

RASD: Retrieval-Augmented Speculative Decoding

Anonymous ACL submission

Abstract

Speculative decoding accelerates inference in large language models (LLMs) by generating draft tokens for target model verification. Current approaches for obtaining draft tokens rely on lightweight draft models or additional model structures to generate draft tokens and retrieve context from databases. Due to the draft model’s small size and limited training data, model-based speculative decoding frequently becomes less effective in out-of-domain scenarios. Additionally, the time cost of the drafting phase results in a low upper limit on acceptance length during the verification step, limiting overall efficiency. This paper proposes RASD (Retrieval-Augmented Speculative Decoding), which adopts retrieval methods to enhance model-based speculative decoding. We introduce tree pruning and tree fusion to achieve this. Specifically, we develop a pruning method based on the draft model’s probability distribution to construct the optimal retrieval tree. Second, we employ the longest prefix matching algorithm to merge the tree generated by the draft model with the retrieval tree, resulting in a unified tree for verification. Experimental results demonstrate that RASD achieves state-of-the-art inference acceleration across tasks such as DocQA, Summary, Code, and In-Domain QA. Moreover, RASD exhibits strong scalability, seamlessly integrating with various speculative decoding approaches, including both generation-based and retrieval-based methods.

1 Introduction

Transformer-based Large Language Models (LLMs) (Vaswani et al., 2017; Brown et al., 2020) exhibit remarkable capabilities and are extensively applied across diverse domains. However, autoregressive generation in LLMs produces tokens sequentially, resulting in slow inference speeds. To address this issue, an innovative approach called Speculative Decoding has been introduced (Chen

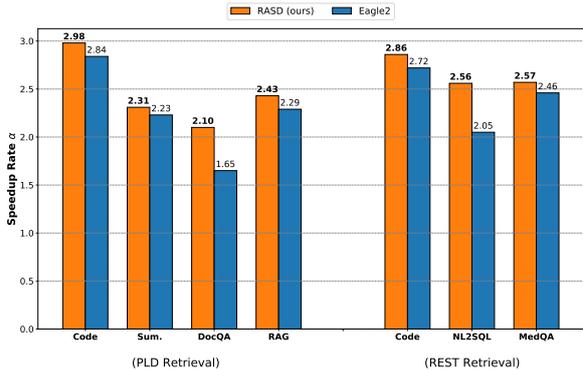


Figure 1: Speedup Performance of EAGLE2 vs. RASD with PLD and REST Retrieval Methods on Qwen2.5-14B.

et al., 2023; Miao et al., 2023). Speculative decoding modifies the inference process by dividing the LLM’s task into two phases: a cost-efficient draft phase and a parallel verification phase. This strategy significantly improves computational parallelism in LLM inference. By enabling LLMs to generate multiple tokens simultaneously and minimizing the time spent in the drafting and verification phases, speculative decoding reduces the overall inference time (Leviathan et al., 2023).

In speculative decoding, if the majority of draft tokens are correct, the overall decoding steps can be significantly reduced. Therefore, obtaining draft tokens with a high acceptance rate is essential. Some researchers employ small draft models to predict draft tokens (Leviathan et al., 2023; Chen et al., 2023; Li et al., 2024b,a). These draft models are trained to match the distribution of the target model and then generate draft tokens. Other researchers have proposed using a parameter-efficient model structure to generate the next k candidate tokens in a single forward pass of the target model (Cai et al., 2024; Xiao et al., 2024). This approach also requires training the additional structure. An alternative approach is retrieval-based speculative

068 decoding, which is a train-free method (Saxena, 069 2023; He et al., 2024). Here, a retrieval library 070 pre-defines tokens to follow the suffix of the cur- 071 rent content as draft tokens. This approach is fast 072 and eliminates the need for additional model train- 073 ing. Although it is generally less effective than 074 generation-based methods, it performs well in spe- 075 cific knowledge-intensive scenarios.

076 In draft model-based speculative decoding, the 077 acceptance rate of draft tokens hinges on the draft 078 model’s capabilities. However, these models typi- 079 cally have simple structures with limited param- 080 eters, which restricts their ability to effectively retain 081 training data. Typically, a draft model is trained 082 on a general dataset to perform well in common 083 scenarios, such as the ShareGPT dataset. As a re- 084 sult, in specific scenarios, the draft model’s token 085 acceptance rate tends to be low. This issue arises 086 because the draft model **lacks capabilities for out-** 087 **of-domain scenarios**. There has been work on 088 improving the generalizability of the draft model 089 through feature-level supervision, which closes the 090 representation gap between the draft model and the 091 target model without requiring large amounts of 092 data (Li et al., 2024b; Du et al., 2024). However, 093 we have evaluated these methods on our domain 094 datasets and discover that the draft model still does 095 not perform well in the specific downstream do- 096 main after our test of business data. Moreover, 097 generating the draft token sequence requires multi- 098 ple forward passes through the draft model, making 099 the process time-intensive. Consequently, previous 100 methods limit the length of candidate draft tokens. 101 This restriction leads to a **low upper limit on the** 102 **acceptance length** during a single verification step.

103 In this paper, we propose the RASD (Retrieval- 104 Augmented Speculative Decoding) method to ad- 105 dress the draft model’s shortcomings mentioned 106 above. First, we design an algorithm to select suit- 107 able retrieval results and then combine the draft 108 tokens using a tree-fusion approach. This enables 109 the generation of draft tokens that incorporate in- 110 formation from the language model and retrieval. 111 Finally, the target model verifies the draft tokens to 112 achieve acceleration.

