Assumption 1: The Assumption 1 holds for tasks where the next-token distribution depends solely on the current token and the input tokens rather than the full history of output tokens. This is satisfied by many reasoning tasks, where the aim is to keep track of an intermediate state (e.g., the current sum) and update this state based only on the current state and the input, independently of the earlier trajectory.

For example, in the MNNS task, the model generates a token representing the current partial sum at each step. To compute the distribution over the next possible sums, the model adds or subtracts the selected number from the input context X to the current sum, without needing to remember the sequence of previous sums explicitly. Thus, the next-state distribution at each step is only determined by the current state and it naturally satisfies the Assumption 1.

Proposition 2 (Consistency of CoT and CoT2 inference). Under Assumption 1 and given X, the output of base CoT2 is $z_m = \sum_{j=1}^{v} \alpha_{m,j} e_j$ where $\alpha_m = \alpha_1 \prod_{t=1}^{m-1} M_t$. Discrete CoT and CoT2-MTS have the same output once we take the expectation over their stochastic sampling. *Proof.* Let $\hat{\alpha}_t^{(\text{disc})}$, $\hat{\alpha}_t^{(\text{MTS})}$ denote the empirical output token distributions at step *t* under one trajectory obtained by the discrete CoT, and CoT2 with MTS models, respectively. We define $\alpha_t^{(\text{disc})} = \mathbb{E}[\hat{\alpha}_t^{(\text{disc})}]$ and $\alpha_t^{(\text{MTS})} = \mathbb{E}[\hat{\alpha}_t^{(\text{MTS})}]$ to be

Proof. Let $\hat{\alpha}_{t}^{(\text{disc})}, \hat{\alpha}_{t}^{(\text{disc})}$ denote the empirical output token distributions at step *t* under one trajectory obtained by the discrete CoT, and CoT2 with MTS models, respectively. We define $\alpha_{t}^{(\text{disc})} = \mathbb{E}[\hat{\alpha}_{t}^{(\text{disc})}]$ and $\alpha_{t}^{(\text{MTS})} = \mathbb{E}[\hat{\alpha}_{t}^{(\text{MTS})}]$ to be corresponding expected distributions. The discrete CoT model at each step picks exactly 1 token from α_{t} . On the other hand, CoT2-MTS samples *K* i.i.d. tokens at every step independently according to their probabilities from $\hat{\alpha}_{t}$. We denote them i_{1}, \ldots, i_{K} , and average their embeddings to produce a single query.

We will use induction in our argument. For the base case, all models start with the same initial distribution, so we trivially have $\alpha_1^{(\text{disc})} = \alpha_1^{(\text{MTS})} = \alpha_1^{(\text{CoT2})}$. For the inductive step, assume that we have $\alpha_{t-1}^{(\text{disc})} = \alpha_{t-1}^{(\text{MTS})} = \alpha_{t-1}^{(\text{CoT2})}$. We will show that $\alpha_t^{(\text{disc})} = \alpha_t^{(\text{MTS})} = \alpha_t^{(\text{CoT2})}$. On the other, for the discrete CoT model, the model samples one token $e_{i_{t+1}}$ from the row of M_t for a token i_{t+1} . Therefore, we need to condition on the token at step *t*. We have:

$$\mathbb{E}\left[\hat{\boldsymbol{\alpha}}_{t+1}^{(\text{disc})}\right] = \sum_{j=1}^{\nu} \mathbb{P}(\boldsymbol{z}_{t} = \boldsymbol{e}_{j}) \mathbb{E}\left[\hat{\boldsymbol{\alpha}}_{t+1}^{(\text{disc})} \mid \boldsymbol{z}_{t} = \boldsymbol{e}_{j}\right] = \sum_{j=1}^{\nu} \mathbb{P}(\boldsymbol{z}_{t} = \boldsymbol{e}_{j}) \operatorname{LM}_{\theta}^{(t)}(\cdot \mid \boldsymbol{e}_{j}, \boldsymbol{X})$$
$$\sum_{j=1}^{\nu} \boldsymbol{\alpha}_{t,j}^{(\text{disc})} \operatorname{LM}_{\theta}^{(t)}(\cdot \mid \boldsymbol{e}_{j}, \boldsymbol{X})$$
$$\stackrel{(a)}{=} \sum_{j=1}^{\nu} \boldsymbol{\alpha}_{t,j}^{(\text{CoT2})} \operatorname{LM}_{\theta}^{(t)}(\cdot \mid \boldsymbol{e}_{j}, \boldsymbol{X})$$
$$\stackrel{(b)}{=} \boldsymbol{\alpha}_{t+1}^{(\text{CoT2})}$$
(6)

where (a) follows from the induction argument and (b) follows from Assumption 1. Therefore, we obtain $\alpha_{t+1}^{(\text{disc})} = \mathbb{E}\left[\hat{\alpha}_{t+1}^{(\text{disc})}\right] = a_{t+1}^{(\text{CoT2})}$. For the CoT2-MTS model, the argument will be similar. Using the decoupling of trajectories by Assumption 1, the next distribution is:

