# INTERPRETABLE ANALYSIS AND REASONING EN HANCEMENT FOR LLMS VIA CROSS-GENERATION REASONING TREE

Anonymous authors

Paper under double-blind review

#### ABSTRACT

Generating diverse reasoning paths by varying the context, such as demonstrations, prompts, instructions, etc, or sampling methods, such as top-k, top-p, beamsearch, etc, and then selecting appropriate paths via majority voting or verifierbased strategies to enhance the reasoning capabilities of large language models (LLMs) is a commonly recognized approach. Although both different contexts and sampling techniques can generate diverse contents, using sampling methods alone does not significantly enhance the diversity of generation. Context variation, however, while fostering greater diversity in reasoning, can also introduce negative effects. It causes that switching contexts can not necessarily lead to proportional improvements in performance. Therefore, there is a need to investigate how context influences LLM generation and mitigate any adverse impacts. The primary challenge lies in the inability to conduct comparative studies once divergences occur in reasoning paths generated under different contexts. Specifically, once the predicted tokens at a given step differ, it becomes unclear whether subsequent tokens in the inference path are influenced by the context or the content already generated. In this paper, we propose a Cross-Generation Reasoning Tree (CGRT) algorithm for studying the impact of different contexts on LLM generation and enhancing LLMs' reasoning performance. Experimental findings reveal that, beyond enhancing interpretability, CGRT integrates the positive effects of both context and sampling strategies more effectively than previous approaches, leading to more rational inference paths. Experiments conducted on Llama2, Llama3, and Qwen demonstrate that, when generating an equivalent number of diverse inference paths, those produced via the "reasoning tree" method exhibit higher accuracy.

035 036 037

038

039

006

008 009 010

011 012 013

014

015

016

017

018

019

021

023

025

026

027

028

029

031

032

033

034

#### 1 INTRODUCTION

Large Language Models (LLMs) have emerged with impressive performance in generative reasoning, often approaching that of human experts. OpenAI et al. (2023); Anil et al. (2023); Dubey et al. (2024); cla (2024) However, when confronted with more complex problems, LLMs may still make unexpected mistakes. A commonly recognized solution to this issue is generating diverse inference paths and then deciding the final answer via appropriate judgment strategies, such as "selfconsistency" Wang et al.; Aggarwal et al. (2023) or "best-of-N" Liu et al. (2020), *etc.* 

Typically, diversifying generated content can be achieved through appropriate sampling methods Fan et al. (2018); Holtzman et al.; Freitag & Al-Onaizan (2017); Chuang et al. (2023). Nevertheless, the diversity of sampling via decoding strategies remains constrained by the model's own distribution for answering the query question, which means that the diversity of sampled paths is limited when the model is highly confident. Besides, varying the context can also produce more diverse reasoning. However, inappropriate contexts can overly change the model's original distribution, sometimes resulting in significant negative effects. Therefore, it is necessary to investigate how context influences LLM generation and to correct any adverse impacts, which is a key topic that requires urgent attention in the field of LLM reasoning. 054 Often, conditioned with different contexts, once the LLMs' generation diverges at a certain token, 055 the subsequent generated tokens cannot be compared. Therefore, previous interpretability analysis 056 methods can only study the first inconsistent token or a few unexpected tokens but cannot track 057 the entire reasoning process, *i.e.*, they cannot analyze the entire generation process under different 058 contextual scenarios. Starting from the first inconsistent token, the subsequent predicted tokens are influenced not only by the context but also by the previously generated contents. To address this issue, we propose the "Cross-Generation Reasoning Tree" (CGRT). For each reasoning problem, 060 CGRT constructs the LLMs generation into a tree structure. Each node in the tree represents a 061 token generated by the LLM at a certain step. CGRT determines the next step of generation by 062 placing different contexts before and after the problem, *i.e.* few-shot chain-of-thought (CoT) Wei 063 et al. (2022) demonstrations, prompts Kojima et al. (2022); Sahoo et al. (2024), instructions Efrat & 064 Levy (2020); Mishra et al. (2022), etc, and then selects the child nodes of the current step node via 065 certain sampling strategies, such as top-k, top-p, beam-search decoding, etc. When the generation 066 at the current step diverges, the tree structure branches out. For each branch path of the tree, when 067 determining the child nodes of a certain node, *i.e.*, the next step of generation by the LLM, the full 068 combination of contexts and sampling strategies is applied. An example of the CGRT is illustrated in Figure 1. Thus, each branching point in the CGRT structure represents a divergence caused solely 069 by the context, as the preceding generated reasoning path is controlled.

Through experiments using the Cross-Generation Reasoning Tree (CGRT), we further confirmed the widely acknowledged fact that context has both positive and negative effects on the LLMs' reasoning. Additionally, we conducted further investigations and found the following:

- 1. In CGRT constructed entirely from good contexts, a significant number of erroneous reasoning paths exist. Similarly, in CGRT constructed entirely from bad contexts, there are also correct reasoning paths.
- 2. The critical branching nodes in CGRT that determine the correctness of reasoning are predominantly composed of tokens with strong semantic information, such as numbers, operators, nouns *etc*. In contrast, words with weaker semantic significance, such as prepositions, conjunctions, punctuation, pronouns *etc*, are often not the critical branching nodes that determine the correctness of reasoning.
- 082 083

075

076

077 078

079

081

where good context refers to the context that, when placed before the question and decoded using
 greedy decoding, results in a correct reasoning answer; bad context, conversely, leads to an incorrect
 reasoning answer under the same conditions. Critical branching nodes represent nodes in a CGRT
 where all paths following the branching point lead to correct (or incorrect) reasoning answers, while
 their sibling nodes lead to the opposite, i.e., all paths following its sibling nodes lead to incorrect (or
 correct) reasoning answers.

However, for real-world reasoning problems, it is challenging to determine which path to take at a 090 branching point. From the perspective of a branching point, although some critical branching nodes 091 determine the correctness of subsequent reasoning paths, this determination is based on the CGRT 092 constructed by the current context combinations. Exhaustively enumerating all possible contexts to generate an exceptionally large CGRT is obviously impractical. Furthermore, from a semantic 094 perspective, most critical branching nodes cannot be fully understood as to why they determine the 095 correctness of reasoning when viewed only from the current step. Human experts would find that, 096 from the token of critical branching nodes that lead to wrong answers, it is possible to reason towards 097 the correct answer. However, the performance of LLMs does not reflect this. The most likely 098 explanation is that LLMs remain biased causal language models. LLMs do not possess genuine reasoning capabilities but rather model human f. Due to factors such as training data and parameter 099 scales, they are unable to model the reasoning embedded within human language completely. 100

To address the challenge of selecting the appropriate path at branching nodes, we propose the inference version of the Cross-Generation Reasoning Tree (CGRT), referred to as iCGRT. iCGRT draws inspiration from the majority voting strategies but applies it at both the token-level and path-level. When studying the interpretability of LLMs, CGRT typically selects a limited number of context combinations (usually 2 - 4), and a small number of sampling tokens for each context combination (typically one token decoded via greedy strategy or 1 - 2 tokens sampled via top-k/p strategy). In contrast, iCGRT selects a larger number of context combinations ( $\leq 8$ ) and usually generates 4 - 8samples for each context combination. At each node, if the generated tokens across all samples



Figure 1: An Example of Cross-Generation Reasoning Tree (CGRT).