113 We tested two retrieval methods PLD (Saxena, 114 2023) and REST (He et al., 2024) to enhance 115 model-based speculative decoding. As shown in 116 Figure 1, the PLD method performs well in tasks 117 where the output includes input content. There- 118 fore, we conducted comprehensive experiments 119 using RASD (PLD) on datasets with these charac-

120 teristics across four tasks: code generation, doc- 121 ument question answering, summarization, and 122 retrieval-augmented generation. We have con- 123 ducted experiments on HumanEval (Chen et al., 124 2021), CNN/Daily Mail (Nallapati et al., 2016), 125 MFQA (Bai et al., 2024), and DPR (Karpukhin 126 et al., 2020). RASD (PLD) achieved the best re- 127 sults in all four tasks.

128 To extend the applicability of the RASD method, 129 we also tested the REST retrieval method. REST 130 is theoretically beneficial for any dataset when the 131 database’s data distribution closely matches the 132 test set, excelling in knowledge-intensive tasks. 133 We experimented with RASD (REST) on three 134 tasks: HumanEval, MedQA ¹, and NL2SQL. The 135 results demonstrate that RASD (REST) effectively 136 improves the acceleration ratio of speculative de- 137 coding.

138 2 Related Work

139 2.1 EAGLE

140 Model-based speculative decoding is widely re- 141 garded as the most effective approach for achieving 142 acceleration. As the current state-of-the-art method, 143 the EAGLE series method (Li et al., 2024b,a; Gao 144 et al., 2024) is designed to provide feature super- 145 vision signals. The core structure of the EAGLE 146 draft model includes one layer of the target model, 147 sharing parameters with the embedding layer and 148 the language model head. During training, only the 149 parameters of this layer are updated, while the em- 150 bedding layer and the language model head remain 151 frozen.

152 Before training, EAGLE processes the training 153 data through the target model to extract the output 154 features from its last layer. During draft model 155 training, EAGLE combines the current embedding 156 input with the previous output features of the target 157 model to form a new input. It introduces a loss 158 function based on the output features of the last 159 layer, aiming to closely align the draft model’s 160 output features with those of the target model.

161 EAGLE-2 (Li et al., 2024a) follows the same 162 design but replaces the static tree structure with dy- 163 namic drafting structures during decoding to gener- 164 ate higher-quality candidate trees.

165 In RASD, we employ EAGLE-2 as the model- 166 based speculative decoding method. Any specula- 167 tive decoding approach that utilizes tree attention in

¹<https://huggingface.co/datasets/lavita/medical-qa-datasets>

168 drafting and verifying can be enhanced with RASD,
 169 as we do not change the method itself.

170 Before training, EAGLE processes the training
 171 data through the target model to extract the output
 172 features from its last layer. During draft model
 173 training, EAGLE combines the current embedding
 174 input with the previous output features of the target
 175 model to form a new input. It introduces a loss func-
 176 tion based on the output features of the last layer,
 177 aiming to closely align the draft model’s output
 178 features with those of the target model. EAGLE-2
 179 (Li et al., 2024a) follows the same design but re-
 180 places the static tree structure with dynamic draft-
 181 ing structures during decoding to generate higher-
 182 quality candidate trees. In RASD, we employ
 183 EAGLE-2 as the model-based speculative decoding
 184 method. Any speculative decoding approach that
 185 utilizes tree attention in drafting and verifying can
 186 be enhanced with RASD, as we do not change the
 187 method itself.

188 2.2 PLD

189 Retrieval-based speculative sampling methods per-
 190 form well in knowledge-intensive scenarios. PLD
 191 (Prompt Lookup Decoding) is a simple and low-
 192 cost method for retrieving sequences from the input.
 193 However, it cannot predict new tokens or their com-
 194 binations. It relies on the last n tokens of the input
 195 for n -gram matching. The original implementa-
 196 tion² terminates at the first match. In our version,
 197 we return the first k matching results instead.

198 2.3 REST

199 REST is a novel algorithm that leverages retrieval
 200 to generate draft tokens. Unlike PLD, which re-
 201 trieves directly from the input, REST retrieves from
 202 a pre-defined context database. It utilizes existing
 203 knowledge, fetching relevant tokens based on the
 204 current context without relying on draft models.
 205 REST converts existing corpora into a retrieval li-
 206 brary and organizes the results into a draft tree for
 207 verification by the target model. As a plug-and-
 208 play speculative sampling method, REST does not
 209 require additional draft model structures. While
 210 user-friendly, it is less effective compared to draft
 211 model-based methods.

212 We employ REST as the retrieval method to
 213 demonstrate the excellent scalability of RASD.
 214 All speculative decoding methods with external
 215 retrieval can be integrated into RASD.

²[https://github.com/apoorvumang/
 prompt-lookup-decoding/](https://github.com/apoorvumang/prompt-lookup-decoding/)

216 3 Retrieval-Augmented Speculative 217 Decoding

218 In this section, we first present the background of
 219 speculative decoding and then introduce our pro-
 220 posed RASD framework.

221 3.1 Background: Speculative Decoding

222 Speculative decoding with the draft model com-
 223 bines drafting and verification processes. In a sin-
 224 gle step of drafting, we use s to denote the input,
 225 which contains n tokens, including those from user
 226 input and tokens generated by previous speculative
 227 decoding processes:

$$228 \quad s = (x_1, \dots, x_{n-1}, x_n). \quad (1)$$

229 The target model generates the first token y_0 based
 230 on the target model probability distribution com-
 231 puted by s :

$$232 \quad y_0 \sim P(x \mid s; \theta_{\text{target}}), \quad (2)$$

233 where θ_{target} represents the parameters of the target
 234 model. If the current step is not the first step of
 235 speculative decoding, the token y_0 will be gen-
 236 erated together with other tokens in the verification
 237 phase.