 $\mathrm{LM}_{\theta}^{(t)}\left(\cdot \mid \frac{1}{K} \sum_{i=1}^{K} \boldsymbol{e}_{i_{r}}, \boldsymbol{X}\right) = \frac{1}{K} \sum_{i=1}^{K} \mathrm{LM}_{\theta}^{(t)}\left(\cdot \mid \boldsymbol{e}_{i_{r}}, \boldsymbol{X}\right).$

1045 Therefore, we write:

1076 1077

$$\mathbb{E}\left[\hat{a}_{i+1}^{(MTS)}\right] = \sum_{(i_1,...,i_K)\in[v]^K} \mathbb{P}(e_{i_1},...,e_{i_K}) \mathbb{E}\left[\hat{a}_{i+1}^{(MTS)} \mid e_{i_1},...,e_{i_K}\right]$$

$$= \sum_{(i_1,...,i_K)\in[v]^K} \mathbb{P}(e_{i_1},...,e_{i_K}) \mathbb{L}M_{\theta}^{(i)}\left(\cdot\mid\frac{1}{K}\sum_{r=1}^{K}e_{i_r},X\right)$$

$$= \sum_{(i_1,...,i_K)\in[v]^K} \mathbb{P}(e_{i_1},...,e_{i_K}) \frac{1}{K}\sum_{r=1}^{K} \mathbb{L}M_{\theta}^{(i)}\left(\cdot\mid e_{i_r},X\right)$$

$$= \sum_{(i_1,...,i_K)\in[v]^K} \left[\left(\prod_{r=1}^{K} \mathbb{P}(e_{i_r})\right)\frac{1}{K}\sum_{r=1}^{K} \mathbb{L}M_{\theta}^{(i)}\left(\cdot\mid e_{i_r},X\right)$$

$$= \sum_{(i_1,...,i_K)\in[v]^K} \left[\left(\prod_{r=1}^{K} \mathbb{P}(e_{i_r})\right)\frac{1}{K}\sum_{r=1}^{K} \mathbb{L}M_{\theta}^{(i)}\left(\cdot\mid e_{i_r},X\right)$$

$$= \sum_{(i_1,...,i_K)\in[v]^K} \mathbb{I}M_{\theta}^{(i)}\left(\cdot\mid e_{i_r},X\right)$$

$$= \sum_{(i_1,...,i_K)\in[v]^K} \mathbb{I}M_{\theta}^{(i)}\left(\cdot\mid e_{i_r},X\right)$$

$$= \frac{1}{K}\sum_{r=1}^{K}\sum_{j=1}^{v} \mathbb{L}M_{\theta}^{(i)}\left(\cdot\mid e_{j_r},X\right) \sum_{(i_1,...,i_{r-1},i_{r+1},...,i_K)\in[v]^{K-1}} \prod_{s\neq r}^{K} \alpha_{i,i_s}^{(MTS)}$$

$$= \frac{1}{K}\sum_{r=1}^{K}\sum_{j=1}^{v} \alpha_{i,j}^{(MTS)} \mathbb{L}M_{\theta}^{(i)}\left(\cdot\mid e_{j_r},X\right) \sum_{(i_1,...,i_{r-1},i_{r+1},...,i_K)\in[v]^{K-1}} \prod_{s\neq r}^{K} \alpha_{i,i_s}^{(MTS)}$$

$$= \frac{1}{K}\sum_{r=1}^{K}\sum_{j=1}^{v} \alpha_{i,j_s}^{(MTS)} \mathbb{L}M_{\theta}^{(i)}\left(\cdot\mid e_{j_s},X\right) = \sum_{j=1}^{v} \alpha_{i,j_s}^{(MTS)} \mathbb{L}M_{\theta}^{(i)}\left(\cdot\mid e_{j_s},X\right)$$

$$= \sum_{j=1}^{v} \alpha_{i,j_s}^{(CoT2)} \mathbb{L}M_{\theta}^{(i)}\left(\cdot\mid e_{j_s},X\right) = \alpha_{i+1}^{(CoT2)}.$$

$$(7)$$

1074 Thus, combining (6), and (7) completes the induction and our argument:

$$\boldsymbol{\alpha}_{t+1}^{(\text{disc})} = \mathbb{E}\left[\hat{\boldsymbol{\alpha}}_{t+1}^{(\text{disc})}\right] = \boldsymbol{\alpha}_{t+1}^{(\text{CoT2})} = \mathbb{E}\left[\hat{\boldsymbol{\alpha}}_{t+1}^{(\text{MTS})}\right] = \boldsymbol{\alpha}_{t+1}^{(\text{MTS})}.$$

Proposition 3. Let α_m be the final expected output distribution after *m* steps of CoT according to Proposition 2. Let $\hat{\alpha}_m$ be the distribution resulting from averaging the outputs of *i.i.d.* CoT2-MTS traces with parallelism *K*. Then, to guarantee $\|\hat{\alpha}_m - \alpha_m\|_1 \le \epsilon$ with high probability, the total number of samples (traces) required scales as $\Theta\left(\frac{\nu}{K\epsilon^2}\right)$.