129 are inconsistent, a majority voting strategy is employed to select the top 2 predicted tokens as the 130 nodes for that step. This approach has two benefits: firstly, the tokens selected through majority 131 voting reduce the probability of leading to incorrect answers at that node; secondly, by restricting 132 iCGRT to a binary tree, the computational burden of reasoning is reduced. Finally, the reasoning 133 answer is chosen from all the reasoning paths in the iCGRT through a majority voting strategy. Experiments conducted on mainstream open-source LLMs, such as Llama2 Touvron et al. (2023), 134 Llama3 AI@Meta (2024), and Qwen Team (2024), demonstrate that under the same number of 135 context combinations and sampling quantities, iCGRT can produce more accurate reasoning paths. 136

#### 137 138 139

144 145

146

150 151

152

127 128

#### 2 METHODOLOGY

In this section, we first elaborate on the construction algorithm of the Cross-Generation Reasoning
 Tree (CGRT) after defining the mathematical notation. Subsequently, we present the interpretability
 conclusions of LLMs derived from CGRT experiments. Finally, we show the algorithm details of
 the inference-version CGRT, iCGRT.

#### 2.1 MATHEMATICAL NOTATION

147 A N-ary tree can be represented as the pair of node set and edge set  $\mathscr{T} = \{\mathcal{V}, \mathcal{E}\}$ , where  $\mathcal{V}$  is the 148 node set and  $\mathcal{E}$  is the edge set. The element of edge set,  $(\mathbf{u}, \mathbf{v}) \in \mathcal{E}$ ,  $\mathbf{u} \in \mathcal{V}$ ,  $\mathbf{v} \in \mathcal{E}$ , represents that 149 the node  $\mathbf{v}$  is one of the children node of the node  $\mathbf{u}$ . We define the insert operation as:

$$\mathrm{Insert}(\mathscr{T},\,\mathbf{w},\,\mathbf{p})=(\mathcal{V}=\mathcal{V}\cup\{\mathbf{w}\},\;\mathcal{E}=\mathcal{E}\cup(\mathbf{p},\,\mathbf{w}))$$

For a node v, we use the corresponding non-bold italic v to denote the token it represents. With the notation conventions for an N-ary tree established, we define the Cross-Generation Reasoning Tree (CGRT) for the language model. For a query question q and a C is the set of contexts, define the CGRT as  $\mathcal{T}_{q,C} = \{\mathcal{V}, \mathcal{E}, \mathcal{C}\}$ . Each node represents a token generated at a certain step, and each node can have up to  $|\mathcal{C}|$  child nodes, where  $|\mathcal{C}|$  is the number of elements in the context set.

158 159

160

2.2 CROSS-GENERATION REASONING TREE

When generating the  $i^{\text{th}}$  token, the LLMs' prediction distribution is influenced by three factors: context (such as demonstrations, prompts, instructions, *etc.*), the query question, and the content



Figure 2: Figure 1: Illustration of the iCGRT construction algorithm. For each node in the iCGRT, 177 when generating the next step's tokens, it considers  $n_c$  contexts, and takes  $n_s$  samplings for each 178 context, generating  $n_c \cdot n_s$  candidate tokens. If the candidate tokens are not entirely consistent, the 179 top 2 tokens among them are selected as the child nodes of this node, thus limiting iCGRT to a 180 binary tree. In practice, there are additional implementation details of iCGRT not illustrated in this figure; for a full explanation, please refer to section 2.3. 182

181

- 185

187 already generated. Under different contexts or samplings, or a combination of both, when diver-188 gence occurs in two inference paths from a certain step, the subsequent reasoning paths will appear 189 different, although semantically they may convey the same meaning. Previous analyzing methods to 190 seek LLMs' interpretability, such as probing Belinkov (2022), information flow Abnar & Zuidema 191 (2020); Wang et al. (2023a), could only study the first branching position or manually identify locations in the reasoning path where unexpected errors occur. In order to delve deeper into the impact of 192 different contexts on content generation, we propose the CGRT (Cross-Generation Reasoning Tree) 193 structure. 194

195 For each query, given two contexts, CGRT constructs the model's generation into a tree structure. 196 Specifically, CGRT represents each step of the LLM's prediction as a node in the tree, with each branch corresponding to a reasoning path. If the predicted token is <|end of sentence|>, 197 the node is marked as a leaf node. For each step, one of the non-leaf nodes serves as the parent 198 node. If the model predicts the same token for the current step under both contexts, a single child 199 node is inserted for the parent node. If the predictions are inconsistent, two child nodes are inserted, 200 indicating a bifurcation in the tree at this step. This process is recursively iterated until all terminal 201 nodes in the tree are leaf nodes. In the case of binary CGRTs constructed from two contexts, only 202 the paths of the leftmost and rightmost branches are generated based on information from a single 203 context. Most other paths in the CGRT combine information from both contexts, which is the most 204 significant difference from traditional tree structures and is the reason why we name the algorithm 205 as a "Cross-Generation" Tree.

206 The formalized algorithm for constructing CGRT (Cross-Generation Reasoning Tree) is presented 207 in Algorithm 1. It is important to note that the flexibility of language permits the same concept 208 to be articulated in numerous varied manners. Hence, if no limitations are applied to the CGRT, 209 this variability can lead to an extraordinarily high number of branches in some instances, poten-210 tially extending the construction time of a CGRT to several days. In practice, therefore, we opt 211 to disregard branch point tokens that possess lesser semantic significance. The precise definitions 212 of these non-significant tokens are delineated in Appendix A.1.4. If, at a branching point, some 213 branches are determined to be non-significant, these will be disregarded; if all branches are deemed non-significant, a single branch will be chosen at random. However, even disregarding these tokens, 214 the number of branches in CGRT might still grow uncontrollably for certain queries. Therefore, in 215 practical implementation, it may be necessary to impose a limit on the maximum width of the tree.

