

N-Gram Trie Speculative Decoding for Faster LLM In-Context Inference

Anonymous ACL submission

Abstract

As an important method of prompt engineering, In-Context Learning (ICL) provides generalization and knowledge enhancement capabilities for Large Language Models (LLMs) (Dong et al., 2024). However, the extensive length of retrieved contexts and the limited token throughput in autoregressive models constrain model’s reasoning speed. To resolve this constraint, we propose **N-Gram-Trie**, a novel approach that exploits the potential overlap between the context and model output. The strategy utilizes context to construct an n-gram trie. The trie will be used to construct drafts that benefit the LLM to increase the speed of token generation. We evaluate our method on the summarization, Retrieval-Augmented Generation (RAG) and context Question Answering (context QA) tasks. Experiment results on Vicuna-7B, Llama2-7B-Chat, and Llama3-8B-Instruct all demonstrate a significant increase in speed without compromising accuracy. In comparison experiments, our method achieves the best mean speedup among various baselines.

1 Introduction

In-Context Learning (ICL) has emerged as a transformative paradigm in the field of prompt engineering, fundamentally reshaping how Large Language Models (LLMs) adapt to and perform on diverse tasks. By leveraging context information provided within the input prompt, ICL enables LLMs to generalize across tasks and domains without requiring task-specific fine-tuning. This capability has profound implications for the scalability and versatility of LLMs, allowing them to excel in various applications, such as context question answering, summarization and Retrieval-Augmented Generation (RAG). The ability to dynamically incorporate contextual knowledge has made ICL a cornerstone of modern LLM deployment, driving advancements in both academic research and industrial applications.

Despite its remarkable success, ICL faces a significant challenge: the extensive length of retrieved contexts and the inherent limitations of autoregressive token generation will result in slow reasoning speeds. As the complexity and length of context information increase, the computational overhead grows substantially, leading to delays in token generation and reduced efficiency. This bottleneck is particularly problematic in real-time applications, such as interactive systems or large-scale retrieval-augmented tasks, where speed is critical. Addressing this issue is essential to unlocking the full potential of ICL and enabling its broader adoption in time-sensitive scenarios.

Speculative decoding (Leviathan et al., 2023; Cai et al., 2024; Li et al., 2024; He et al., 2023; Luo et al., 2024) can effectively accelerate model inference. This approach employs a smaller, faster draft model to predict potential token sequences, which are then verified by the larger target model in parallel. By reducing the number of sequential decoding steps required by the target model, speculative decoding achieves significant speedups while maintaining output quality. However, this method often requires additional computational resources and careful tuning to balance the trade-off between speed and accuracy. REST (He et al., 2023) employs an external corpus to generate draft tokens, where the output tokens serve as prefixes to search for matching suffixes within the corpus. However, the excessive reuse of nodes and the global corpus tire reduce the acceptance rate of draft tokens. Lookahead Decoding (Fu et al., 2024) utilizes n -gram token histories as drafts for verification. While this method shows promise, its utility is primarily confined to scenarios where output tokens exhibit repetitive patterns, restricting its applicability in more diverse or dynamic contexts.

We propose N-Gram-Trie, a novel approach designed to accelerate token generation by exploiting the overlap between the context and the model’s

output. Then a trie is constructed by using the set of prefixes and suffixes. Build a trie from the prefix and suffix sets. In the model prediction stage, the draft is constructed through the nodes in the trie, which significantly improves reasoning speed without compromising output quality.

We evaluate our approach on summarization (Nallapati et al., 2016), Retrieval-Augmented Generation (Xia et al., 2024; Joshi et al., 2017) and context Question Answering (context QA) (Kamalloo et al., 2023) tasks. Multiple base models including Vicuna-7B (Zheng et al., 2023), Llama2-7B-Chat (Touvron et al., 2023) and Llama3-8B-Instruct (AI@Meta, 2024) are selected to be tested. Experiment results show that our method exhibits remarkable speedups on multiple models (mean 2.27x on Vicuna-7B, 2.10x on Llama2-7B-Chat and 1.56x on Llama3-8B-Instruct). Through the experiment comparison of the inference effect of the model, we prove that our method can accelerate the model in the process of context prompt inference without affecting the inference ability of the base model. We also conduct many further experiments around the speedup effect. This work not only addresses a critical limitation of ICL but also provide a effective method for more efficient and scalable deployment of LLMs in real-world applications.

The contribution of this paper can be summarized as follows:

- We propose an n-gram trie speculative decoding method. It can effectively use the potential overlap of the context and output tokens to accelerate model inference speed.
- We design a novel n-gram trie construction method. The trie constructed by n-gram sampling can effectively improve the acceptance rate of the draft.
- We conduct extensive experiments on several models. It shows our excellent acceleration effect on summarization, RAG and context QA tasks.

2 Related work

2.1 In-Context Learning

In-Context Learning (ICL) is an approach which makes LLMs perform better on specific-domain task. By giving only a few examples or hints, LLMs can find the underlying patterns of the context and answer the question correctly. (Dong et al., 2024).

There are many approaches that can be applied to ICL. (Gu et al., 2023) extract the context by pre-training in a large corpus that contains long context. (Wei et al., 2023) propose symbol tuning, which uses tagged symbols as fine-tuned data for LLMs to study. (Wei et al., 2022) leverages instruction tuning in LLMs to enhance the zero-shot learning in LLMs.