Proof. We will utilize the empirical distributions $\hat{\alpha}_t^{(\text{disc})}, \hat{\alpha}_t^{(\text{MTS})}$ that are defined in the proof of Proposition 3 and show that i.i.d. sampling *K* discrete CoT trajectories and averaging their results at the last step is distributionally equivalent to CoT2-MTS using *K* tokens, under Assumption 1. We will first argue the results for m = 2 and then we will show it for any m = t + 1. As discussed in the previous proposition, using the decoupling of trajectories by Assumption 1, the next distribution $\hat{\alpha}_1^{(\text{MTS})}$ when e_{i_1}, \ldots, e_{i_K} are drawn from α_1 is:

$$\hat{\alpha}_{2}^{(\text{MTS})} = \text{LM}_{\theta}^{(1)} \left(\cdot \mid \frac{1}{K} \sum_{r=1}^{K} e_{i_{r}}, X \right)$$

$$\hat{\alpha}_{2}^{(\text{MTS})} = \text{LM}_{\theta}^{(1)} \left(\cdot \mid e_{i_{r}}, X \right)$$

$$= \frac{1}{K} \sum_{r=1}^{K} \text{LM}_{\theta}^{(1)} \left(\cdot \mid e_{i_{r}}, X \right)$$

$$= \frac{\sum_{r=1}^{K} \text{LM}_{\theta}^{(1)} \left(\cdot \mid e_{i_{r}}, X \right)$$

$$= \frac{\sum_{r=1}^{K} \text{LM}_{\theta}^{(1)} \left(\cdot \mid e_{i_{r}}, X \right)$$

$$= \frac{\sum_{r=1}^{K} \frac{\sum_{r=1}^{K} \frac{1}{K} \frac{1}{K}}{K}$$

$$= \frac{\sum_{r=1}^{K} \frac{1}{K} \frac{1}{K}}{K},$$

$$= \frac{\sum_{r=1}^{K} \frac{1}{K} \frac{1}{K}}{K},$$

$$= \frac{\sum_{r=1}^{K} \frac{1}{K} \frac{1}{K}}{K},$$

$$= \frac{1}{K} \sum_{r=1}^{K} \frac{1}{K} \sum_{r=1}^{K} \frac{1}{K} \sum_{r=1}^{K} \frac{1}{K} \sum_{r=1}^{K} \sum_{r=1}^{K}$$

which is the empirical mean of K i.i.d. discrete CoT draws. Thus, under our assumption, drawing K tokens with the MTS approach is distributionally equivalent to combining outcomes of K independent discrete CoT draws. Now, for any t, we know by Assumption 1 that output at step t + 1 does not depend on the previous history of tokens. Because of this, we only focus on the tokens e_{i_1}, \ldots, e_{i_k} drawn at step t from $\alpha_t^{(MTS)}$, where $\alpha_t^{(MTS)} = \alpha_t^{(disc)}$ according to Proposition 2. Following the same steps as in (8), we obtain that:

$$\hat{\boldsymbol{\alpha}}_{t+1}^{(\mathrm{MTS})} = \frac{\sum_{r=1}^{K} \mathrm{LM}_{\theta}^{(t)} \left(\cdot \mid \boldsymbol{e}_{i_r}, \boldsymbol{X}\right)}{K} = \frac{\sum_{r=1}^{K} \hat{\boldsymbol{\alpha}}_{t+1,r}^{(\mathrm{disc})}}{K},$$

which is again the empirical mean of K i.i.d. discrete CoT trajectory output token distributions. To finish our argument, we will benefit from a standard result in multinomial estimation that $\Theta\left(\frac{v}{\epsilon^2}\right)$ i.i.d. samples are necessary and sufficient to learn a v-category distribution in $\|\cdot\|_1$ -distance $\leq \epsilon$ (Kamath et al., 2015). In our MTS setting, each sample uses K i.i.d. draws internally. This reduces the variance by a factor of K and achieves the same estimation guarantee with $\Theta\left(\frac{v}{K\epsilon^2}\right)$ aggregated samples. Hence the total sample complexity in terms of MTS samplings is $\Theta\left(\frac{v}{\kappa\epsilon^2}\right)$, as claimed. This completes the argument.

D.1. Construction for Minimum Non-Negative Sum (MNNS) Task

We describe a single-layer transformer with an attention block followed by a mixture-of-experts (MoE) feed-forward block. Let n be the length of the input sequence of integer tokens. Denote the tokenized input numbers as z_1, z_2, \ldots, z_n ; and let the arrow (\rightarrow) token be denoted as z_{n+1} . We also have a dummy input token z_{n+2} , which is the embedding corresponding to number 0, so that we have n + 2 tokens initially. We will construct the transformer with n + 1 MLPs in the mixture of experts layer, where the first n are partial-sum MLPs and the last one is the MLP that reads off the answer from among the all stored partial sums after m steps. We start with the following assumption on the structure of the tokens.

Assumption 2. Let $d = d_e + d_p$ be the embedding size where $d_e = 2^{n+1}$ and $d_p = n + 2$. The token embeddings are on the first d_e coordinates, while the positional encodings are on the last d_p coordinates and are one-hot encoded. where each

 $z_i = \begin{pmatrix} e_i \\ \mathbf{p}_i \end{pmatrix} \in \mathbb{R}^{d_e+d_p}$ is formed by vertically concatenating a content embedding $e_i \in \mathbb{R}^{d_e}$ and a positional encoding $\mathbf{p}_i \in \mathbb{R}^{d_p}$.