216 Algorithm 1: CGRT (Cross-Generation Reasoning Tree) 217 **Input:** a language model M, a query question q, the set of contexts C218 **Output:** CGRT  $\mathscr{T}_{q,\mathcal{C}} = \{\mathcal{V}, \mathcal{E}\}$ 219 1 Initialize: a non-leaf root node **r** with empty data,  $\mathscr{T}_{q,\mathcal{C}} \leftarrow \{\mathcal{V} = \{\mathbf{r}\}, \mathcal{E} = \varnothing\}$ 220 <sup>2</sup> while there are non-leaf nodes without child nodes in the tree  $\mathcal{T}_{a,C}$  do 221 Select a non-leaf node  $\mathbf{v} \in \mathcal{V}, v \notin \mathcal{W}_{eos}$  without child nodes 3 222 Get the path from the root node  $\mathbf{r}$  to node  $\mathbf{v}$ : Sequence( $\mathbf{v}$ ) 4 For each context  $c_i \in C$ , given the input  $c_i + q + \text{Sequence}(\mathbf{v})$ , obtain the model M's 5 224 predicted token of the next step:  $u_i = M(c_i + q + \text{Sequence}(\mathbf{v}))$ 225 Remove duplicates from the predicted tokens across different contexts, resulting in j6 (where  $1 \le j \le |\mathcal{C}|$ ) unique tokens: 226 227  $\{w_1, w_2, \ldots, w_j\} = \operatorname{Set}(u_1, u_2, \ldots, u_{|\mathcal{C}|})$ 228 229 for  $w \in \{w_1, w_2, \ldots, w_i\}$  do 7 230 if  $w \in \mathcal{W}_{eos}$  then 8 The node w is marked as a leaf node. 231 9 end 10 else 11 233 The node w is marked as a non-leaf node. 12 end 13 235 Insert the node w as a child node of the node v: 14 236  $\mathscr{T}_{q,\mathcal{C}} \leftarrow \operatorname{Insert}(\mathscr{T}_{q,\mathcal{C}},\mathbf{v},\mathbf{w}) = (\mathcal{V} \cup \{\mathbf{w}\}, \mathcal{E} \cup \{(\mathbf{v},\mathbf{w})\})$ 237 238 end 15 239 16 end 240 241 242 243 2.3 INFERENCE-VERSION CGRT 244 245 Theoretically, after obtaining a complete CGRT, an advanced selector could be employed to make 246 247 248 249

choices at the branching nodes of the tree, thereby navigating the correct reasoning path. However, in practical application, it poses an almost insurmountable challenge. Firstly, at the step of branching nodes, although some branching nodes within a given CGRT may lead to entirely correct or incorrect subsequent reasoning paths, it is a phenomenon limited to the current CGRT. Given that 250 CGRTs can vary infinitely due to the countless combinations of contexts, exhaustively enumerating 251 all possible contexts to construct an extraordinarily large CGRT is impractical. Secondly, from a semantic perspective, most critical branching nodes cannot be discerned by human experts to de-253 termine what factors decide the correctness or incorrectness of subsequent reasoning paths. Often, even when a branching node leads to all incorrect reasoning paths within a certain CGRT, human 254 experts are still capable of reasoning correctly from the current step. We enumerate many such cases 255 in Appendix A.3. Numerous studies Jin et al.; Valmeekam et al. (2024); Kambhampati (2024) indi-256 cate that LLMs lack genuine reasoning capabilities - they merely model human language, and this 257 modeling is biased. Due to factors such as training data and parameter scale, LLMs are unable to 258 fully model the reasoning embedded within human language. This explains why the performance of 259 LLMs diverges from that of human experts, influenced as it is by contextual factors. Consequently, 260 making the right decision at the branching nodes of a CGRT, viewed solely from the step at the 261 branching node, remains a challenging task. 262

Therefore, during the practical inference process, a solution to this problem is required. To address
 this, we propose the inference version of CGRT, referred to as iCGRT (inference-version ). The
 algorithm of building iCGRT is elaborated on the Algorithm 2. Compared to the standard building
 CGRT algorithm, iCGRT introduces the following modifications to adapt to the inference demands:

267 268

1. Employing a larger number of context combinations and more diverse sampling than used in interpretability studies to avoid biases caused by insufficient sampling. This adjustment is reflected in Line 5 of Algorithm 2. 270 Algorithm 2: iCGRT (inference-version Cross-Generation Reasoning Tree) 271 **Input:** a language model M, a query question q, the set of contexts C, the decoding function S, 272 the number of sampling steps per context  $n_s$ , the maximum number of iCGRT branches 273  $n_p$ 274 **Output:** iCGRT  $\mathscr{T}_{q,\mathcal{C},S}^{(i)} = \{\mathcal{V}, \mathcal{E}\}$ 275 1 Initialize: a non-leaf root node **r** with empty data,  $\mathscr{T}_{q,\mathcal{C},S}^{(i)} \leftarrow \{\mathcal{V} = \{\mathbf{r}\}, \mathcal{E} = \varnothing\}$ 276 277 **2 while** there are non-leaf nodes without child nodes in the tree  $\mathcal{T}_{q,\mathcal{C},S}^{(i)}$  **do** 278 Select a non-leaf node  $\mathbf{v} \in \mathcal{V}, v \notin \mathcal{W}_{eos}$  without child nodes 3 279 Get the path from the root node r to node v: Sequence(v) 4 For each context  $c_i \in C$ , given the input  $c_i + q + \text{Sequence}(\mathbf{v})$ , obtain the model M's  $n_s$ 5 281 predicted tokens of the next step via the decoding method S:  $\mathcal{U} \leftarrow \{u_{i,1}, \ldots, u_{i,n_s}\} = S(M(c_i + q + \text{Sequence}(\mathbf{v})), n_s)$ 283 Select the top 2 tokens from  $|C| \cdot n_s$  predicted tokens: 284 6  $\{w_1, w_2\} = \text{Top-}2(u_{1,1}, \dots, u_{1,n_s}, u_{2,1}, \dots, u_{2,n_s}, u_{|\mathcal{C}|,1}, \dots, u_{|\mathcal{C}|,n_s})$ 287 if  $|\{w \in \mathcal{U} \mid w = w_1\}| > \lambda |\{w \in \mathcal{U} \mid w = w_2\}|$  or PathNumber $(\mathcal{T}_{a, \mathcal{C}, S}^{(i)}) \ge n_p$  then 7  $\mid \mathcal{W}_{\text{kept}} \leftarrow \{w_1\}$ 8 289 end 290 10 else 291  $| \mathcal{W}_{\text{kept}} \leftarrow \{w_1, w_2\}$ 11 292 end 12 293 for  $w \in \mathcal{W}_{kept}$  do 13 if  $w \in \mathcal{W}_{eos}$  then 14 The node w is marked as a leaf node. 295 15 16 end 296 else 17 297 The node w is marked as a non-leaf node. 18 298 19 end 299 Insert the node w as a child node of the node v: 20 300 301  $\mathscr{T}_{a\,\mathcal{C}\,S}^{(i)} \leftarrow \operatorname{Insert}(\mathscr{T}_{a\,\mathcal{C}\,S}^{(i)}, \mathbf{v}, \mathbf{w}) = (\mathcal{V} \cup \{\mathbf{w}\}, \, \mathcal{E} \cup \{(\mathbf{v}, \mathbf{w})\})$ 302 end 303 21 22 end 305 306 307 2. Introducing a majority voting strategy during the inference at each node (each step gener-308 ated by LLMs). This serves two purposes: firstly, it maximizes the avoidance of undesirable token generation; secondly, since a large number of tokens (often  $\leq 32$ ) are generated per 310 step, leading to an overly extensive tree structure, retaining only the top-2 tokens limits the tree structure to a binary tree, thus significantly reducing the number of reasoning paths in 311 the iCGRT. This modification is reflected in Line 6 of Algorithm 2. 312 313 3. Building upon the previous improvement, if the quantity of the top-1 token far exceeds 314 that of the top-2 token (for example, if the top-1 token count is more than 4 times that 315 of the top-2 token), we consider it as a non-branching node. It is controlled by the input hyperparameter  $\lambda$  and the first condition in Line 6 of Algorithm 2. 316 317 4. iCGRT imposes limitations on the size of the tree to prevent issues arising from excessive 318 branching in rare cases. Specifically, when the total number of paths in the tree exceeds 319 a predefined maximum, only the top-1 token is selected during node generation. This 320 limitation is governed by the input hyperparameter  $n_p$  and the second condition in Line 6 of Algorithm 2. 321 322 In addition to the methods mentioned above for accelerating and preventing excessive computational 323 costs, iCGRT can also process only half the number of tokens at each prediction step to approxiTable 1: *maj-vote acc.* represents the accuracy of the predicted answer via majority voting strategy, and *prop. of correct paths* represents the proportion of correct reasoning paths among the total reasoning paths. The experiment uses the Llama3-8B model on the GSM8K test set via the 1-shot CoT setting. For each context, the CGRT generates predictions via greedy decoding.