Also, there are also large variety of downstream applications in the In Context learning. Prompt engineering is one of them. We can write an accurate prompt to make LLMs easier to understand the downstream tasks and give a satisfying answer. Prompt engineering are widely used in downstream tasks, such as Context QA, RAG, Few-shot Learning and Summary. Context QA (Kamalloo et al., 2023) tasks need LLMs to read the context and find the potential answers. Concatenating the context and question as prompts, LLMs can read them and give an answer in a efficient way. Like Context QA, RAG (Li et al., 2023) also needs retrieved context to carry out user’s query. In Few-shot Learning, some examples about downstream tasks are usually given. LLMs can study the potential patterns between them and complete the task based on the given pattern. Summarization also needs the ability of context-reading.

2.2 Speculative Decoding

Speculative decoding (Leviathan et al., 2023) has been first proposed to ease the problem of throughput in LLM generation. Using a small draft model to *explore the token way*, target LLM just need to verify in one step without calculating repeatedly for getting these tokens.

Now, many speculative methods are based on the *guess and verify* approach. For example, Medusa (Cai et al., 2024) uses some trained Medusa head to predict the next n -tokens, but the prediction is not continuous and it degrades accept rate. Based on Medusa (Cai et al., 2024), Hydra (Ankner et al., 2024) take the continuation of the draft into consideration. The draft head can predict tokens with However, both Medusa and Hydra need extra training cost for draft models. Also, some works focus on the reusing of the former tokens or the external corpus. For instance, REST (He et al., 2023) uses an external corpus as draft. The output tokens is used as prefix to search for the suffix in the corpus. But the smaller corpus decreases accept rate while the bigger corpus makes REST harder to find the right suffix.

Lookahead decoding (Fu et al., 2024) uses n -gram token history as draft to verify. But it is useful only when the output token is repeatedly generated. LLMA (Yang et al., 2023) also try to use the overlap between the input and output, but it simply use all the retrieved tokens to make evaluation without clipping the drafts. PLD (Saxena, 2023) copy prompts as the drafts, but it simply copy the first-match suffix without matching all the potential suffixes in the prompt.

2.3 Tree Attention

Tree attention (Miao et al., 2023) is proposed to solve the problem how a tree-structured token sequences can be decoding in parallel. By using an attention mask, the drafts can be easily integrated in one mask in inference. In the attention mask, Now tree attention is widely used in multi-draft verification.

SpecInfer (Miao et al., 2024) uses some small draft models to independently predict the potential tokens sequences, the tokens will then be clipped and put in the attention masks. Medusa (Cai et al., 2024) uses some positional draft heads to predict the top-k tokens in the next i place. It uses attention mask to integrate the top-k tokens into token sequences for prediction. REST (He et al., 2023) retrieved many tokens in a big suffix-array dataset. After clipping the tokens, He et al. also use tree attention mask to make a trie tree for faster decoding.

3 Proposed Method

The structure of N-Gram-Trie is shown in Figure 1. In the in-context prompt tasks, we first build an n -gram trie based on the context. The tree records the dependencies between preceding and following tokens of context. Subsequently, in the process of model inference, the draft of model inference is constructed by speculative decoding through the dependencies of n -gram trie, which can accelerate model inference speed.

3.1 N-Gram Trie Construction

Trie is a tree structure used to store and retrieve strings efficiently by organizing tokens in a prefix-based hierarchy. Its key advantage is faster suffix finding, which makes it suitable for speculative decoding (He et al., 2023). However, traditional Trie relies on massive documents to build for higher acceptance rate. It is difficult to construct an effective

Algorithm 1 Trie Generation

Input: T : Token list
 D : Collected n -gram sample results
 L_p : Maximum prefix length

Output: τ

- 1: Init $root$ as τ \triangleright an empty root node
- 2: **for** $\langle P_i, S_i, f \rangle \in D$ **do**
- 3: **for** $j \in (0, L_p)$ **do**
- 4: $subprefix \leftarrow P_i[j : L_p]$
- 5: $key \leftarrow subprefix + S_i$
- 6: $node \leftarrow root$
- 7: **for** $t \in key$ **do**
- 8: **for** $child$ in $node.children$ **do**
- 9: **if** $t = child.token$ **then**
- 10: $node \leftarrow child$
- 11: $node.frequency.update(f)$
- 12: **end if**
- 13: **end for**
- 14: **if** $t \notin node.children$ **then**
- 15: $new \leftarrow Node(t, f)$
- 16: $node.children.insert(new)$
- 17: $node \leftarrow new$
- 18: **end if**
- 19: **end for**
- 20: **end for**
- 21: **end for**
- 22: **return** τ

retrieval scheme in the case of a small amount of corpus. To this end, we design n -gram trie, sampled by n -gram sliding window, and then used the sampling results to build the trie. This method can effectively improve the efficiency and accuracy of suffix retrieval by constructing additional dependency chains.