We assume each \mathbf{p}_i is a one-hot vector in \mathbb{R}^{d_p} , so that $\mathbf{p}_i^{\mathsf{T}} \mathbf{p}_i = 0$ for $i \neq j$, and $||\mathbf{p}_i|| = 1$.

We now state the following proposition, which helps us to attend and select the input tokens z_1, \ldots, z_{n+1} one by one by the attention block.

Proposition 4. Suppose we have n + 2 tokens $\{z_1, z_2, \ldots, z_{n+2}\}$ in \mathbb{R}^d , each of the form $z_i = \begin{pmatrix} e_i \\ p_i \end{pmatrix}$, where $e_i \in \mathbb{R}^{d_e}, p_i \in \mathbb{R}^{d_e}$.

 \mathbb{R}^{d_p} , $d = d_e + d_p$. Let $p_1, p_2, \ldots, p_{n+2} \in \mathbb{R}^{d_p}$ be orthonormal set of positional vectors according to Assumption 2. Then, there exists a rotation matrix $\mathbf{R} \in \mathbb{R}^{d_p \times d_p}$ satisfying $\mathbf{R}\mathbf{p}_j = \mathbf{p}_{j-1 \mod (n+2)}$ for all $j \in [n+2]$, and the block matrices

$$\boldsymbol{W} = \begin{pmatrix} \boldsymbol{0}_{d_e \times d_e} & \boldsymbol{0}_{d_e \times d_p} \\ \boldsymbol{0}_{d_p \times d_e} & c \cdot \boldsymbol{R} \end{pmatrix} \in \mathbb{R}^{d \times d} \quad and \quad \boldsymbol{W}_{v} = \begin{pmatrix} \boldsymbol{I}_{d_e} & \boldsymbol{0}_{d_e \times d_p} \\ \boldsymbol{0}_{d_p \times d_e} & \boldsymbol{I}_{d_p} \end{pmatrix} \in \mathbb{R}^{d \times d}$$

with $c \rightarrow \infty$, ensure that the attention block

$$\operatorname{Attn}(\boldsymbol{z},\boldsymbol{Z}) = \mathbb{S}\left(\boldsymbol{z}^{\top}\boldsymbol{W}\boldsymbol{Z}^{\top}\right)\boldsymbol{Z}\boldsymbol{W}_{\boldsymbol{v}},$$

performs a cyclic next-index selection: if the query is z_i , it selects column $j^* \equiv (i + 1) \pmod{n + 2}$ from \mathbb{Z} and returns z_{i^*} .

Proof. Definition of Matrix W. We will first construct a rotation matrix. We have n + 2 orthonormal position vectors $p_1, \ldots, p_{n+2} \in \mathbb{R}^{d_p}$. Then, **R** is the following $(n+2) \times (n+2)$ permutation matrix

1147 1148		\int_{0}^{0}	1	0		0	0	
1149 1150	n	0	0 0	1 0	· · · · · · ·	0 0	0	
1151 1152	<i>K</i> =	:	:	:	·	:	:	
1153		$\begin{pmatrix} 0\\ 0 \end{pmatrix}$	0	0	•••	0	$\begin{pmatrix} 1\\ 0 \end{pmatrix}$	