238 The draft model concatenates y_0 to s to get the
 239 input s' .

$$240 \quad s' = (x_1, \dots, x_{n-1}, x_n, y_0). \quad (3)$$

241 Given the input, the draft model generates multi-
 242 ple candidate tokens by autoregressive decoding.
 243 Assuming the m -th token is the last token of the
 244 current draft phase:

$$245 \quad \hat{p}_m = P(x \mid s', \hat{y}_1, \dots, \hat{y}_{m-1}; \theta_{\text{draft}}), \quad (4)$$

$$246 \quad \hat{y}_m \sim \hat{p}_m, \quad (5)$$

$$247 \quad \hat{s}' = (s', \hat{y}_1, \dots, \hat{y}_m), \quad (6)$$

248 where \hat{p}_m represents the probability distribution of
 249 the m -th draft token, \hat{y}_m represents the m -th draft
 250 token. \hat{s}' represents the output of the draft model
 251 in this step, and θ_{draft} represents the parameters of
 252 the draft model.

253 In the verification process, the target model veri-
 254 fies \hat{s}' in a single forward pass. The target model
 255 probability of the k -th token is p_k . The draft token
 256 \hat{y}_k has an acceptance probability $\min\left(1, \frac{p_k}{\hat{p}_k}\right)$.
 257 If the draft token \hat{y}_k is rejected, all subsequent tokens
 258 are discarded, and this token is resampled from a
 259
 260

distribution $\|\max(0, p_k - \hat{p}_k)\|$. If \hat{y}_k is accepted, \hat{y}_k is transferred to y_k and the target model continues to verify the next token \hat{y}_{k+1} until the last token. Finally, the accepted sequence and the re-sampled token compose s' for the next speculative decoding step. The probability distribution generated by the verification method is equivalent to sampling directly from the target model (Leviathan et al., 2023).

We employ tree attention for the simultaneous verification of multiple candidates. Traditional causal attention masks are designed for linear sequences, where each token attends to all previous tokens, limiting speculative decoding to verify only one sequence at a time. Tree attention alters the attention mask to compress multiple sequences into a single merged sequence while preserving a tree structure. Each child node attends only to its parent nodes. In summary, speculative decoding, through guess-and-verify and tree attention, robustly and efficiently improves inference latency compared to autoregressive decoding. RASD relies on tree attention to achieve retrieval-augmented speculative decoding.

3.2 Our Approach: RASD

Our method RASD consists of three main steps: retrieval process, retrieval tree construction, and tree fusion. An overview of these steps is shown in Figure 2.

3.2.1 Retrieval Process

In RASD, retrieval is employed to enhance the quality of candidate draft tokens by leveraging knowledge sources. We evaluate two retrieval settings: PLD and REST. To optimize retrieval efficiency and accuracy, an exact-match approach, inspired by PLD, is utilized to identify continuation candidates.

Given the context s' , we retrieve context-continuation pairs from the datastore D , generating a set of continuation candidates S :

$$S = \text{Retrieve}(D, s'), \quad (7)$$

where $\text{Retrieve}(D, s')$ is a retrieve method based on suffix matching, which returns a set of retrieval results using s' as the query to find contexts in D that match the longest suffix of s' . We illustrate a detailed process in Algorithm 1. After this process, we obtain the candidate sequences.

For PLD, the retrieval process is designed to be iterative. If no results are initially found, we

Algorithm 1 Retrieve Method of Suffix Match

Require: x_n is the n -th token of the draft model input, and y_0 is the first token generated by the target model. Datastore D , result set S , suffix length limits n_{max} and n_{min} , number of candidates n , and length of candidates l .

```

1: for  $i = n_{max}$  downto  $n_{min}$  do
2:   if  $\text{len}(S) \geq n$  then
3:     break
4:   end if
5:    $\text{suffix} \leftarrow (x_{n-i+1}, \dots, x_n, y_0)$ 
6:   for each  $i\_gram$  in  $D$  do
7:     if  $i\_gram == \text{suffix}$  then
8:        $c \leftarrow$  next  $l$  tokens after  $i\_gram$  in  $D$ 
9:       if  $c$  not in  $S$  then
10:        Add  $c$  to  $S$ 
11:       end if
12:       if  $\text{len}(S) \geq n$  then
13:         break
14:       end if
15:     end if
16:   end for
17: end for

```

recursively select the top- k tokens generated by the draft model and append them to the context s' . This process continues until retrieval results are successfully obtained. This iterative approach improves the success rate of retrieval, even when the initial query fails to yield matches.

For REST, the retrieval process is limited to a single execution within each draft phase. This restriction is imposed due to the method's higher time overhead. This approach balances the richness of retrieved candidates with the practical constraints of runtime performance.

3.2.2 Retrieval Tree Pruning

We then construct the retrieval tree using the set of candidates S . For PLD, each candidate t_i can be treated as a linked list. We merge nodes in S that share the same prefix, transforming these linked lists into one or more trees. The last token y_0 of input s' is used as the root node to merge these trees, resulting in the retrieval tree T_r . For REST, the number of retrieval results can be substantial. To manage this, we construct the retrieval tree T_r by prioritizing candidates t_i using high-frequency prefixes as filters, following the same approach as REST.