N-Gram Sampling Specifically, for the context token list $T = \{t_1, t_2, \dots, t_l\}$, we set a sliding window of n -grams for sampling. The sampling length is n . The sliding window moves token by token from the beginning to the end over T . In the sliding window workspace, we set a maximum prefix length L_p to split tokens in the window. The split part will be the prefix part and the suffix part of the segment tokens. The prefix P_i and suffix S_i can be expressed as follows:

$$\left. \begin{aligned} P_i &= \{t_i, t_{i+1}, \dots, t_{i+L_p-1}\} \\ S_i &= \{t_{i+L_p}, t_{i+L_p+1}, \dots, t_{i+n-1}\} \end{aligned} \right\} i \in [1, l], \quad (1)$$

where i denotes the start index of the window. We establish the dependency between the prefix and

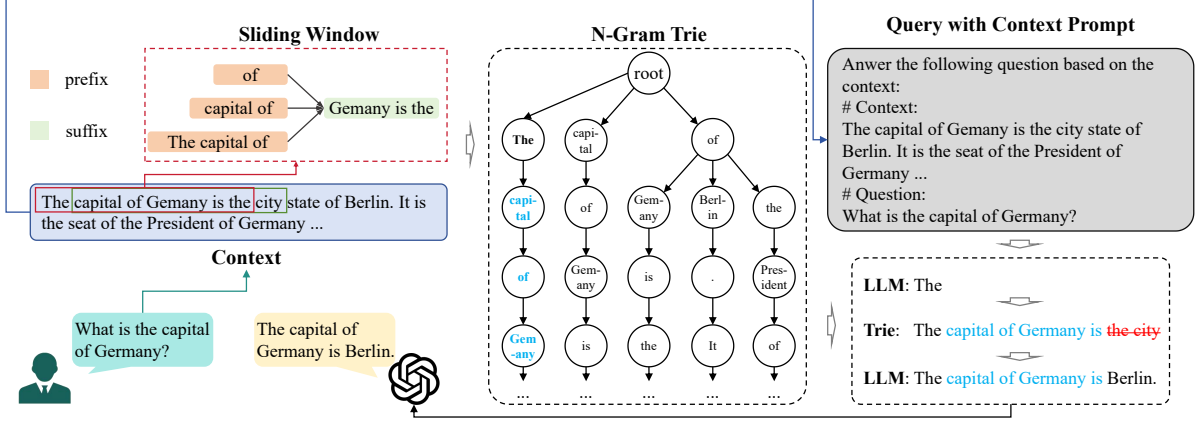


Figure 1: The structure of N-Gram-Trie. We sample through a sliding window of n-grams and get the prefixes and suffixes from the documents in that window. A trie can be constructed based on the set of prefixes obtained by window sliding sampling. In the process of model inference, the trie is used for speculative decoding to quickly predict the model output. The n in the n-gram sampling in the example of the figure is 6 and the maximum prefix length L_p is 3.

suffix for each tokens group, and obtain the dependency set D by sliding window sampling. D can be defined in the following form:

$$D = \{ \langle P_i, S_i, f \rangle \mid i \in [1, l] \}, \quad (2)$$

where f is the frequency of dependency $\langle P_i, S_i \rangle$ during the sampling process.

Trie Construction We build trie τ based on the sample results D and the construction process is as shown in Algorithm 1.

Specifically, for the prefix P_i in the sample set D , we traverse and split it according to the maximum prefix length to obtain its sub-prefixes SP_i . The process can be defined as:

$$\begin{aligned} SP_i &= \{ SP_{i,j} \mid j \in (0, L_p) \} \\ &= \{ P_i[j : L_p] \mid j \in (0, L_p) \}, i \in [1, l], \end{aligned} \quad (3)$$

where $j \in (0, L_p)$ denotes the cut length of the sub-prefix. By constructing additional prefix nodes, the corresponding prefix can be effectively found according to the model output in the retrieval process.

We take the dependency of each subprefix and its suffix as the basic unit for trie insertion. During insertion, the token t is used as the basic units of the tree nodes. We iterate from the root node, sharing a node for the same token. If there is no corresponding token in the current nodes, insert an additional token. The insertion logic is as follows:

$$node = \begin{cases} child, & \text{if } t \in node.children \\ node_t, & \text{if } t \notin node.children \end{cases}, \quad (4)$$

where $child$ is the child of $node$ and $child.token = t$, $node_t$ is a new node built by t and inserted into the children of the original $node$. In this way, we let suffix nodes with the same prefix share the same prefix.

Note that we also record the frequency f of each node as it is inserted, in order to provide a priority reference for subsequent retrieval. Finally, by exploiting the samples in D , we can construct an efficient and accurate n-gram trie τ .

3.2 Draft Collecting and Matching

As shown in the gray area in Figure 1, in-context learning combines context with user query through templates in the prompt engineering. The query with context will serve as the reasoning basis for the target model. We define the tokens that have been generated by the s time step target model as $T_s = \{t_1, t_2, \dots, t_k\}$. We will build the draft after s time step through the n-gram trie τ constructed in the former subsection that stores prefix and suffix dependency of context. Then, the target model will verify and revise the draft.