1155 which cyclically shifts the basis vectors p_i backward by one index, i.e., $Rp_i = p_{i-1 \mod (n+2)}$. Then, we specify 1156 $\boldsymbol{W} = \begin{pmatrix} \boldsymbol{0}_{d_e \times d_e} & \boldsymbol{0}_{d_e \times d_p} \\ \boldsymbol{0}_{d_p \times d_e} & c \cdot \boldsymbol{R} \end{pmatrix} \in \mathbb{R}^{d \times d}.$ 1158 1159 Hence for $z_i = (e_i; p_i)$, we have $(e_i^{\top}, p_i^{\top}) W = (0, p_i^{\top} R)$. Thus the dot-product with z_i is 1160 1161 $(0 \mathbf{p}_i^{\mathsf{T}} \mathbf{R}) \begin{pmatrix} \mathbf{e}_j \\ \mathbf{p}_j \end{pmatrix} = \mathbf{p}_i^{\mathsf{T}} \mathbf{R} \mathbf{p}_j.$ 1162 1164 Since positional encodings are orthogonal, we know that: $\boldsymbol{p}_i^{\mathsf{T}} \boldsymbol{R} \boldsymbol{p}_j = \begin{cases} 1, & j \equiv i + 1 \pmod{(n+2)}, \\ 0, & \text{else.} \end{cases}$ So row-wise softmax $S(x^TWX^T)$ places all probability mass at column $i^* \equiv i + 1 \pmod{(n+2)}$ by saturating softmax at 1169 position *j* as $c \to \infty$. 1171 Definition of Matrix W_{y} . In this case, we simply set $W_{y} = I_{d}$, and thus, once the row-wise softmax selects column j^{*} with probability 1, we have 1173 $\boldsymbol{z}_{i^*}^\top \boldsymbol{W}_{v} = \boldsymbol{z}_{i^*},$ 1174 1175 so the final output is precisely the chosen z_{i^*} . This completes the construction. 1176 1177 Having defined the attention block, we state the following proposition that helps selecting different MLPs for the tokens z_1, \ldots, z_{n+1} outputted by the attention block. 1179 1180 Having defined the attention block, we now show how a mixture-of-experts layer can exclusively select MLP_i for each token 1181 z_i , i = 1, ..., n + 1 outputted by the attention block. 1182 **Proposition 5.** Let MLP_1, \ldots, MLP_{n+1} be n + 1 experts in a mixture-of-experts (MoE) module. Suppose we have n + 1 fixed 1183 token embeddings $\{z_1, z_2, \ldots, z_{n+1}\} \subset \mathbb{R}^d$, where each token is formed according to Assumption 2. Given routing parameters 1184 $W = [w_1 \ldots w_{n+1}]^{\mathsf{T}}$, define the MoE feed-forward block as 1185 MoEBlock(z) = $\sum_{i=1}^{n+1} \left[\text{Softmax}(Wz)_j \cdot MLP_j(z) \right],$ 1186 1187 1188 1189 where 1190 Softmax $(Wz)_j = \frac{\exp\left(\mathbf{w}_j^\top z\right)}{\sum_{k=1}^{n+1} \exp\left(\mathbf{w}_k^\top z\right)}, \quad j = 1, \dots, n+1.$ 1192 1193 There exist routing matrix $\mathbf{W} \in \mathbb{R}^{(n+1)\times d}$ such that the distribution $\operatorname{Softmax}(c \cdot \mathbf{W}\mathbf{z}_i)$ as $c \to \infty$ assigns a weight of 1 on MLP_i 1194 when z_i is given as input. 1196 *Proof.* We partition w_j to ignore the content embedding e_i and match the positional block p_j . Concretely, write $w_j = \begin{pmatrix} \mathbf{0}_{d_e} \\ \mathbf{p}_j \end{pmatrix}$. Then, for each token $z_i = (e_i; p_i)$, 1199 1200 $\boldsymbol{w}_{j}^{\top}\boldsymbol{z}_{i} = \begin{pmatrix} \boldsymbol{0}_{d_{e}}^{\top} & \boldsymbol{p}_{j}^{\top} \end{pmatrix} \begin{pmatrix} \boldsymbol{e}_{i} \\ \boldsymbol{p}_{i} \end{pmatrix} = \boldsymbol{p}_{j}^{\top}\boldsymbol{p}_{i}.$ Since $p_i^{\top} p_i = \delta_{ij}$, we have $w_j^{\top} z_i = \delta_{ij}$. Therefore, the softmax evaluates to 1203 1204 1205

 $\lim_{c \to \infty} \operatorname{Softmax}(c \cdot W z_i)_j \to \frac{\exp(\delta_{ij})}{\sum_{k=1}^{n+1} \exp(\delta_{ik})} = \delta_{ij}.$

1207 In other words, Softmax $(c \cdot Wz_i)$ places all mass on expert j = i. Thus each token z_i (for i = 1, ..., n + 1) deterministically selects the *i*-th expert MLP_{*i*}. 1209