Next, we propose a method to prune the retrieval

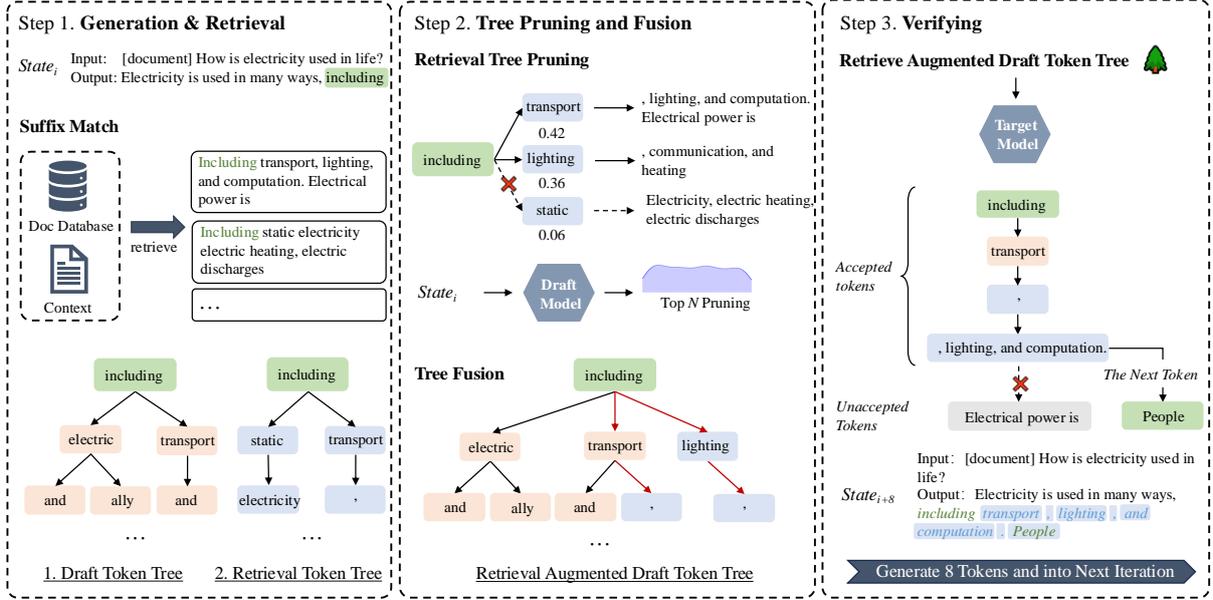


Figure 2: An overview of RASD. We obtain the draft token tree and retrieval results through the draft model generation and retrieval methods, respectively. In the next step, we construct and prune the retrieval tree. Then, we fuse the two trees, resulting in the retrieval-augmented draft token tree. Finally, the retrieval-augmented draft token tree is verified recursively. In the figure, green tokens denote y_0 in the current turn, and red tokens are accepted by the target model.

tree using the draft model’s output distribution. There is a strong positive correlation between the draft model’s confidence score and the token acceptance rate. We leverage this confidence score to prune the retrieval tree.

In the first layer of target model output, the probability distribution is P_1 :

$$P_1 = P(x | s'; \theta_{\text{draft}}). \quad (8)$$

We hypothesize that the first token of a high-quality retrieval result should have a high probability of appearing in P_1 . We reject the search results whose first token is not in the top- k of P_1 .

3.2.3 Tree Fusion

Thus far, we have obtained both the retrieval tree and the tree generated by the draft model. To produce the final output tree, we combine the draft model’s generation tree with the retrieval tree.

Given that the forward propagation time of a large language model theoretically scales quadratically with input length, tree fusion can significantly reduce this time by eliminating redundant input tokens. Specifically, we perform a simple longest prefix match on each branch of both trees. Branches with identical prefixes in the two trees are merged,

using the last node of the shared prefix as the parent node. This classic algorithm can be efficiently implemented using a trie tree.

After fusing the two trees, the attention matrix and position embeddings for the nodes of the new tree must be updated accordingly. The final combined tree integrates information from both the draft language model and knowledge base retrieval. Knowledge-based retrieval effectively addresses the out-of-domain problem and increases the upper limit of the length of candidate draft sequences.

3.2.4 Draft Tree Verification

We adopt a recursive verification strategy to accommodate tree attention. Utilizing tree attention, the target LLM calculates the probability of each token in the tree-structured draft within a single forward pass. Each node of the tree is arranged in one dimension, following the order of level-order traversal. The position embedding and attention matrices respectively represent the current node’s level in the tree and its relationship with the parent node.

At each node of the draft tree, we recursively apply speculative decoding algorithms to sample or adjust the probability distribution, in line with

model	method	HumanEval		CNN/DM		MultiFieldQA		Qasper		DPR	
		SR	τ	SR	τ	SR	τ	SR	τ	SR	τ
Temperature = 0											
Q 14B	PLD	1.49	0.55	1.59	0.70	1.98	1.60	1.62	1.58	1.57	0.70
	EAGLE2	2.84	3.10	2.23	2.41	1.65	1.73	1.48	1.83	2.39	2.60
	RASD(PLD)	2.98	3.44	2.31	2.68	2.10	2.78	1.66	2.58	2.43	2.92
L 8B	PLD	1.38	0.51	1.34	0.53	1.90	1.39	1.46	1.04	1.55	0.75
	EAGLE2	2.76	3.37	2.15	2.51	1.68	2.28	1.52	2.36	2.37	2.97
	RASD(PLD)	2.79	3.55	2.21	2.72	1.98	2.86	1.59	2.84	2.44	3.38
Temperature = 1											
Q 14B	PLD	1.53	0.58	1.42	0.49	1.72	1.25	1.54	1.53	1.74	0.59
	EAGLE2	2.73	3.09	1.95	2.03	1.56	1.61	1.39	1.76	2.40	2.36
	RASD(PLD)	3.00	3.39	1.98	2.17	1.93	2.29	1.60	2.48	2.49	2.63
L 8B	PLD	1.32	0.45	1.37	0.43	1.81	1.45	1.53	0.95	1.42	0.65
	EAGLE2	2.57	3.38	2.04	2.26	1.60	2.35	1.53	2.61	2.19	2.76
	RASD(PLD)	2.59	3.48	2.13	2.42	1.82	2.88	1.66	3.16	2.24	3.11

Table 1: The speedup performance of our proposed RASD with PLD retrieval method and baselines in different datasets. We test on Qwen2.5-14B and LLaMA3-instruct-8B with temperatures of generation 0 and 1. The best results among all methods are in bolded.

the approach of SpecInfer (Miao et al., 2023).