Draft Construction When searching for the draft, we firstly extract the suffix of new tokens T_s for prefix matching. At first, the length of the prefix token will be set to L_p . If T_s matches the prefix chain in τ , we can extract the suffix of this prefix and break matching. If not found, we subtract one token from prefix tokens until match the or prefix tokens length is 0. Then, we can obtain a suffix tree τ_s that matches the generated tokens

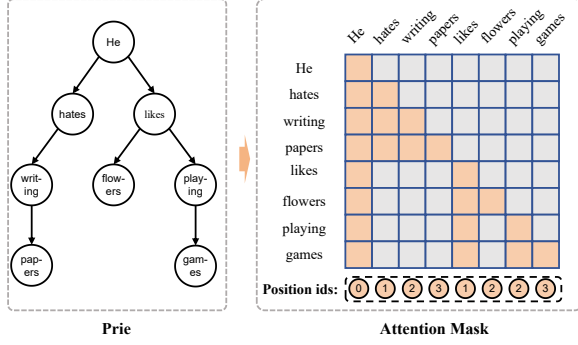


Figure 2: An Example of Tree Attention. The tokens in the orange part of the attention mask are visible to each other, and the tokens in the gray part are invisible to each other

T_s of the target model. To improve the acceptance rate of the draft, we prune the suffix tree according to the frequency f of nodes and extract nodes with lower f . The draft tree is not always very big, so sometimes the pruning is not used.

Specifically, referring to (He et al., 2023) and (Cai et al., 2024), we set a min-heap for storage of suffix chains. For each node v_k in τ_s , we build a draft d_k based on its path chain with the root node of τ_s . The priority of the draft is determined by the frequency of v_k . This process can be expressed as:

$$d_k = \langle \text{Path}(v_r, v_k), f_k \rangle \\ = \langle \{v_r, v_1, \dots, v_i, \dots, v_k\}, f_k \rangle, \quad (5) \\ i \in [1, k], v_i \in \tau_s,$$

where $\text{Path}(v_r, v_k)$ means the nodes from node v_r to node v_k . v_r is the root node of τ_s and f_k is the frequency of v_k .

Then, v_k will be placed in the min-heap in order of priority. Finally, alternative drafts are retained according to the length of min-heap. In this way, redundant nodes can be effectively removed and the pruning of suffix tree τ_s can be realized.

Model Verification Figure 2 shows an example of tree attention verifying the draft trie. For the draft trie τ_s , deep traverse it to obtain a linear list of tokens. In order to realize the tree attention, we set the same position id for the nodes of the same level. The specific form can be expressed as:

$$p_i = \text{Level}(v_i) + h, v_i \in \tau_s, \quad (6)$$

where p_i is the position id of t_i . $\text{Level}(v_i)$ is the level of v_i . h is the length of the preceding model tokens. This makes the tokens in each chain of the trie continuous.

Then, following tree attention meathod, we use attention mask to convert the draft tree into a 2-dimension mask m . For any tokens t_i and t_j , $m_{i,j}$ is 0 if there is a relationship between v_i and v_j in τ_s , otherwise it is 1. By matching the mask and position ids. The target model can verify multiple branches of trie simultaneously.

4 Experiments

4.1 Experiment Setting

We implement all the experiments on one NVIDIA RTX 4090 with python version 3.9. All the experiments are run on greedy decoding. The pytorch version is 2.5.1 with CUDA version is 12.2.

4.1.1 Baselines

We choose the baselines provided on the Speculate Bench (Xia et al., 2024): vanilla inference without any speculative methods, speculative Sampling (Chen et al., 2023), Medusa (Cai et al., 2024), SPACE (Yi et al., 2024), Hydra (Ankner et al., 2024), Lookahead (Fu et al., 2023) and REST (He et al., 2023). For Speculative Sampling, we used Llama-68m (Miao et al., 2024) as draft model to match Llama2-7b and use Vicuna-68m to match Vicuna-7b-v1.3. For Lookahead and REST, we simply use the same experiment setup in SpecBench(Xia et al., 2024).

4.1.2 Datasets

For datasets, we choose RAG, summary in Specbench (Xia et al., 2024). The RAG dataset contains 80 data from Natural Questions. Five retrieved documents from Wikipedia (Li et al., 2023) are concatenated. (Kwiatkowski et al., 2019) and the summary dataset is randomly chosen by CNN/Daily Mail (Nallapati et al., 2016). In addition, we make TriviaQA (Joshi et al., 2017) dataset for additional RAG task and make Hagrid (Kamalloo et al., 2023) dataset for context QA task. For TriviaQA task, we use bge-m3 (Chen et al., 2024a) and bge-reranker-v2-m3 (Chen et al., 2024b) to search for 5 relevant documents in Wikipedia corpus. For Hagrid task, we simply concatenate the given context and the question.

4.1.3 Base models

To conduct the experiments, we use three models for validation. One is Vicuna-7B-v1.3 (Zheng et al., 2023), One is Llama-2-7B-chat (Touvron et al., 2023) and the other is Llama-3-8B-Instruct (AI@Meta, 2024).