In the next proposition, we show how to iteratively expand the partial sums by adding and subtracting the digit obtained 1211 from the attention block and write each resulting sum to a distinct spot in the output vector. 1212 **Proposition 6** (Partial-Sum MLPs). Suppose that the embedding dimension d satisfies $d \ge 2^{j+1} + d_p$. Let z_{prev} contain the 1213 2^{j-1} partial sums s_k each encoded by a pair $(\cos(\omega s_k), \sin(\omega s_k))$ of coordinates such that: 1214 1215 $z_{prev} = \left[\cos(\omega s_1) \sin(\omega s_1) \dots \cos(\omega s_{2^{j-1}}) \sin(\omega s_{2^{j-1}}) 0 \dots 0\right]^{\mathsf{T}} \in \mathbb{R}^d,$ 1216 1217 and let z_{curr} contain the input digit d_i encoded in the first two coordinates: 1218 1219 $z_{curr} = \left[\cos(\omega d_i) \sin(\omega d_i) 0 \dots 0\right]^\top \in \mathbb{R}^d.$ 1220 Then, for any $1 \leq j \leq n$, there exist $MLP_j : \mathbb{R}^d \times \mathbb{R}^d \to \mathbb{R}^d$ such that when (z_{prev}, z_{curr}) is given as input, it outputs the vector 1221 $z_{out} \in \mathbb{R}^d$ so that its first 2^j coordinate-pairs store the trigonometric encodings of $(s_k + d_j)$, and the next 2^j coordinate-pairs 1223 store those of $(s_k - d_i)$. Formally, first 2^j coordinates are $[\cos(\omega(s_k + d_i)), \sin(\omega(s_k + d_i))]$ for all partial sums s_k , and the 1224 next 2^{j} coordinates are $[\cos(\omega(s_{k} - d_{j})), \sin(\omega(s_{k} - d_{j}))]$ for all partial sums s_{k} , with any remaining coordinates set to zero. 1225 1226 *Proof.* Each expert MLP_i (for $1 \le j \le n$) adds j-th integer d_i in both its positive and negative form to all previously computed partial sums. For simplicity, let's say that j-th integer to add is d_j . By trigonometric identities, we know that 1229 $\cos(\omega(s_k + d_j)) = \cos(\omega s_k)\cos(\omega d_j) - \sin(\omega s_k)\sin(\omega d_j),$ 1230 $\sin(\omega(s_k + d_i)) = \sin(\omega s_k)\cos(\omega d_i) + \cos(\omega s_k)\sin(\omega d_i),$ 1231 1232 and similarly, 1234 $\cos(\omega(s_k - d_i)) = \cos(\omega s_k)\cos(\omega d_i) + \sin(\omega s_k)\sin(\omega d_i),$ $\sin(\omega(s_k - d_i)) = \sin(\omega s_k)\cos(\omega d_i) - \cos(\omega s_k)\sin(\omega d_i).$ 1236 1237 Using the above identities, we will obtain the sum by introducing matrices that do shift/swap operations. Concretely, for $k = 1, ..., 2^m$, the k-th 2 × 2 block acts on $\binom{\cos(\omega s_k)}{\sin(\omega s_k)}$ in z_{prev} . We define: 1239 1240 1241 $W_{\sin}^{+} = \operatorname{diag}\left[\underbrace{\begin{pmatrix} 0 & -1\\ 1 & 0 \end{pmatrix}, \dots, \begin{pmatrix} 0 & -1\\ 1 & 0 \end{pmatrix}}_{n+1}, 0, \dots, 0\right],$ 1243 1244 1245 $W_{\sin}^{-} = \operatorname{diag}\left[\underbrace{\begin{pmatrix} 0 & 1 \\ -1 & 0 \end{pmatrix}, \dots, \begin{pmatrix} 0 & 1 \\ -1 & 0 \end{pmatrix}}, 0, \dots, 0\right].$ 1247 1248 1249 1251 The above constructions of W_{\sin}^+ and W_{\sin}^- satisfy, 1252 1253 $\boldsymbol{W}_{\sin}^{+}\boldsymbol{z}_{\text{prev}} = \left[-\sin(\omega s_{1})\,\cos(\omega s_{1})\,\cdots\,-\,\sin(\omega s_{2^{j-1}})\,\cos(\omega s_{2^{j-1}})\,0\,\ldots\,0\right]^{\top} \in \mathbb{R}^{d}$ 1254 and 1256 $\boldsymbol{W}_{\sin}^{-}\boldsymbol{z}_{\text{prev}} = \left[\sin(\omega s_{1}) - \cos(\omega s_{1}) \dots \sin(\omega s_{2^{j-1}}) - \cos(\omega s_{2^{j-1}}) 0 \dots 0\right]^{\mathsf{T}} \in \mathbb{R}^{d}.$ 1257 1258 1259 Each of these act blockwise on the first 2^{j} coordinates of z_{prev} and zeroes out everything else in dimension d. We also have 1260 $z_{\text{curr}} \in \mathbb{R}^d$ with two designated coordinates $z_{\text{curr},1} = \cos(\omega d_j)$, and $z_{\text{curr},2} = \sin(\omega d_j)$, with all other coordinates being zero. We multiply z_{prev} by $\cos(\omega d_j)$ and $\sin(\omega d_j)$ elementwise. Formally, the sum $z_{\text{curr},1} \cdot z_{\text{prev}} + z_{\text{curr},2} \cdot (M_{\text{sin}}^+ z_{\text{prev}})$ 1264 23