4 Experiments

4.1 Models and Tasks

To evaluate the effectiveness of RASD in accelerating large language models, we use EAGLE2 as our draft model and conduct a series of experiments with two different target models across various tasks. We tested RASD on the LLaMA3-Instruct-8B (Dubey et al., 2024) and Qwen2.5-Instruct-14B (Yang et al., 2024) models to assess its acceleration capabilities. The RASD (PLD) retrieval enhancement method is anticipated to perform well on tasks where the input contains potential output n-grams, such as in RAG (Karpukhin et al., 2020), DocQA, Summary, and Code tasks. Therefore, we selected the HumanEval, CNN/Daily Mail, MultiFieldQA, Qasper, and DPR datasets for our experiments. To push the boundaries of the RASD method, we utilize the REST-based retrieval method RASD (REST) to enhance tasks in more scenarios, choosing the HumanEval, NL2SQL, and Med-QA datasets for RASD (REST).

Both greedy sampling and non-greedy sampling are considered in all experiments to comprehensively evaluate speculative decoding performance. All evaluations are conducted on an NVIDIA A100 80G GPU.

4.2 Metrics

In our experiments, we adopt two primary metrics: average acceptance length τ and speedup ratio (SR). The average acceptance length τ evaluates the average number of tokens accepted per forward pass by the target large language models, excluding any overhead associated with retrieving or constructing draft tokens. This metric indicates the maximum possible acceleration. The second metric, speedup ratio (SR), measures the relative improvement in decoding speed compared to vanilla auto-regressive decoding.

4.3 Baseline

In this study, we focus solely on lossless speculative decoding approaches for LLMs. Among the methods that do not rely on draft models, we examine Prompt Lookup Decoding (PLD) (Saxena, 2023), REST (He et al., 2024), and EAGLE-2 (Li et al., 2024a), the latter being regarded as the state-of-the-art method for lossless speculative decoding tasks. Collectively, these baseline methods provide a robust framework for evaluating the efficiency of RASD in the LLM decoding process.

4.4 Training

We use the SharedGPT dataset, which comprises 68,000 dialogues from the Vicuna (Chiang et al., 2023) series models’ supervised fine-tuning dataset, as our training corpus. Given the substantial

model	method	HumanEval		NL2SQL		MedQA	
		SR	τ	SR	τ	SR	τ
Temperature = 0							
Q 14B	REST	1.69	0.86	2.36	2.39	1.66	0.84
	EAGLE2	2.72	3.10	2.05	2.36	2.46	2.57
	RASD(REST)	2.86	3.46	2.56	2.45	2.57	2.94
L 8B	REST	1.70	0.84	3.59	1.51	1.85	0.89
	EAGLE2	2.76	3.36	3.78	2.76	3.13	2.76
	RASD(REST)	2.86	3.75	4.07	3.41	3.26	3.00
Temperature = 1							
Q 14B	REST	1.76	0.83	2.32	2.19	1.42	0.50
	EAGLE2	2.73	3.09	2.07	2.35	2.24	2.06
	RASD(REST)	2.90	3.48	2.40	3.30	2.35	2.36
L 8B	REST	1.53	0.85	3.58	1.14	1.30	0.51
	EAGLE2	2.57	3.38	4.05	2.77	2.24	2.47
	RASD(REST)	2.59	3.44	4.24	3.26	2.35	2.75

Table 2: The speedup performance of our proposed RASD with REST retrieval method and baselines in different datasets. We test on Qwen2.5-14B and LLaMA3-instruct-8B with temperatures of generation 0 and 1. The best results among all methods are in bolded.

time and computational resources required, we choose not to regenerate responses for each dialogue turn using the target LLMs. Conducting training without re-generated data across all comparative methods remains equitable, although previous work (Li et al., 2024b) suggests that such an approach could slightly enhance the performance of the draft model. The learning rate is set to $5e-5$, with $(\beta_1 = 0.9, \beta_2 = 0.95)$ for the AdamW (Loshchilov and Hutter, 2019) optimizer, and we implement gradient clipping at 0.5. We utilized eight NVIDIA A100 80G GPUs for the training process.

4.5 Experimental Results

Table 1 shows the performance of RASD (PLD) compared to other methods. RASD (PLD) consistently achieves the highest speedup across all tasks and models. Unlike the PLD method, RASD can generate draft tokens that reflect the language model distribution. Compared to EAGLE2, RASD retrieves more accurate draft tokens. When the PLD score is high, the improvement with RASD (PLD) is more significant. RASD (PLD) shows the most improvement in DocQA tasks (MultiFieldQA and Qasper) because these tasks usually contain repeated paragraphs in the input. In tasks with fewer repeated input segments, such as code and summary tasks, the improvement with RASD (PLD) is less pronounced.