Model	Method	Spec-Bench		TriviaQA	Hagrid	Mean Speedup
		Summary	RAG			
Vicuna-7B	Vanilla	1.00× (1.00)	1.00× (1.00)	1.00× (1.00)	1.00× (1.00)	1.00x
	SpS	1.69× (2.44)	1.59× (2.30)	1.74× (2.49)	1.40× (2.46)	1.61x
	Medusa	1.48× (2.01)	1.45× (2.08)	1.45× (2.03)	<u>1.56× (2.17)</u>	1.49x
	SPACE	1.69× (2.26)	1.47× (1.91)	1.57× (2.26)	<u>1.26× (2.05)</u>	1.50x
	Hydra	1.86× (2.70)	<u>1.88× (2.90)</u>	<u>1.86× (2.84)</u>	<u>1.52× (2.98)</u>	<u>1.78×</u>
	Lookahead	1.29× (1.54)	1.19× (1.48)	1.27× (1.46)	0.95× (1.47)	1.18x
	REST	1.13× (1.65)	1.32× (1.89)	1.31× (1.71)	1.33× (1.82)	1.27x
	Ours	<u>1.75× (2.39)</u>	3.48× (5.19)	1.92× (2.36)	1.94× (3.07)	2.27×
Llama2-7B-Chat	Vanilla	1.00× (1.00)	1.00× (1.00)	1.00× (1.00)	1.00× (1.00)	1.00x
	SpS	1.25× (1.54)	<u>1.47× (1.91)</u>	1.35× (1.86)	1.38× (1.62)	1.36x
	Lookahead	1.44× (1.59)	1.40× (1.63)	<u>1.53× (1.71)</u>	1.38× (1.97)	<u>1.44×</u>
	REST	1.03× (1.54)	1.14× (1.91)	1.22× (1.68)	<u>1.42× (1.47)</u>	1.20x
	Ours	<u>1.28× (1.76)</u>	3.62× (5.00)	1.89× (2.88)	1.61× (2.34)	2.10×
Llama3-8B-Instruct	Vanilla	1.00× (1.00)	1.00× (1.00)	1.00× (1.00)	1.00× (1.00)	1.00x
	Lookahead	1.25× (1.60)	<u>1.18× (1.51)</u>	<u>1.58× (1.54)</u>	<u>1.38× (1.73)</u>	<u>1.50×</u>
	REST	0.93× (1.54)	1.14× (1.91)	1.13× (1.61)	1.02× (1.69)	1.05x
	Ours	<u>1.06× (1.42)</u>	1.77× (2.11)	1.75× (1.86)	1.68× (2.34)	1.56×

Table 1: Speedup Ratio and Accept Length Comparison. The data on the left means speedup and the data on the right means average accept length. The best performance for each metric is highlighted in **bold** font, while the second-best performance is indicated with an underline.

4.1.4 Hyperparameters

In the experiment, there are two hypermeter that need to be tuned: matched prefix L_p and gram-length n . So we conduct the experiment to test the efficiency. The details can be seen in table3. Also, we use FAISS (Douze et al., 2025) to store the embedding of the corpus using IVF-PQ method. The parameter of the number of clusters is 4096 and the parameter that the vector will be separated is 64. The clusters that will be searched is set to 16. We firstly encode all the corpus text using bge-m3 (Chen et al., 2024a), and search top-100 relevant texts for questions in triviaQA (Joshi et al., 2017). Then we rerank the texts using bge-reranker-v2-m3 (Chen et al., 2024b) to get the top-5 relevant contexts.

4.1.5 Metrics

Like other speculative decoding, we use *average accept length*, *mean speedup* and sentence F1-score in our evaluation. Average accept length shows the length that the drafts are accepted in every decoding step. Usually the accept length is higher, the speedup can be higher. Mean speedup indicates the speedup of tokens throughput compared with decoding without any speculative method (baseline).

Sentence F1-score is used to evaluate whether the output result is the same. Because of the ran-

dom possibilities, the results may vary.

4.2 Main Results

The experiment results can be seen in Table 1. We can see that our method achieves optimal acceleration results compared to the baseline for all tasks except the summarization task. On average, the mean speedup of our method has achieved a meaningful improvement over the baselines (2.27× on Vicuna-7B, 2.10× on Llama2-7B-Chat and 1.56× on Llama3-8B-Instruct). It is worth noting that our method performs better than REST (He et al., 2023) on each model and task in the same speculative decoding with trie, demonstrating the superiority of our n-gram trie.

4.2.1 RAG Task

The experiment result on Spec-Bench RAG dataset show that the accept length of the drafts achieves 5.19 on Vicuna-7B, 5.00 on Llama2-7B-Chat and 2.11 on Llama3-8B-Instruct, making the speedup rate achieve 3.48×, 3.62× and 1.77×. Its acceleration performance is much better than that of the basic method REST. Experiment results in multiple models show that our method has the strongest speedup effect. It is far ahead of second place on all models.

On TriviaQA dataset, the speedup rate of our

Hillary Clinton's security detail has added a second " Scooby " van to her motorcade , raising questions about the need for such an elaborate security measure . The second van , a GMC , is mechanically identical to the first van , a Chevrolet , but has different license plates. The Secret Service has employed de co y vehicles to confuse and discourage would-be attackers , but the use of two identical vans has raised eyebrows . The vans were seen driving separately to Clinton's appointed location before leaving together in a seven-car motorcade. The Secret Service has declined to comment on the security arrangements for dignitaries. The use of two Scooby vans has been observed in other instances , including when President Barack Obama returns to the White House after long trips . The Secret Service frequently deploys duplicates of aircraft and cars it uses to transport VIPs , including Marine One, the president's customized helicopter, which usually travels with two decoys . </s>

Figure 3: A Case on Summary Dataset. The inference exploration in one step. In the figure, the words are separated in tokens. Yellow text is generated by n-gram trie.

method achieves $1.92\times$ on Vicuna-7B, $1.89\times$ on Llama2-7B-Chat and $1.75\times$ on Llama3-8B-Instruct. The speed-up performance is also the best. Even though the accept length of our approach is smaller than Hydra (Ankner et al., 2024) on Vicuna-7B, we still have a better throughput performance in this task.