num of context	2			4	8		
matric	maj-vote	prop. of	maj-vote	prop. of	maj-vote	prop. of	
meure	acc.	correct paths	acc.	correct paths	acc.	correct paths	
hard level 1	20.83	14.57	23.31	15.28	28.65	19.81	
hard level 2	0.00	1.34	7.54	3.63	8.01	3.34	
hard level 3	0.00	2.75	6.89	3.22	5.04	3.19	
hard level 4	0.00	2.65	3.01	4.13	2.89	1.34	

mately to accelerate the inference. Specifically, in Step 6 of Algorithm 2, only half of the  $n_s \cdot n_c$  tokens that originally should be generated are predicted. If the first half  $n_s \cdot n_c/2$  tokens are consistent or the number of top-1 tokens far exceeds the number of top-2 tokens, the generation of the remaining half is skipped. In practice, there is a considerable number of prediction steps where all  $n_s \cdot n_c$  tokens are consistent. This approach can reduce the computational load by approximately 1/3. Throughout this paper, unless otherwise specified, all iCGRT experiments employ this method to accelerate inference.

#### 3 EXPERIMENTS

#### 3.1 MINING CORRECT REASONING PATHS FROM BAD CONTEXT

By constructing the CGRT, we found that even when using bad contexts, *i.e.*, contexts that cause LLMs to answer query questions incorrectly, the CGRT still generates some correct reasoning paths. To quantify this phenomenon, we first define four criteria to measure the difficulty levels of samples:

- A. Randomly select N combinations of demonstrations, and generate one inference path per combination using greedy decoding. The answer derived from the majority voting of these N inference paths is incorrect.
  - B. In addition to satisfying criterion A, all N inference paths generated must have incorrect answers.
- C. From the N combinations of demonstrations mentioned in criteria A and B, randomly select M (where  $M \leq N$ ) combinations. Generate S inference paths for each of these M combinations using top-k or top-p sampling. The answer derived from the majority voting of the  $M \cdot S$  generated inference paths is incorrect.
  - D. In addition to satisfying criterion C, all  $M \cdot S$  generated inference paths must have incorrect answers.

With these four criteria, we categorize more challenging samples into four levels, ordered from less difficult to most difficult:

- 1. Hard Level 1 (Not Very Hard): Meets criterion A
- 2. Hard Level 2 (Normal Hard): Meets criterion B
  - 3. Hard Level 3 (Very Hard): Meets both criteria B and C
  - 4. Hard Level 4 (Extremely Hard): Meets both criteria B and D

Conventional self-consistency/majority voting strategies do not effectively address hard-level samples arbitrarily. Furthermore, criteria B and D involve even stricter constraints on the difficulty, requiring that none of the generated paths are correct. Therefore, for samples classified within Hard
 Levels 2-4, particularly those at Hard Level 4, it is challenging for models to produce correct inferences based on ordinary diverse sampling methods. We leverage CGRT to generate reasoning paths for these difficult samples, with the results provided in Table 1.



Figure 3: Part-of-Speech Analysis at Branching Nodes in CGRT, conducted on the GSM8K test set under 1-shot CoT setting via greedy decoding. For each query, we use two contexts to construct a binary CGRT. In the figure, good (or bad) context represents the model can reason correctly (or incorrectly) when we place it before the query.

405 From the experimental results, it can be observed that for challenging samples where majority vot-406 ing strategy fails, CGRT is capable of predicting some of these samples correctly. Furthermore, 407 for extremely difficult samples, *i.e.*, those for which correct inferences are not generated through 408 normal sampling methods even after numerous attempts, although CGRT cannot predict a majority of correct inference paths, it does include some correct paths among its predictions. The reason can 409 be roughly inferred to be that, CGRT employs a 'cross-generation' paradigm, allowing certain paths 410 to benefit from multiple contexts. Due to the limitations of the model's capability, extremely diffi-411 cult samples cannot be resolved with appropriate contexts or sampling methods, and in such cases, 412 CGRT is similarly unable to generate correct answers. 413

In the experiments, to avoid excessive computational costs due to an unmanageable number of branches generated by CGRT for certain samples, we limited the maximum number of paths for CGRT to 16, 32, and 48 when using 2, 4, and 8 contexts, respectively. We observed that as the number of contexts increased, the accuracy of the paths generated by CGRT using majority voting also generally improved (although the cost of constructing CGRT also increased). However, it is noted that the metrics for 8 contexts is not always better than those for 4 contexts. This is because while an increase in the number of contexts indeed raises the probability of uncovering correct inference paths, it also leads to an increase in the number of incorrect inference paths generated.

421 422 423

#### 3.2 PART-OF-SPEECH ANALYSIS AT BRANCHING NODES IN CGRT

424 We use two contexts for each query to construct a binary tree form of CGRT, using Llama3-8B 425 model on the GSM8K test set with a 1-shot CoT (Chain of Thought) setting. Since the goal is to 426 perform Part-of-Speech Analysis, we did not ignore the non-significant tokens defined in section 2.2 427 during the construction of the CGRT. The statistical results are shown in Figure 3. It can be seen 428 that in mathematical reasoning tasks, the proportion of numbers and operators in critical branching nodes is higher, intuitively. However, because the figure only shows proportions and the number of 429 non-critical branching nodes far exceeds that of critical ones, even if a branching node is a number 430 or an operator, it is difficult to determine whether it is a critical one. In fact, if a branching node is a 431 number or an operator, it is still more likely to be non-critical. In the appendix A.2.1, we illustrate

Models Detect	GSM8K		MAWPS		SVAMP		AQuA	
WIOUEIS \Dataset	SC	iCGRT	SC	iCGRT	SC	iCGRT	SC	iCGRT
Llama2-7B	18.59	19.98	53.97	52.94	40.00	39.00	12.20	11.81
Llama2-13B	49.09	50.95	67.65	69.32	70.00	70.00	14.17	14.57
Llama3-8B	67.26	69.29	78.99	79.83	68.00	71.00	18.11	18.50
Qwen2.5-7B-Instruct	42.03	44.21	84.45	85.29	75.00	76.00	16.53	15.74

 Table 2: Performance Comparison with SC Baseline

this point more clearly. In mathematical reasoning, parts of mathematical equations can adjust the order of computation, or the numbers or operators in the equation can be rearranged; for example, writing 16 - 3 - 4 as 16 - 4 - 3 is also acceptable. The part-of-speech analysis elucidates the complexity of branching nodes in the CGRT.