gives the 2^{j-1} partial sums $\{s_k + d_j\}_{k=1}^{2^{j-1}}$ stored in the coordinates from 1 to 2^j . We define $W_{\text{shift}} \in \mathbb{R}^{d \times d}$ in a block form with three row blocks and two column blocks: 1265 1266 1267 $W_{\text{shift}} = \begin{pmatrix} \mathbf{0}_{2^{j} \times 2^{j}} & \mathbf{0}_{2^{j} \times (d-2^{j})} \\ I_{2^{j}} & \mathbf{0}_{2^{j} \times (d-2^{j})} \\ \mathbf{0}_{(d-2^{j+1}) \times 2^{j}} & \mathbf{0}_{(d-2^{j+1}) \times (d-2^{j})} \end{pmatrix}.$ 1268 1270 1271 When applied, the above matrix shifts the first 2^{j} entries of z_{prev} by 2^{j} coordinates. Now, also define 1272 $z_{\text{curr},2} \cdot (W_{\text{shift}} W_{\text{sin}}^{-} z_{\text{prev}}) + z_{\text{curr},1} \cdot (W_{\text{shift}} z_{\text{prev}}).$ 1273 1274 This way, the above sum gives us the 2^{j-1} partial sums $\{s_k - d_j\}_{k=1}^{2^{j-1}}$ stored in the coordinates from $2^j + 1$ to 2^{j+1} encoded in trigonometric format. Then, we normalize this output of the model by 1/2 and obtain the following output: 1275 1276 $\left(z_{\text{curr},1} \cdot z_{\text{old}} + z_{\text{curr},2} \cdot \left(M_{\sin}^{+} z_{\text{old}}\right) + z_{\text{curr},2} \cdot \left(W_{\text{shift}} W_{\sin}^{-} z_{\text{prev}}\right) + z_{\text{curr},1} \cdot \left(W_{\text{shift}} z_{\text{prev}}\right)\right)$ $= \left[\cos\left(\omega\left(s_1+d_j\right)\right), \sin\left(\omega\left(s_1+d_j\right)\right), \dots, \cos\left(\omega\left(s_{2^{j-1}}+d_j\right)\right), \sin\left(\omega\left(s_{2^{j-1}}+d_j\right)\right), \dots\right]\right]$ 1279 1280 $\cos\left(\omega\left(s_{1}-d_{i}\right)\right), \sin\left(\omega\left(s_{1}-d_{i}\right)\right), \ldots, \cos\left(\omega\left(s_{2^{j-1}}-d_{i}\right)\right), \sin\left(\omega\left(s_{2^{j-1}}-d_{i}\right)\right), \sin\left(\omega\left(s_{2^{j-1}}-d_{i}$ 1281 $[0,\ldots,0]^{\top} \in \mathbb{R}^d.$ 1282 1283 Thus, this is exactly the representation of 2^{j} partial sums. This completes the argument. We should remark that, the above 1284 argument utilizes a gated MLP which explicitly multiplies the elements of the input features, namely, z_{curr} with the partial 1285 sums z_{prev} . On the other hand, we don't require any nonlinear activation function, so our MLP constructions have the form 1286 $MLP(z) = W_3(W_1z \odot W_2z)$ for suitable choices of W_1, W_2, W_3 where \odot denotes the Hadamard product. The use of gated 1287 MLPs is a standard practice in transformer architectures (Shazeer, 2020). 1288 1289 **Proposition 7** (Read-Off MLP). Suppose that every partial sum s_k is in the range [-S, S] and let $\omega < \pi/2S$. Assume that 1290 the vector $z = [\cos(\omega s_1), \sin(\omega s_1), \ldots, \cos(\omega s_{2^n}), \sin(\omega s_{2^n}), 0, \ldots, 0]^\top \in \mathbb{R}^d$ 1292 contains 2^n partial sums $\{s_1, \ldots, s_{2^n}\}$ encoded in trigonometric form, where $d = 2^{n+1} + n + 2$. Then there exists a single 1293 feed-forward network $MLP_{n+1} : \mathbb{R}^d \to \mathbb{R}^d$ such that, given input z, it selects the smallest nonnegative s_ℓ from $\{s_1, \ldots, s_{2^n}\}$ 1294 and outputs the embedding $e_{s_{\ell}} \in \mathbb{R}^{d}$, where s_{ℓ} is that minimal nonnegative partial sum. 1296 Remark: Our construction relies on gated MLP, rather than standard MLP, as in Proposition 6. 1298 *Proof.* We know that the input embedding z represents 2^n pairs, each pair $(\cos(\omega s_i), \sin(\omega s_i))$ stored consecutively. That is, 1299 1300 $z = [\cos(\omega s_1), \sin(\omega s_1), \ldots, \cos(\omega s_{2^n}), \sin(\omega s_{2^n}), 0, \ldots, 0]^\top \in \mathbb{R}^d,$ 1302 We will identify the smallest $s_{\ell} \ge 0$ and output an embedding $e_{s_{\ell}}$ denoting that integer. We are given that ω is small enough such that when $s_{\ell} \in [0, S]$, we ensure $S\omega < \pi/2$. This guarantees $\sin(\omega s_{\ell}) \ge 0$ if and only if $s_{\ell} \ge 0$. First, we wish to 1303 collapse z into a single vector of size 2^n , keeping $\cos(\omega s_\ell)$ only when $\sin(\omega s_\ell) \ge 0$ and zeroing it out otherwise. We define 1304 two matrices $W_{cos}, W_{sin} \in \mathbb{R}^{d \times d}$ by 1305 1306 $(\mathbf{W}_{\cos})_{i,(2i-1)} = 1, \quad (\mathbf{W}_{\cos})_{i,j} = 0 \text{ for } j \neq 2i - 1,$ 1307 $(\mathbf{W}_{\sin})_{i,(2i)} = 1, \quad (\mathbf{W}_{\sin})_{i,j} = 0 \text{ for } j \neq 2i.$ 1308 1309 for $1 \le i \le 2^n$ and all other rows/columns of W_{sin} , W_{cos} are zero. Hence each matrix picks out alternate coordinates: 1310 $z_{\cos} = W_{\cos} z = \begin{vmatrix} \cos(\omega s_1) \\ \cos(\omega s_2) \\ \vdots \\ \cos(\omega s_{2^n}) \\ 0 \\ \vdots \end{vmatrix} \in \mathbb{R}^d, \quad z_{\sin} = W_{\sin} z = \begin{vmatrix} \sin(\omega s_1) \\ \sin(\omega s_2) \\ \vdots \\ \sin(\omega s_{2^n}) \\ 0 \\ \vdots \end{vmatrix} \in \mathbb{R}^d.$ 1312 1313 1314 1317 1318 1319

1320 In order to find the minimum non-negative number, we need to find the number *s* such that it maximizes $\cos(\omega s)$ and satisfies 1321 $\sin(\omega s) \ge 0$. For this, we utilize a sigmoid activation function in the following way:

$$\mathbf{z}_{\text{filter}} = \mathbf{z}_{\cos} \odot \sigma \left(c \, \mathbf{z}_{\sin} \right),$$

where $\sigma(x) = \frac{1}{1 + \exp(-x)}$ is element-wise sigmoid function, and $c \to \infty$ is a large constant. With this choice of *c*, the sigmoid output will be 1 when $s_{\ell} \ge 0$ and 0 otherwise. Therefore, the resulting vector z_{filter} contains $\cos(\omega s)$ values at indices where sin(ωs) is positive. Now, for $0 \le s_{\ell} \le S$ with $S \omega \le \frac{\pi}{2}$, the ordering of s_{ℓ} from smallest to largest is the same as the ordering of $\cos(\omega s_{\ell})$ from largest to smallest. Thus, to find the minimum nonnegative sum, we find the partial sum ℓ^* that maximizes $\cos(\omega s_{\ell})$. Utilizing another gating, we calculate