4.2.2 Context QA Task

Our approach achieves the best results on all models ($1.94\times$ on Vicuna-7B, $1.61\times$ on Llama2-7B-Chat and $1.68\times$ on Llama3-8B-Instruct) on Hagrid dataset, which outperforms other approaches by $0.19\times$ - $0.99\times$. Compared to the basic method REST, we have more speedup on all models. This fully demonstrates the advantages brought by n-gram trie.

4.2.3 Summary Task

The performance of our method on Spec-Bench Summary dataset is not the best. This is mainly due to the fact that in this task, the target drafts are often put in a certain place, which usually don't need to get all the drafts in global. Global draft-getting method will result in wrongly draft clipping and unnecessarily draft-searching. But our method still ranks second ($1.75\times$ on Vicuna-7B, $1.28\times$ on Llama2-7B-Chat and $1.06\times$ on Llama3-8B-Instruct) in terms of speedup and outperform REST.

4.2.4 Model Response Quality

In order to evaluate the effect of accelerated inference on model accuracy, we conduct precision tests on three base models. The experiment results are shown in Table 2. We can observe that the responses from our method are almost identical to those of the original model. This is because the draft of the speculative decoding is verified by the target model. When inconsistent tokens appear, the output will be corrected based on the output of the target model. Therefore, this method does not affect the accuracy of the model.

Dataset	Spec-Bench	Hagrid	TriviaQA
Vicuna-7B			
Vanilla	1.0000	1.0000	1.0000
SpS	0.9871	0.9936	0.9954
Medusa	0.9667	0.9905	0.9922
Hydra	0.9726	0.9894	0.9972
Lookahead	0.9857	0.9948	0.9939
REST	0.9695	0.9915	0.9922
Ours	0.9756	0.9942	0.9891
Llama2-7B-Chat			
Vanilla	1.0000	1.0000	1.0000
SpS	0.9838	0.9811	0.9866
Lookahead	0.9787	0.9754	0.9769
REST	0.9702	0.9792	0.9815
Ours	0.9739	0.9546	0.9743
Llama3-8B-Instruct			
Vanilla	1.0000	1.0000	1.0000
Lookahead	0.9837	0.9947	0.9913
REST	0.9888	0.9961	0.9951
Ours	0.9682	0.9619	0.9724

Table 2: The Sentence F1-score between Every Method and Baseline. Because of many random possibilities, the output may not be always the same, the results vary between 95 percent and 100 percent.

4.3 Case Study

To fully demonstrate the speedup effect of our method, we conduct a case study on the Summary dataset. The example is shown in Figure 3. As can be seen from the figure, our speculative decoding method based on n-gram trie correctly predict the large model output many times. A large number of useful drafts provide an effective speedup scheme for in-context based model inference.

4.4 Hyperparameter Analysis

In this section, we will will conduct experiments on the hyperparameters n and L_p in our method to get the best hyperparameter configuration. We use the RAG task of the Spec-bench (Xia et al., 2024) on Llama2-7B-Chat to test the performance

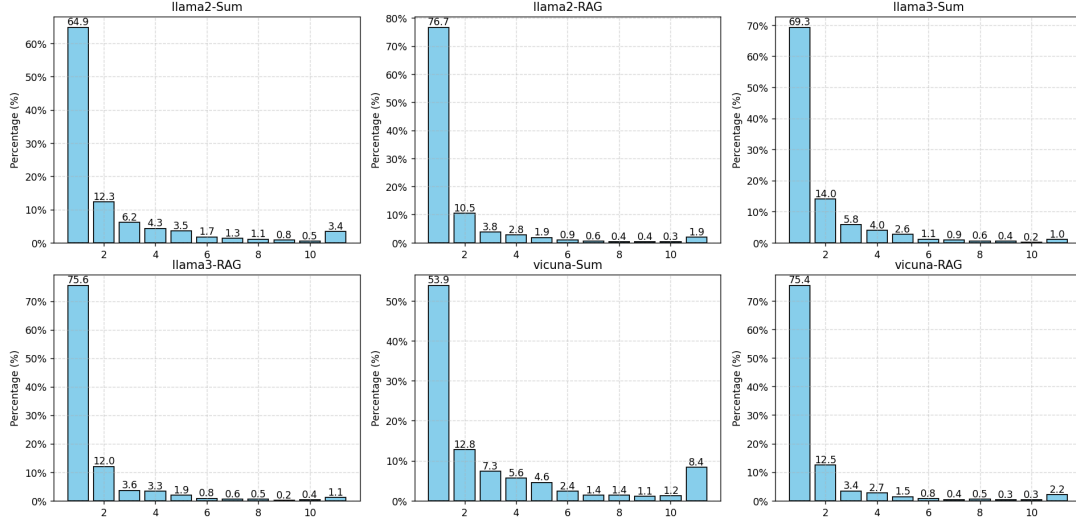


Figure 4: The accept length percentage in summary task and RAG tasks between models. The frequency of smaller accept length is usually larger except in the longest accept length

n/L_p	2	3	4	5
8	75.75	71.48	63.49	54.06
9	70.22	68.96	72.02	68.34
10	83.05	80.05	75.65	74.56
11	70.69	82.68	76.46	74.86
12	74.21	82.49	72.55	74.64
13	87.45	92.46	88.02	85.52
14	91.52	78.64	84.80	74.06
15	80.25	75.09	77.41	73.99
16	73.72	87.51	83.77	89.34

Table 3: Token Output Speed. The value in the table is token output number per second with different n and L_p

to 3.