444 445 446

442

443

432

447 448

#### 3.3 PERFORMANCE COMPARISON WITH TRADITIONAL SC BASELINE

449 We evaluate iCGRT against SC Wang et al. on four mathematical reasoning datasets: GSM8K Cobbe et al. (2021), SVAMP Patel et al. (2021), MAWPS Koncel-Kedziorski et al. (2016), 450 and AQuA Ling et al. (2017), using four open-source LLMs: Llama2-7B Touvron et al. (2023), 451 Llama2-13B, Llama3-8B AI@Meta (2024), and Qwen2.5-7B-Instruct Team (2024). The test results 452 are shown in Table 2. All experiments are conducted using a 2-shot CoT setting. For each sample, 453 two different examples are randomly selected as demonstrations. In the baseline experiments for 454 SC, eight sets of 2-shot demonstrations are randomly selected for each query to serve as context. 455 For each context, top-3 sampling is employed to generate four inference paths, resulting in a total 456 of 32 inference paths. For iCGRT, the same eight sets of context are selected, and the same top-3 457 sampling method is applied to generate 32 tokens at each step (if the first 16 tokens is consistent or 458 if the count of the top 1 token exceeded 12, the remaining 16 tokens will not be computed). The 459 maximum number of paths for iCGRT is capped at  $n_p = 32$ .

460 461

#### 4 RELATED WORKS

462 463

464 The paradigm of Chain-of-Thought (CoT) Wei et al. (2022) is commonly employed today to en-465 hance the reasoning capabilities of models. Whether the manually or automatically constructed CoT demonstrations, or prompts Sahoo et al. (2024) or instructions Efrat & Levy (2020); Mishra et al. 466 (2022) used to inform the model to think step-by-step, these can all be considered as context. The 467 response given by LLMs to the same question can vary depending on the context. However, most 468 current research on how context influences LLM generation remains at a relatively macro level, such 469 as the length of the context Li et al. (2024), sequence order Chen et al.; Pezeshkpour & Hruschka 470 (2024), domain conflicts Wang et al. (2023b), etc. A barrier that limits researchers from examining 471 this issue from a more microscopic angle, such as at the token level, is that from the first inconsis-472 tent output token, subsequent generations will be jointly influenced by the context and the already 473 generated content. Moreover, as language sequences progress, deviations become increasingly pro-474 nounced. Due to the polysemy of language, the same semantics can be expressed in many different 475 ways. Therefore, finding a method to compare the reasoning content generated by LLMs under 476 different contexts poses an exceptionally challenging problem.

477 This paper addresses the issue by constructing CGRT, a specialized tree structure with a cross-478 generative approach. Other methods that combine tree structures with Chain-of-Thought (CoT), 479 such as ToT Yao et al. (2024), GoT Besta et al. (2024), XoT Ding et al. (2023), LLM+MCTS Zhang 480 et al. (2024a;b), etc., differ from CGRT in two main aspects: Firstly, these methods use "thoughts" 481 as units for tree nodes, whereas CGRT generates a tree structure at the token level. Secondly, the 482 core objective of these methods is to enhance the model's reasoning performance, requiring the introduction of additional verifiers or evaluators to validate and assess the thoughts. In contrast, CGRT 483 is a pure inference method that does not require an additional verifier. Instead, it leverages "cross-484 generation" to integrate positive influences from multiple contexts, thereby generating diversified 485 reasoning paths with a higher probability of correctness.

Due to the frequent occurrence of hallucinations and logical errors in LLMs Huang et al. (2023);
Xu et al. (2024), self-consistency Wang et al. is a widely adopted method to enhance the accuracy of LLMs' reasoning. The underlying insight is that LLMs may produce occasional errors in a single reasoning; thus, correct reasoning can be more possibly achieved through a majority voting on multiple diversely-sampled reasoning paths (by changing different contexts or utilizing varied decoding strategies Shi et al. (2024)).

492 493

494

### 5 DISCUSSIONS & FUTURE WORKS

495 **Inference Cost.** The inference speed of iCGRT is influenced by two factors: the number of tokens 496 predicted per step ( $n_c \cdot n_s$  in Algorithm 2), and the width of the tree (with the maximum width con-497 strained by  $n_p$  in Algorithm 2). Under all experiment settings presented in this paper, it is generally 498 set that  $n_c \cdot n_s = n_p$ , and the SC baseline for comparison is also configured to generate an equiva-499 lent number of inference paths. If the number of tokens generated at each step matches that of the baseline inference path, and considering the cumulative length of the paths within the tree structure 500 is typically longer than the average path length generated by the baseline, the computational cost 501 would be higher than that of the SC baseline. However, iCGRT can avoid excessive computational 502 burden by halving the number of tokens predicted at each step (as detailed in Section 2.3), which is 503 satisfied during most generation steps. Additionally, the actual width of iCGRT often does not reach 504 the predefined maximum width ( $n_n$  in Algorithm 2) in many scenarios. Therefore, the computational 505 overhead introduced by iCGRT is minimal. 506

Future Works: More In-depth Research on the Impact of Context on LLMs' Generation. As 507 a tool of analyzing the impact of context on LLMs' generation at the token level, CGRT utilizes 508 a "cross-generation" approach to eliminate the influence of previously generated content on the 509 current token, attributing the branching nodes in reasoning paths solely to variations in context. 510 However, the interpretability conclusions derived from this study remain preliminary, due to the 511 following two challenges: Firstly, while CGRT can generate all possible reasoning paths for LLMs 512 under different contexts for the same question, the polysemy of language means there are numerous 513 ways to express correct or incorrect reasoning. This multiplicity leads to a large number of tokens 514 non-decisive for the correctness of the reasoning. Future research should focus on identifying critical 515 branching nodes within CGRT. Secondly, the topic of how context affects LLMs' generation is still 516 acknowledged as a challenging subject, influenced by factors such as model size, pre-training, and post-training data, etc. Furthermore, whether LLMs merely represent powerful modeling of human 517 language or possess genuine cognitive abilities remains a widely debated topic without definitive 518 resolution Jin et al.; Valmeekam et al. (2024); Kambhampati (2024). Therefore, future works can 519 focus on two aspects: On one hand, we can gain more profound insights into the interpretability of 520 LLMs via CGRT. On the other hand, for the critical branching nodes of CGRT, efforts can be made 521 to develop a robust token-level decision-maker to improve LLMs' reasoning performance.