Table 2 shows the performance of RASD (REST) against other methods. RASD (REST) outperforms in terms of speedup across all tasks and models. Compared with PLD, REST has a larger retrieval space and applies to more scenarios. Therefore, the improvement of RASD (REST) verifies RASD’s compatibility with more retrieval methods. For tasks in the field of medical question answering, RASD can also achieve better results. We consider that for most knowledge-dependent tasks, RASD can effectively improve the speed of speculative decoding by building a contextual database in the domain.

5 Ablation Study

Method	SR	τ
EAGLE2	2.73	3.09
RASD(REST) w/o p	2.82	3.62
RASD(REST) w/o tf	2.87	3.48
RASD(REST)	2.90	3.48

Table 3: The impact of pruning and tree fusion operations in RASD on the speedup performance in HumanEval.

5.1 Pruning and Tree Fusion

We investigated how retrieval tree pruning and tree fusion impact the results of RASD experiments.

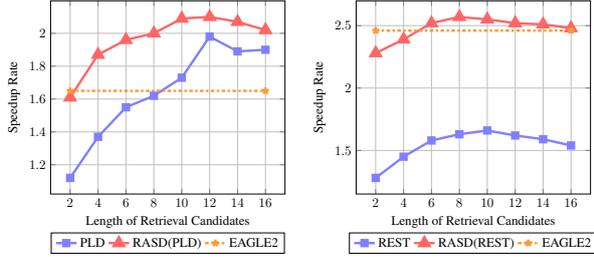


Figure 3: RASD performance with different lengths of retrieval candidates compared with the baselines

We conducted experiments on the qwen2.5-14b model using the HumanEval dataset with the REST method. As shown in Table 3, both methods enhance the speedup ratio, with RASD incorporating pruning yielding a more pronounced effect. Pruning reduces the average acceptance length by eliminating unnecessary retrieval results early, thereby saving time required for the target model’s forward pass. Tree fusion has a relatively smaller impact. It does not alter the average acceptance length but reduces redundant draft tokens, consequently saving time for the target model’s forward pass.

5.2 Length of Retrieval Candidates

In the retrieval phase, the length of retrieval candidates l is a crucial variable. If l is small, the advantage of retrieval may not be significant, and if l is large, the verification phase will require more time. We construct experiments on qwen2.5-14b, MultiFieldQA dataset used RASD(PLD) and qwen2.5-14b, MedQA dataset used RASD(REST) with different length of retrieval candidates l . As shown in 3, On the left, since PLD is suitable for DocQA tasks, it can be seen that when only the PLD method is used, as l increases, it can be inferred that the acceptance length of PLD is increasing. When $l > 0$, the performance exceeds EAGLE2. However, when l continues to increase, the increase in verification time cost will limit the effect. RASD (PLD) also shows the same trend. When l is only 2, the performance decreases compared to EAGLE2. This is because PLD requires additional time cost and does not bring enough length of retrieval candidates. On the right, since Medqa’s database is relatively small, the search results using only REST are not as good as those of EAGLE2. However, RASD (REST) can still surpass EAGLE2 with the help of REST. RASD achieves the best results when l is 8.

5.3 Effect of Datastore size

Method	Size	Time	SR	τ
EAGLE2	-	-	2.73	3.09
RASD(REST)	0.9 GB	0.2 ms	2.76	3.22
RASD(REST)	8.7 GB	0.6 ms	2.85	3.40
RASD(REST)	27 GB	0.7 ms	2.90	3.48

Table 4: The speedup with different datastore sizes in RASD(REST).

We explored how datastore size affects RASD performance in REST. Increasing the datastore size improves the accuracy of retrieved draft tokens, which significantly boosts generation speed. As shown in Table 4, both average acceptance length and speedup ratio improve with larger datastore sizes. However, the speedup growth is less pronounced than the increase in mean generated length. This difference may be due to the overhead of retrieving draft tokens. We assume that in industrial applications, there will be enough disk storage to build large data stores and enough CPU cores for fast retrieval. Therefore, there is still potential to achieve even faster speeds with a larger datastore.

6 Conclusion

In this work, we propose RASD: Retrieval-Augmented Speculative Decoding. We use retrieval methods to improve the quality of candidate sequences from draft models. Concretely, RASD improves the speedup of speculative sampling on the out-of-domain datasets that are difficult for the draft model to handle and improves the maximum output length of draft models. We develop a pruning method to select appropriate retrieval sequences and we fuse the sequence tree from the draft model with the retrieval tree, creating a final combined tree for verification. Experiments have proven that our method is effective and exhibits strong scalability.

7 Limitations

Based on our experiments and conclusions, we conclude some limitations of our work as follows:

- In our experiments, we only performed the draft model generation and retrieval process serially. parallelism of the two processes is not considered, which will result in the acceleration ratio not reaching the theoretical maximum value.

563 • We only consider pruning the retrieval tree by
564 controlling hyperparameters, which is not au-
565 tomatic and using the same hyperparameters
566 for different conversation types is not the best
567 choice.

568 Ethics Statement

569 We hereby declare that all authors of this article are
570 aware of and adhere to the provided ACL Code of
571 Ethics and honor the code of conduct.

572 **Use of Human Annotations** Human annotations
573 are only used in methodological research at the
574 beginning of the work, to assist in analyzing the
575 feasibility of the proposed solution. Annotators
576 consented to the use of data for research purposes.
577 We ensure that the privacy of all annotators is pro-
578 tected throughout the annotation process, and all of
579 them are adequately paid according to local stan-
580 dards. Human annotations are not applied during
581 the evaluation of our method.

582 **Risks** In this paper, all datasets are obtained from
583 official sources. The datasets adopted have been
584 anonymized and do not contain offensive infor-
585 mation. However, we cannot guarantee that the
586 datasets do not contain socially harmful or toxic
587 language.