4.5 Further Study

In order to explore the distribution of acceptance length of different models. We test Vicuna-7B (Zheng et al., 2023), Llama2-7B-Chat (Touvron et al., 2023), and Llama3-8B-Instruct (AI@Meta, 2024) on the Spec-Bench (Xia et al., 2024) dataset. The experimental results are shown in Figure 4. In this figure, it can be seen that the accept length concentrates in 1 (which means that no tokens are accepted). Besides, most of accept length is smaller than 4. And the percentage of the accept length decreases except in accept length = 11.

5 Conclusion

In this paper, we propose N-Gram-Trie, which reuse the context to build drafts to enhance in-context based model inference. By n-gram sampling from the context, we can obtain the prefix and suffix dependency set. We use trie to modeling the context based on the dependency set and retrieval on the tire to build the speculative decoding draft. By utilizing the overlap of context and model output, our approach effectively accelerates in-context based model inference. We test the summary, RAG, and context QA tasks on multiple large language models. Experiment results on multiple datasets show that our method can greatly improve the model inference speed in the context prompt domain without affecting the output quality.

of N-Gram-Trie. We try the value between 8-16 for n-gram length n and 2-5 for maximum prefix length L_p . The performance of the N-Gram-Trie with different hyperparameter is shown in the Figure 3. We can find that when n is small, the speedup effect will gradually deteriorate with the increase of L_p . We think this is because the excessively long L_p limits the length of the suffix, which in turn reduces the acceleration ability. When n is large, the token generation speed first speeds up and then slows down as L_p increases. This is mainly because when the suffix is not short, the longer prefix can better match the token of the model inference. Furthermore, it can be proved that when n value is large, appropriate redundant nodes can effectively improve the acceleration effect of speculative decoding. Statistically, we can see that the best choice of n is 13 and the maximum prefix length L_p is set

6 Limitations

This approach presents several limitations. First, while the trie-based generation and search mechanism offers efficiency advantages, its current implementation has suboptimal aspects. A key issue arises when multiple suffix candidates share identical frequency scores, which may lead to the premature elimination of potentially useful draft outputs due to the fixed threshold imposed by the `num_draft` parameter. Second, the method exhibits strong dependency on the quality of external retrieved corpora - performance degradation becomes inevitable when processing noisy or irrelevant retrieval results. To address these challenges, our future work will focus on developing enhanced trie construction algorithms that incorporate more sophisticated frequency weighting schemes and context-aware candidate selection strategies.

References

- AI@Meta. 2024. [Llama 3 model card](#).
- Zachary Ankner, Rishab Parthasarathy, Aniruddha Nrusimha, Christopher Rinard, Jonathan Ragan-Kelley, and William Brandon. 2024. [Hydra: Sequentially-dependent draft heads for medusa decoding](#). *Preprint*, arXiv:2402.05109.
- Tianle Cai, Yuhong Li, Zhengyang Geng, Hongwu Peng, Jason D Lee, Deming Chen, and Tri Dao. 2024. Medusa: Simple llm inference acceleration framework with multiple decoding heads. *arXiv preprint arXiv:2401.10774*.
- Charlie Chen, Sebastian Borgeaud, Geoffrey Irving, Jean-Baptiste Lespiau, Laurent Sifre, and John Jumper. 2023. [Accelerating large language model decoding with speculative sampling](#). *Preprint*, arXiv:2302.01318.
- Jianlv Chen, Shitao Xiao, Peitian Zhang, Kun Luo, Defu Lian, and Zheng Liu. 2024a. [Bge m3-embedding: Multi-lingual, multi-functionality, multi-granularity text embeddings through self-knowledge distillation](#). *Preprint*, arXiv:2402.03216.
- Jianlv Chen, Shitao Xiao, Peitian Zhang, Kun Luo, Defu Lian, and Zheng Liu. 2024b. [Bge m3-embedding: Multi-lingual, multi-functionality, multi-granularity text embeddings through self-knowledge distillation](#). *Preprint*, arXiv:2402.03216.
- Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Jingyuan Ma, Rui Li, Heming Xia, Jingjing Xu, Zhiyong Wu, Tianyu Liu, Baobao Chang, Xu Sun, Lei Li, and Zhifang Sui. 2024. [A survey on in-context learning](#). *Preprint*, arXiv:2301.00234.
- Matthijs Douze, Alexandr Guzhva, Chengqi Deng, Jeff Johnson, Gergely Szilvasy, Pierre-Emmanuel Mazaré, Maria Lomeli, Lucas Hosseini, and Hervé Jégou. 2025. [The faiss library](#). *Preprint*, arXiv:2401.08281.
- Yichao Fu, Peter Bailis, Ion Stoica, and Hao Zhang. 2023. [Breaking the sequential dependency of llm inference using lookahead decoding](#).
- Yichao Fu, Peter Bailis, Ion Stoica, and Hao Zhang. 2024. Break the sequential dependency of llm inference using lookahead decoding. *arXiv preprint arXiv:2402.02057*.
- Yuxian Gu, Li Dong, Furu Wei, and Minlie Huang. 2023. [Pre-training to learn in context](#). *Preprint*, arXiv:2305.09137.
- Zhenyu He, Zexuan Zhong, Tianle Cai, Jason D Lee, and Di He. 2023. [Rest: Retrieval-based speculative decoding](#). *Preprint*, arXiv:2311.08252.
- Mandar Joshi, Eunsol Choi, Daniel Weld, and Luke Zettlemoyer. 2017. [triviaqa: A Large Scale Distantly Supervised Challenge Dataset for Reading Comprehension](#). *arXiv e-prints*, arXiv:1705.03551.
- Ehsan Kamalloo, Aref Jafari, Xinyu Zhang, Nandan Thakur, and Jimmy Lin. 2023. HAGRID: A human-llm collaborative dataset for generative information-seeking with attribution. *arXiv:2307.16883*.
- Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, Kristina Toutanova, Llion Jones, Matthew Kelcey, Ming-Wei Chang, Andrew M. Dai, Jakob Uszkoreit, Quoc Le, and Slav Petrov. 2019. [Natural questions: A benchmark for question answering research](#). *Transactions of the Association for Computational Linguistics*, 7:452–466.
- Yaniv Leviathan, Matan Kalman, and Yossi Matias. 2023. Fast inference from transformers via speculative decoding. In *International Conference on Machine Learning*, pages 19274–19286. PMLR.
- Chaofan Li, Zheng Liu, Shitao Xiao, and Yingxia Shao. 2023. [Making large language models a better foundation for dense retrieval](#). *Preprint*, arXiv:2312.15503.
- Yuhui Li, Fangyun Wei, Chao Zhang, and Hongyang Zhang. 2024. EAGLE: Speculative sampling requires rethinking feature uncertainty. In *International Conference on Machine Learning*.
- Xianzhen Luo, Yixuan Wang, Qingfu Zhu, Zhiming Zhang, Xuanyu Zhang, Qing Yang, Dongliang Xu, and Wanxiang Che. 2024. [Turning trash into treasure: Accelerating inference of large language models with token recycling](#). *Preprint*, arXiv:2408.08696.
- Xupeng Miao, Gabriele Oliaro, Zhihao Zhang, Xinhao Cheng, Zeyu Wang, Zhengxin Zhang, Rae Ying Yee Wong, Alan Zhu, Lijie Yang, Xiaoxiang Shi, Chunan Shi, Zhuoming Chen, Daiyaan Arfeen, Reyna