522 523

## 6 CONCLUSIONS

524 525

In this paper, we introduced a novel algorithm, Cross-Generation Reasoning Tree (CGRT), aimed 526 at enhancing the reasoning capabilities and interpretability of large language models (LLMs). The 527 key challenge addressed is the impact of context variation on LLM generation, particularly the dif-528 ficulties posed by diverging reasoning paths. CGRT tackles this by constructing a reasoning tree 529 that tracks token generation across different contexts, enabling a thorough analysis of the role that 530 context plays in model outputs. Our experimental results demonstrate that CGRT integrates the 531 strengths of different context variations and samplings, generating more accurate and diverse rea-532 soning paths. Moreover, CGRT reveals that even bad contexts can lead to correct reasoning paths 533 and provides insights into the critical branching points that determine the success of reasoning. We 534 also introduced an inference version, iCGRT, which uses majority voting at both token and path levels to improve decision-making at branching nodes, leading to more efficient and accurate reasoning in LLMs. In conclusion, CGRT offers a powerful tool for understanding and improving LLMs by 536 combining multiple contexts and sampling techniques, enhancing both reasoning performance and 537 interpretability. Future work could focus on deeper investigations into critical branching nodes and 538 further optimizing the balance between reasoning accuracy and computational costs.

## 540 REFERENCES

547

551

552

553

554

555

556

559

566

567

568

- 542 Claude 3.5 sonnet model card addendum. 2024. URL https://www-cdn.anthropic. 543 com/fed9cc193a14b84131812372d8d5857f8f304c52/Model\_Card\_Claude\_ 3\_Addendum.pdf.
- Samira Abnar and Willem Zuidema. Quantifying attention flow in transformers. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 4190–4197, 2020.
- Pranjal Aggarwal, Aman Madaan, Yiming Yang, et al. Let's sample step by step: Adaptiveconsistency for efficient reasoning and coding with llms. In *Proceedings of the 2023 Conference* on Empirical Methods in Natural Language Processing, pp. 12375–12396, 2023.
  - AI@Meta. Llama 3 model card. 2024. URL https://github.com/meta-llama/ llama3/blob/main/MODEL\_CARD.md.
  - Rohan Anil, Andrew M Dai, Orhan Firat, Melvin Johnson, Dmitry Lepikhin, Alexandre Passos, Siamak Shakeri, Emanuel Taropa, Paige Bailey, Zhifeng Chen, et al. Palm 2 technical report. *arXiv preprint arXiv:2305.10403*, 2023.
- Yonatan Belinkov. Probing classifiers: Promises, shortcomings, and advances. *Computational Linguistics*, 48(1):207–219, 2022.
- Maciej Besta, Nils Blach, Ales Kubicek, Robert Gerstenberger, Michal Podstawski, Lukas Gianinazzi, Joanna Gajda, Tomasz Lehmann, Hubert Niewiadomski, Piotr Nyczyk, et al. Graph of thoughts: Solving elaborate problems with large language models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pp. 17682–17690, 2024.
- Xinyun Chen, Ryan Andrew Chi, Xuezhi Wang, and Denny Zhou. Premise order matters in reason ing with large language models. In *Forty-first International Conference on Machine Learning*.
  - Yung-Sung Chuang, Yujia Xie, Hongyin Luo, Yoon Kim, James R Glass, and Pengcheng He. Dola: Decoding by contrasting layers improves factuality in large language models. In *The Twelfth International Conference on Learning Representations*, 2023.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.
- Ruomeng Ding, Chaoyun Zhang, Lu Wang, Yong Xu, Minghua Ma, Wei Zhang, Si Qin, Saravan Rajmohan, Qingwei Lin, and Dongmei Zhang. Everything of thoughts: Defying the law of penrose
  triangle for thought generation. *arXiv preprint arXiv:2311.04254*, 2023.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
- Avia Efrat and Omer Levy. The turking test: Can language models understand instructions? *arXiv preprint arXiv:2010.11982*, 2020.
- Angela Fan, Mike Lewis, and Yann Dauphin. Hierarchical neural story generation. arXiv preprint arXiv:1805.04833, 2018.
- Markus Freitag and Yaser Al-Onaizan. Beam search strategies for neural machine translation. In
   *Proceedings of the First Workshop on Neural Machine Translation*, pp. 56–60, 2017.
- Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. The curious case of neural text degeneration. In *International Conference on Learning Representations*.
- Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong
   Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, et al. A survey on hallucination in large language
   models: Principles, taxonomy, challenges, and open questions. *arXiv preprint arXiv:2311.05232*, 2023.

- Zhijing Jin, Jiarui Liu, LYU Zhiheng, Spencer Poff, Mrinmaya Sachan, Rada Mihalcea, Mona T Diab, and Bernhard Schölkopf. Can large language models infer causation from correlation? In *The Twelfth International Conference on Learning Representations*.
- Subbarao Kambhampati. Can large language models reason and plan? Annals of the New York
   Academy of Sciences, 1534(1):15–18, 2024.
- Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large language models are zero-shot reasoners. *Advances in neural information processing systems*, 35:22199–22213, 2022.
- Rik Koncel-Kedziorski, Subhro Roy, Aida Amini, Nate Kushman, and Hannaneh Hajishirzi.
  MAWPS: A math word problem repository. In Kevin Knight, Ani Nenkova, and Owen Rambow (eds.), *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp. 1152–1157, San Diego, California, June 2016. Association for Computational Linguistics. doi: 10.18653/v1/N16-1136. URL https://aclanthology.org/N16-1136.
- Tianle Li, Ge Zhang, Quy Duc Do, Xiang Yue, and Wenhu Chen. Long-context llms struggle with
   long in-context learning. *arXiv preprint arXiv:2404.02060*, 2024.
- Wang Ling, Dani Yogatama, Chris Dyer, and Phil Blunsom. Program induction by rationale generation: Learning to solve and explain algebraic word problems. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 158–167, 2017.
- Fei Liu et al. Learning to summarize from human feedback. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, 2020.
- Swaroop Mishra, Daniel Khashabi, Chitta Baral, Yejin Choi, and Hannaneh Hajishirzi. Reframing instructional prompts to gptk's language. In *Findings of the Association for Computational Linguistics: ACL 2022*, pp. 589–612, 2022.
- <sup>623</sup> R OpenAI et al. Gpt-4 technical report. ArXiv, 2303:08774, 2023.

619

624

- Arkil Patel, Satwik Bhattamishra, and Navin Goyal. Are NLP models really able to solve simple
   math word problems? In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp. 2080–2094, Online, June 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.
   naacl-main.168. URL https://aclanthology.org/2021.naacl-main.168.
- Pouya Pezeshkpour and Estevam Hruschka. Large language models sensitivity to the order of options in multiple-choice questions. In *Findings of the Association for Computational Linguistics:* NAACL 2024, pp. 2006–2017, 2024.
- Pranab Sahoo, Ayush Kumar Singh, Sriparna Saha, Vinija Jain, Samrat Mondal, and Aman Chadha.
  A systematic survey of prompt engineering in large language models: Techniques and applications. *arXiv preprint arXiv:2402.07927*, 2024.
- 637 Chufan Shi, Haoran Yang, Deng Cai, Zhisong Zhang, Yifan Wang, Yujiu Yang, and Wai Lam. A
  638 thorough examination of decoding methods in the era of llms. *arXiv preprint arXiv:2402.06925*, 2024.
- Qwen Team. Qwen2.5: A party of foundation models, September 2024. URL https://qwenlm.
   github.io/blog/qwen2.5/.
- Hugo Touvron, Louis Martin, Kevin R. Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei,
  Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Daniel M. Bikel, Lukas
  Blecher, Cristian Cantón Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony S. Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel M. Kloumann, A. V. Korenev, Punit Singh Koura, Marie-Anne Lachaux,

Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mi-haylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, R. Subramanian, Xia Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zhengxu Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. Llama 2: Open foundation and fine-tuned chat models. ArXiv, abs/2307.09288, 2023. URL https://api.semanticscholar.org/ CorpusID:259950998. 