588 References

589 Yushi Bai, Xin Lv, Jiajie Zhang, Hongchang Lyu,
590 Jiankai Tang, Zhidian Huang, Zhengxiao Du, Xiao
591 Liu, Aohan Zeng, Lei Hou, Yuxiao Dong, Jie Tang,
592 and Juanzi Li. 2024. [Longbench: A bilingual, multi-
593 task benchmark for long context understanding](#). In
594 *Proceedings of the 62nd Annual Meeting of the As-
595 sociation for Computational Linguistics (Volume 1:
596 Long Papers), ACL 2024, Bangkok, Thailand, Au-
597 gust 11-16, 2024*, pages 3119–3137. Association for
598 Computational Linguistics.

599 Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie
600 Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind
601 Neelakantan, Pranav Shyam, Girish Sastry, Amanda
602 Askell, Sandhini Agarwal, Ariel Herbert-Voss,
603 Gretchen Krueger, Tom Henighan, Rewon Child,
604 Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu,
605 Clemens Winter, Christopher Hesse, Mark Chen, Eric
606 Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess,
607 Jack Clark, Christopher Berner, Sam McCandlish,
608 Alec Radford, Ilya Sutskever, and Dario Amodei.
609 2020. [Language models are few-shot learners](#). In *Ad-
610 vances in Neural Information Processing Systems 33:
611 Annual Conference on Neural Information Process-
612 ing Systems 2020, NeurIPS 2020, December 6-12,
613 2020, virtual*.

Tianle Cai, Yuhong Li, Zhengyang Geng, Hongwu Peng,
Jason D. Lee, Deming Chen, and Tri Dao. 2024. [Medusa: Simple LLM inference acceleration frame-
work with multiple decoding heads](#). In *Forty-first In-
ternational Conference on Machine Learning, ICML
2024, Vienna, Austria, July 21-27, 2024*. OpenRe-
view.net.

Charlie Chen, Sebastian Borgeaud, Geoffrey Irv-
ing, Jean-Baptiste Lespiau, Laurent Sifre, and
John Jumper. 2023. [Accelerating large language
model decoding with speculative sampling](#). *CoRR*,
abs/2302.01318.

Mark Chen, Jerry Tworek, Heewoo Jun, Qiming
Yuan, Henrique Pondé de Oliveira Pinto, Jared Ka-
plan, Harri Edwards, Yuri Burda, Nicholas Joseph,
Greg Brockman, Alex Ray, Raul Puri, Gretchen
Krueger, Michael Petrov, Heidy Khlaaf, Girish Sas-
try, Pamela Mishkin, Brooke Chan, Scott Gray,
Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz
Kaiser, Mohammad Bavarian, Clemens Winter,
Philippe Tillet, Felipe Petroski Such, Dave Cum-
mings, Matthias Plappert, Fotios Chantzis, Eliza-
beth Barnes, Ariel Herbert-Voss, William Hebg-
Guss, Alex Nichol, Alex Paino, Nikolas Tezak, Jie
Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain,
William Saunders, Christopher Hesse, Andrew N.
Carr, Jan Leike, Joshua Achiam, Vedant Misra, Evan
Morikawa, Alec Radford, Matthew Knight, Miles
Brundage, Mira Murati, Katie Mayer, Peter Welinder,
Bob McGrew, Dario Amodei, Sam McCandlish, Ilya
Sutskever, and Wojciech Zaremba. 2021. [Evaluat-
ing large language models trained on code](#). *CoRR*,
abs/2107.03374.

Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng,
Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan
Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion
Stoica, and Eric P. Xing. 2023. [Vicuna: An open-
source chatbot impressing gpt-4 with 90%* chatgpt
quality](#).

Cunxiao Du, Jing Jiang, Xu Yuanchen, Jiawei Wu,
Sicheng Yu, Yongqi Li, Shenggui Li, Kai Xu, Liqiang
Nie, Zhaopeng Tu, and Yang You. 2024. [Glide with a
caPE: A low-hassle method to accelerate speculative
decoding](#). In *Forty-first International Conference on
Machine Learning*.

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey,
Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman,
Akhil Mathur, Alan Schelten, Amy Yang, Angela
Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang,
Archi Mitra, Archie Sravankumar, Artem Korenev,
Arthur Hinsvark, Arun Rao, Aston Zhang, Aurélien
Rodriguez, Austen Gregerson, Ava Spataru, Bap-
tiste Rozière, Bethany Biron, Binh Tang, Bobbie
Chern, Charlotte Caucheteux, Chaya Nayak, Chloe
Bi, Chris Marra, Chris McConnell, Christian Keller,
Christophe Touret, Chunyang Wu, Corinne Wong,
Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Al-
lonsius, Daniel Song, Danielle Pintz, Danny Livshits,
David Esiobu, Dhruv Choudhary, Dhruv Mahajan,
Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes,