- Abhyankar, and Zhihao Jia. 2024. [Specinfer: Accelerating large language model serving with tree-based speculative inference and verification](#). In *Proceedings of the 29th ACM International Conference on Architectural Support for Programming Languages and Operating Systems, Volume 3*, ASPLOS '24, page 932–949. ACM.
- Xupeng Miao, Gabriele Oliaro, Zhihao Zhang, Xinhao Cheng, Zeyu Wang, Zhengxin Zhang, Rae Ying Yee Wong, Alan Zhu, Lijie Yang, Xiaoxiang Shi, et al. 2023. Specinfer: Accelerating generative large language model serving with tree-based speculative inference and verification. *arXiv preprint arXiv:2305.09781*.
- Ramesh Nallapati, Bowen Zhou, Cícero Nogueira dos Santos, Çağlar Gülçehre, and Bing Xiang. 2016. [Abstractive text summarization using sequence-to-sequence rnns and beyond](#). In *Proceedings of the 20th SIGNLL Conference on Computational Natural Language Learning, CoNLL 2016, Berlin, Germany, August 11-12, 2016*, pages 280–290. ACL.
- Apoorv Saxena. 2023. [Prompt lookup decoding](#).
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. [Llama 2: Open foundation and fine-tuned chat models](#). *Preprint*, arXiv:2307.09288.
- Jason Wei, Maarten Bosma, Vincent Y. Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M. Dai, and Quoc V. Le. 2022. [Finetuned language models are zero-shot learners](#). *Preprint*, arXiv:2109.01652.
- Jerry Wei, Le Hou, Andrew Lampinen, Xiangning Chen, Da Huang, Yi Tay, Xinyun Chen, Yifeng Lu, Denny Zhou, Tengyu Ma, and Quoc V. Le. 2023. [Symbol tuning improves in-context learning in language models](#). *Preprint*, arXiv:2305.08298.
- Heming Xia, Zhe Yang, Qingxiu Dong, Peiyi Wang, Yongqi Li, Tao Ge, Tianyu Liu, Wenjie Li, and Zhi-fang Sui. 2024. [Unlocking efficiency in large language model inference: A comprehensive survey of speculative decoding](#). In *Findings of the Association for Computational Linguistics ACL 2024*, pages 7655–7671, Bangkok, Thailand and virtual meeting. Association for Computational Linguistics.
- Nan Yang, Tao Ge, Liang Wang, Binxing Jiao, Daxin Jiang, Linjun Yang, Rangan Majumder, and Furu Wei. 2023. [Inference with reference: Lossless acceleration of large language models](#). *Preprint*, arXiv:2304.04487.
- Hanling Yi, Feng Lin, Hongbin Li, Peiyang Ning, Xiaotian Yu, and Rong Xiao. 2024. Generation meets verification: Accelerating large language model inference with smart parallel auto-correct decoding. *arXiv preprint arXiv:2402.11809*.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric. P Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. 2023. [Judging llm-as-a-judge with mt-bench and chatbot arena](#). *Preprint*, arXiv:2306.05685.