Softmax
$$(c z_{\text{filter}})^{\top} z_{\text{filter}}$$

as $c \to \infty$. The softmax vector will be one-hot with 1 at index ℓ^* that has the largest $\cos(\omega s_{\ell})$. A second multiplication with z_{filter} will return this $\cos(\omega s_{\ell^*})$. Therefore, Softmax $(c z_{filter})^{\top} z_{filter} = \cos(\omega s_{\ell^*})$. Next, we retrieve the corresponding sine entry of s_{ℓ^*} by applying the same one-hot selection to z_{sin} . Formally,

Softmax
$$(c z_{\text{filter}})^{\top} z_{\sin} = \sin(\omega s_{\ell^*})$$

1337 as $c \to \infty$. Hence, from these two selected coordinates, $[\cos(\omega s_{\ell^*}), \sin(\omega s_{\ell^*})]$, we produce the final embedding in \mathbb{R}^d by 1338 placing them in the first two coordinates and zeros elsewhere:

$$e_{s_{\ell^*}} = [\cos(\omega s_{\ell^*}), \sin(\omega s_{\ell^*}), 0, \dots, 0]^\top,$$

¹³⁴¹ where s_{ℓ^*} is the minimal nonnegative sum. This completes the argument.

Proposition 1 (Solving MNNS). There exists a 1-layer transformer architecture with a mixture-of-experts MLP layer that solves the MNNS task using CoT2 by storing (sine, cosine) embeddings of all 2^k states at the k-th iteration in a non-overlapping manner.

Proof. We will argue that by combining Propositions 5 to 7, we obtain a single-layer transformer that is formed by an attention block followed by an MoE feed-forward block, which solves the Minimum Non-Negative Sum (MNNS) task.

1349 Suppose that we have n input integers d_1, \ldots, d_n , encoded as z_1, \ldots, z_n , plus an arrow (\rightarrow) token z_{n+1} and a dummy token 1350 z_{n+2} corresponding to the integer 0. In this case, we will output the tokens representing the ground-truth sums s_1, \ldots, s_n , 1351 therefore, the number of output tokens is m = n in MNNS setting. We assume that the inputs are encoded according to 1352 Assumption 2. By Proposition 6, there exist MLP_1, \ldots, MLP_n that perform the following: whenever MLP_j is selected with 1353 input $(z_{\text{prev}}, z_{\text{curr}})$ such that z_{prev} stores 2^{j-1} partial sums and z_{curr} stores the digit d_j , it adds and subtracts d_j to all previously 1354 stored partial sums and stores the resulting 2^{j} partial sums in z_{out} . The dummy token z_{n+2} that corresponds to the integer 0 1355 allows us to initialize the partial sums from zero. If the query token is z_{n+2} , we produce the first partial sums by combining 1356 this dummy 0 with d_1 , which are $(+d_1)$ and $(-d_1)$ encoded in an output token. 1357

We assign positional encodings cyclically to output tokens. That means, the first n + 2 input tokens have positional encodings from p_1 to p_{n+2} , and the output tokens have $p_1, p_2, ..., as$ their positional encodings, in this exact order. This way, by Proposition 4, Attn(z, Z) attends and selects the input digit tokens $z_1, z_2, ..., z_n$ and finally arrow z_{n+1} one by one and feeds to MoEBlock(·).

1362 By Proposition 5, there's a MoEBlock(z) such that if the input is z_i (for $j \le n$), MLP_i is selected with probability 1, and if 1363 the input is arrow token z_{n+1} , MLP_{n+1} is selected with probability 1, which is the MLP to read-off the final answer. In the 1364 input tokens z_1, \ldots, z_n , the first two coordinates store the trigonometric representation of d_1, \ldots, d_n . To allow outputting 1365 the final answer by MLP_{n+1} , the partial sums obtained in the intermediate steps need to be written to separate coordinates. 1366 Therefore, MLP_i takes a vector filled in the first 2^{j} coordinates, adds d_{i} and writes to the first 2^{j} coordinates, subtracts 1367 d_i and writes to the next 2^j coordinates, and finally divides the entire representation by 2 to maintain consistent scaling 1368 since the number of partial sums is doubled. In other words, the first *n* MLPs have some repeated behavior. Finally, by 1369 Proposition 7, MLP_{*n*+1} receives a vector that encodes all 2^n possible partial sums in cos/sin form in 2^{n+1} coordinates and 1370 extracts the embedding of smallest nonnegative number among them. 1371

1372Altogether, this single-layer transformer with an attention module to pass the tokens to the mixture-of-experts MLP solves1373the Minimum Non-Negative Sum task by following CSFT described in 3.

1374

1322 1323

1330

1335

1336

1339 1340

1346