- Karthik Valmeekam, Matthew Marquez, Alberto Olmo, Sarath Sreedharan, and Subbarao Kambhampati. Planbench: An extensible benchmark for evaluating large language models on planning and reasoning about change. *Advances in Neural Information Processing Systems*, 36, 2024.
- Lean Wang, Lei Li, Damai Dai, Deli Chen, Hao Zhou, Fandong Meng, Jie Zhou, and Xu Sun. Label
  words are anchors: An information flow perspective for understanding in-context learning. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pp. 9840–9855, 2023a.
- Kuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V Le, Ed H Chi, Sharan Narang, Aakanksha
   Chowdhery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language
   models. In *The Eleventh International Conference on Learning Representations*.
  - Yike Wang, Shangbin Feng, Heng Wang, Weijia Shi, Vidhisha Balachandran, Tianxing He, and Yulia Tsvetkov. Resolving knowledge conflicts in large language models. *arXiv preprint arXiv:2310.00935*, 2023b.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny
   Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837, 2022.
- Ziwei Xu, Sanjay Jain, and Mohan Kankanhalli. Hallucination is inevitable: An innate limitation of large language models. *arXiv preprint arXiv:2401.11817*, 2024.
- Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik
  Narasimhan. Tree of thoughts: Deliberate problem solving with large language models. *Advances in Neural Information Processing Systems*, 36, 2024.
- Dan Zhang, Sining Zhoubian, Yisong Yue, Yuxiao Dong, and Jie Tang. Rest-mcts\*: Llm self-training via process reward guided tree search. *arXiv preprint arXiv:2406.03816*, 2024a.
  - Di Zhang, Jiatong Li, Xiaoshui Huang, Dongzhan Zhou, Yuqiang Li, and Wanli Ouyang. Accessing gpt-4 level mathematical olympiad solutions via monte carlo tree self-refine with llama-3 8b. *arXiv preprint arXiv:2406.07394*, 2024b.

702 703		A APPENDIX
704 705		A.1 IMPLEMENT DETAILS
706		A.1.1 DATASETS
707 708 709 710 711		GSM8K Cobbe et al. (2021) is a dataset of diverse grade school math word problems created by human problem writers. The test set contains 1,319 problems. These problems need multi-step mathematical reasoning, usually taking between 2 and 8 steps to solve, and solutions primarily involve performing a sequence of elementary calculations using basic arithmetic operations $(+ - \times \div)$ to reach the final answer.
712 713 714		MAWPS Koncel-Kedziorski et al. (2016) dataset is a collection of simple math word problems focused on arithmetic. Its test set contains 238 samples.
715 716 717		SVAMP Patel et al. (2021) is a challenge set for elementary-level Math Word Problems (MWP). An MWP consists of a short Natural Language narrative that describes a state of the world and poses a question about some unknown quantities. The test set contains 100 problems.
718 719 720		AQuA Ling et al. (2017) test set consists 254 algebraic word problems with natural language ratio- nales. Each question provides 5 options and only 1 option is correct.
721 722		A.1.2 MODELS
723 724 725		We perform experiments using the model Llama2-13B Touvron et al. (2023), Llama3-8B AI@Meta (2024) and Qwen2.5-7B Team (2024). We use the base models for Llama2 and Llama3, instead of the -chat or -instruct version. For Qwen2.5-7B, we use the -instruct version.
727		A.1.3 FEW-SHOT COT SETTINGS
728 729		We use the following templates for the few-shot CoT inference:
<ul> <li>730</li> <li>731</li> <li>732</li> <li>733</li> <li>734</li> <li>735</li> <li>736</li> <li>737</li> <li>738</li> </ul>		Question: [QUESTION] Answer: [RATIONALE] The answer is [ANSWER]. Question: [QUESTION] Answer: [RATIONALE] The answer is [ANSWER].  Question: [QUERY] Answer:
739 740 741		where [QUESTION], [RATIONALE], and [ANSWER] represent the question, rationale, and the final answer of the demonstrations respectively. [QUERY] represent the question to be reasoned.
<ul> <li>742</li> <li>743</li> <li>744</li> <li>745</li> <li>746</li> <li>747</li> </ul>		A.1.4 DEFINITION NON-SIGNIFICANT TOKENS Due to that the flexibility of language permits the same concept to be articulated in numerous var- ied manners, we have to apply some limitations on the algorithm of building CGRT to avoid an extraordinarily high number of branches for some query question. we define the following tokens are non significant ones. If of a branching point some branches are non significant these will be
748		disregarded; if all branches are deemed non-significant, a single branch will be chosen at random.
<ul> <li>749</li> <li>750</li> <li>751</li> <li>752</li> <li>753</li> <li>754</li> <li>755</li> </ul>	1 3 5	<pre>non_significant_words = [     'and', 'or', 'but', 'nor', 'for', 'so', 'yet', 'then', '         therefore',     'in', 'on', 'at', 'to', 'by', 'with', 'about',     'a', 'an', 'the',     'of', 'as', 'if', 'that',     ',', '.', '!', '?', ';', ':', '_', '', "'", '"',</pre>
	7	<pre>'what', 'which', 'who', 'whom', 'whose',</pre>



Figure 4: Part-of-Speech Analysis at Branching Nodes

As Figure 4 shows, the quantity of non-critical branching nodes greatly exceeds that of their criti-cal counterparts. Hence, despite the fact that within critical branching nodes, tokens that represent numbers or operators do comprise a significant portion, such an assertion cannot be reciprocally ap-plied. Mathematical reasoning frequently encompasses non-serial logical computations that permit diverse permutations of operational hierarchies. Consequently, even in the CGRT, a branching node that represents a number or an operator cannot be easily identified as critical. 

# A.3 EXAMPLES OF CGRT

819

Here we present some examples from the CGRT of the GSM8K test set. The symbols ·, --, |, `
in the figures are used only for better visualization of the tree structure and have no actual meaning.
Because unrestricted CGRT often generates trees with a large number of branches, we have ignored
the branching nodes representing non-significant tokens.

And for better visualization, we have selected examples with an appropriate number and length of
 reasoning steps. These examples are presented in the form of vector graphics (Figure 5) rather than
 text.