674	Egor Lakomkin, Ehab AlBadawy, Elina Lobanova,	<i>Conference on Learning Representations, ICLR 2019,</i>	732
675	Emily Dinan, Eric Michael Smith, Filip Radenovic,	<i>New Orleans, LA, USA, May 6-9, 2019.</i>	733
676	Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Geor-	OpenRe-	734
677	gia Lewis Anderson, Graeme Nail, Grégoire Mialon,	view.net.	
678	Guan Pang, Guillem Cucurell, Hailey Nguyen, Han-	Xupeng Miao, Gabriele Oliaro, Zhihao Zhang, Xinhao	735
679	nah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov,	Cheng, Zeyu Wang, Rae Ying Yee Wong, Zhuom-	736
680	Imanol Arrieta Ibarra, Isabel M. Kloumann, Ishan	ing Chen, Daiyaan Arfeen, Reyna Abhyankar, and	737
681	Misra, Ivan Evtimov, Jade Copet, Jaewon Lee, Jan	Zhihao Jia. 2023. Specinfer: Accelerating generative	738
682	Geffert, Jana Vranes, Jason Park, Jay Mahadeokar,	LLM serving with speculative inference and token	739
683	Jeet Shah, Jelmer van der Linde, Jennifer Billock,	tree verification.	740
684	Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi,	<i>CoRR</i> , abs/2305.09781.	
685	Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu,	Ramesh Nallapati, Bowen Zhou, Cícero Nogueira dos	741
686	Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph	Santos, Çağlar Gülçehre, and Bing Xiang. 2016.	742
687	Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia,	Abstractive text summarization using sequence-to-	743
688	Kalyan Vasuden Alwala, Kartikeya Upasani, Kate	sequence rnns and beyond.	744
689	Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, and	In <i>Proceedings of the</i>	745
690	et al. 2024. The llama 3 herd of models.	<i>20th SIGNLL Conference on Computational Natural</i>	746
691	<i>CoRR</i> , abs/2407.21783.	<i>Language Learning, CoNLL 2016, Berlin, Germany,</i>	747
		<i>August 11-12, 2016</i> , pages 280–290. ACL.	
692	Xiangxiang Gao, Weisheng Xie, Yiwei Xiang, and	Apoorv Saxena. 2023. Prompt lookup decoding.	748
693	Feng Ji. 2024. Falcon: Faster and parallel inference	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob	749
694	of large language models through enhanced	Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz	750
695	semi-autoregressive drafting and custom-designed	Kaiser, and Illia Polosukhin. 2017. Attention is all	751
696	decoding tree.	you need.	752
	<i>CoRR</i> , abs/2412.12639.	In <i>Advances in Neural Information Pro-</i>	753
697	Zhenyu He, Zexuan Zhong, Tianle Cai, Jason D. Lee,	<i>cessing Systems 30: Annual Conference on Neural</i>	754
698	and Di He. 2024. REST: retrieval-based speculative	<i>Information Processing Systems 2017, December 4-9,</i>	755
699	decoding.	<i>2017, Long Beach, CA, USA</i> , pages 5998–6008.	
700	In <i>Proceedings of the 2024 Conference</i>	Bin Xiao, Chunan Shi, Xiaonan Nie, Fan Yang, Xiang-	756
701	<i>of the North American Chapter of the Association</i>	wei Deng, Lei Su, Weipeng Chen, and Bin Cui. 2024.	757
702	<i>for Computational Linguistics: Human Language</i>	Clover: Regressive lightweight speculative decoding	758
703	<i>Technologies (Volume 1: Long Papers), NAACL 2024,</i>	with sequential knowledge.	759
704	<i>Mexico City, Mexico, June 16-21, 2024</i> , pages 1582–	<i>CoRR</i> , abs/2405.00263.	
	1595. Association for Computational Linguistics.	An Yang, Baosong Yang, Beichen Zhang, Binyuan	760
705	Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick	Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayi-	761
706	S. H. Lewis, Ledell Wu, Sergey Edunov, Danqi Chen,	heng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian	762
707	and Wen-tau Yih. 2020. Dense passage retrieval for	Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang,	763
708	open-domain question answering.	Jiaxi Yang, Jingren Zhou, Junyang Lin, Kai Dang,	764
709	In <i>Proceedings of the 2020 Conference on Empirical Methods in Nat-</i>	Keming Lu, Keqin Bao, Kexin Yang, Le Yu, Mei	765
710	<i>ural Language Processing, EMNLP 2020, Online,</i>	Li, Mingfeng Xue, Pei Zhang, Qin Zhu, Rui Men,	766
711	<i>November 16-20, 2020</i> , pages 6769–6781. Associa-	Runji Lin, Tianhao Li, Tingyu Xia, Xingzhang Ren,	767
712	tion for Computational Linguistics.	Xuanheng Ren, Yang Fan, Yang Su, Yichang Zhang,	768
		Yu Wan, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and	769
713	Yaniv Leviathan, Matan Kalman, and Yossi Matias.	Zihan Qiu. 2024. Qwen2.5 technical report.	770
714	2023. Fast inference from transformers via spec-	<i>CoRR</i> , abs/2412.15115.	771
715	ulative decoding.		
716	In <i>International Conference on</i>		
717	<i>Machine Learning, ICML 2023, 23-29 July 2023,</i>		
718	<i>Honolulu, Hawaii, USA</i> , volume 202 of <i>Proceedings</i>		
719	<i>of Machine Learning Research</i> , pages 19274–19286.		
	PMLR.		
720	Yuhui Li, Fangyun Wei, Chao Zhang, and Hongyang		
721	Zhang. 2024a. EAGLE-2: faster inference of lan-		
722	guage models with dynamic draft trees.		
723	<i>CoRR</i> , abs/2406.16858.		
724	Yuhui Li, Fangyun Wei, Chao Zhang, and Hongyang		
725	Zhang. 2024b. EAGLE: speculative sampling re-		
726	quires rethinking feature uncertainty.		
727	In <i>Forty-</i>		
728	<i>first International Conference on Machine Learning,</i>		
729	<i>ICML 2024, Vienna, Austria, July 21-27, 2024.</i>		
	OpenReview.net.		
730	Ilya Loshchilov and Frank Hutter. 2019. Decoupled		
731	weight decay regularization.		
	In <i>7th International</i>		