820	Quartian: Relative and her research, Resear, share the cast of grounder. In total, they spend about 5400 per sorth. In a four-seek month, has many dollars dues Lessan spend per usek if Relative pays kill of the cost?
821	
020	
022	$-480$ / 4 $\times$ 200forms upper Mill of the test, so that of 200 $\times$ 40 (the means to 40, 100) $\times$ 200 $\times$ 40 (the means to 40, 100) $\times$
823	
824	Quantization Tyrizes theorem is the time receipt time for games one. If the paper and is the paper and and have the paper of the paper and and have the paper of the paper and and have the paper of the
825	4.43 is the state of the sta
826	· van 13 for name very 2 days, this same is 12.
827	
027	
828	Quantum Segme register and method for more and white of allow can its worth bus cents and a plattic battle is worth three sents. One priori three alumines can of sinds and frag plattic battles of water a work, has any cents data Segme are free registing is a fore-set work.
829	1 7 + 7 + 7 + 12 cm 1, Modeging and mark 2 does for a constant multiply that gamma 21 + 4 + 84 cm 14 + 14 + 4 + 84 cm 14, Modeging and 14 + 14 + 14 + 14 + 14 + 14 + 14 + 14
830	- shallow can are warth 3 * 2 + 6 cents. but plantic bettlan are worth 3 * 3 = 33 cents. but html, Grayne worth 5 * 3 = 72 cents a weat, bitter tent, is worth for fare weat, is a dwaren 23 * 4 * 84 cents in
831	- and that much for fur works, in the next $21 + 6 = 60$ cents is $a - bar - abc - abc = 0$ . The instant is $1.4$
832	- NULL, UNL of the set to their furt to pick same place and fund bulf as any huge as arts in the prefers. If there are 98 arts, nailcaines the statel number of inserts in the prefers.
002	80 - m tr = 23 http://dl m tr + 25 http://dl m tr + 26 http://dl m
033	- / 2 = 23(23 + 58 - 73)/The answer is 75.
834	Question have years to be yet jaces to be starting (down 4 which jaces at head) for the heading down, and have been for down 1 by the desce plates and \$2.60 which and head which are been which it are places which are been which are been which are been and been which are been and been are been and been are been and been are been and been are be
835	ament havry will upwed as place writings: 172, 20 + 183, 40 + 143, 40 + 1100, 400 + 1000, 4000, 400 + 1000, 4000,
836	- out of the dimer plane, SLM * D + 172.00/dest first be total and of the back, SLM * D + 182.00/defaily, field be total and of the solal plane, 14.00 * D + 10.40/defaily dimer back field be total and of the plane settings, SZM * SM + 10.40/defaily, field be total and of the solar plane, 14.00 * D + 10.40/defaily, field be total and of the solar plane, 14.00 * D + 10.40/defaily, field be total and of the solar plane, 14.00 * D + 10.40/defaily, field be total and of the solar plane, 14.00 * D + 10.40/defaily, field be total and of the solar plane, 14.00 * D + 10.40/defaily, field be total and of the solar plane, 14.00 * D + 10.40/defaily, field be total and of the solar plane, 14.00 * D + 10.40/defaily, field be total and of the solar plane, 14.00 * D + 10.40/defaily, field be total and of the solar plane, 14.00 * D + 10.40/defaily, field be total and of the solar plane, 14.00 * D + 10.40/defaily, field be total and of the solar plane, 14.00 * D + 10.40/defaily, field be total and of the solar plane, 14.00 * D + 10.40/defaily, field be total and of the solar plane, 14.00 * D + 10.40/defaily, field be total and of the solar plane, 14.00 * D + 10.40/defaily, field be total and of the solar plane, 14.00 * D + 10.40/defaily, field be total and of the solar plane, 14.00 * D + 10.40/defaily, field be total and of the solar plane, 14.00 * D + 10.40/defaily, field be total and field be total and of the solar plane, 14.00 * D + 10.40/defaily, field be total and
837	10.0 + 110.0 to the same is 110.0.
000	- B paper - 3 pice setting-types - 3 pice setting-types - 10.4 m -
030	- + 3 place strategick place strategy - 5.8.4 = DLL.Regick place s
839	<ul> <li>Standard place settings 15.8.</li> <li>Standard place settings 15.8.</li> <li>Standard place settings 15.8.</li> <li>Standard Standard place settings 15.8.</li> </ul>
840	Ampliant Even works for 1 hors on handor. So bindendar be works and handor to be bank on handor. So Thereider be works 2 hors hand then to be bank on blandor handor and hand there been and h
841	to here in Diades, in News in Menning, well there are Menning (the test induce of heres Brace works is 5 + 20 + 1 + 23 heres, (the amount is 23.
842	
0/2	- 2 4 5 5 5 18 2 6 2 7 5 19 2 7 5 19 2 7 5 19 2 7 5 19 2 7 5 10 2
043	<sup>1</sup> -2 + 5 + 0 + 18 + 18 + 2 h
844	Quantizer, 2m/s lag of fulliment conty has 25 shouldes have and 86 candied applies. Each character have only in such candied applie. 2F and desculate have weights Big, have much dama 2m/s lag of conty weight, in grant
845	2 - Charling for weight 2016 2006/2018 -
846	
847	- gravital analysis is 100.
8/8	-glothe ansars is 200g.
040	Quantization formy clays 9 from - list regist. His Fried Tames Light rolp 2/3 of the History Algor. Line may more have, and for any Line base 2. 9 from 16 Line 9 from 2.21 of Brown Line Line Line Line Line Line Line Lin
849	- survey align b hows not have a starts table to hows. Knowy starts 6.6 × 1 hows, more that hows, when starts 1 hows, the starts the hows, when starts have how more than hows, the starts the hows, the h
850	-curry shart blance and anne shart blance. Hanny shart blance anne blance anne blance anne blance anne shart blance.
851	-morey shape 4 have and have shape 4 have. Herey shape 1-6-6-5 have more than Jame, both a many 16 have.
852	Quartiant: Early hight 28 splin for 65 cents and and rectived a 51 discussi. Colly hight 28 splin for 50 cents and and erectived a 54 discussi. The share with the state of the state and the rectived at the state of the state and the rectived at the state of the state at the state at the state of the state at the state at the state of the state at the state of the state at the sta
853	27 - 4 C - cont - B. World - M cont - Stall-World - B. D. S. B. World - D. Wo
054	-000 - 1 + 00 cert/x00 + 30 + 300 cert/x00092 of 1000 + 300 - 000 cert/x-000 + 100 - 000 + cert/x - 000 mer that he/y/x00e means to 1 cert.
004	-23 = 100 control to 100 = 100 + 10 control to a control to control to a control to a control to a control to a control to
855	apples for 60 cents cents = No28 apples for 10 cents cents = Elf-Weidling part f =-1 cent then beingly, MHE ansare is 81. 
856	Queriants for Elign's first report while 24 lines as a new provide as it from 24 of they had B,200 option unbindy, how any option at the total at 12
857	- organize of default by 12 to E_AMM empires (orgAMM empires in the member of experts that hereid shids/offer answer is E_AMM.
858	··· / is * e.vev upons. un/ANY COLLES IN HETERS IN HETY COPIES IN HETERS IN HETY COPIES IN HETERS IN COPIES, UTTER HEMBER IN HERE COPIES, UTTER HERE COPIES, UTTER HEMBER IN HERE COPIES, UTTER HERE COPIES, UTTER HERE IN HERE COPIES, UTTER
850	
000	
860	
861	Figure 5: Examples of CGRT.
862	
863	