Faster In-Context Learning for LLMs via N-Gram Trie Speculative Decoding

Anonymous ACL submission

Abstract

As a crucial method in prompt engineering, In-Context Learning (ICL) enhances the generalization and knowledge utilization capabilities of Large Language Models (LLMs) (Dong 004 et al., 2024). However, the lengthy retrieved contexts and limited token throughput in autoregressive models significantly constrain reasoning speed. To address this challenge, we propose N-Gram Trie Speculative Decoding, a novel approach that leverages the overlap between context and model output. This method constructs an n-gram trie from the context to generate drafts, accelerating token genera-014 tion for LLMs. We evaluate our approach on summarization, Retrieval-Augmented Genera-016 tion (RAG), and context-based Question Answering (QA) tasks. Experimental results on 017 018 Vicuna-7B, Llama2-7B-Chat, and Llama3-8B-Instruct demonstrate substantial speed improvements without compromising accuracy. Compared with various strong baselines, our method achieves the highest mean speedup, showcasing its effectiveness and efficiency.

1 Introduction

In-Context Learning (ICL) has emerged as a transformative paradigm in the field of prompt engineering, fundamentally reshaping how Large Language 027 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 retrievalaugmented 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. 041

042

043

044

045

047

049

052

053

055

059

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

075

076

077

078

081

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

100

101

105 106

108

109

110

111

112

113

114

115

116

117

118

119

121

122

123

124

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 pretraining 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. 130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

152

153

154

155

156

157

158

159

160

161

162

163

164

165

166

167

168

169

170

171

172

173

174

175

176

177

178

179

180

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. Nonetheless REST simply store all the corpus into suffix arrays offline. When inference starts, REST will search for the

suffix in the array corpus. Directly search from the 181 corpus has disadvantage because LLMs can't see 182 the given arrays during output stage. The generation is independent from the external corpus. Also, 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 187 *n*-gram token history as draft to verify. But it is useful only when the output token is repeatedly generated. PLD/LLMA (Yang et al., 2023) (Sax-190 ena, 2023) also try to use the overlap between the input and output, but they simply copy certain num-192 bers of suffixes without matching all the potential 193 suffixes in the prompt. Both of them don't fully 194 make use of the given prompt. 195

2.3 Tree Attention

196

197

198

199 200

201

202

203

207

210

211

212

213

214

215

216

217

218

219

222

229

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 datastore. 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 prefixbased hierarchy. Its key advantage is faster suffix

Algorithm 1 Trie Generation

Inp	ut: T: Token list							
D: Collected n-gram sample results								
	L_p : Maximum prefix length							
Ou	tput: $ au$							
1:	Init root as τ > 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 childinnode.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 $ au$							

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 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. 230

231

232

233

234

236

237

238

239

240

241

242

243

244

245

246

247

248

249

N-Gram Sampling Specifically, for the context token list $T = \{t_1, t_2, ..., t_l\}$, 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



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.

expressed as follows:

$$P_{i} = \{t_{i}, t_{i+1}, \dots, t_{i+L_{p}-1}\} \\ S_{i} = \{t_{i+L_{p}}, t_{i+L_{p}+1}, \dots, t_{i+n-1}\} \} i \in [1, l],$$

$$(1)$$

where i denotes the start index of the window. We establish the dependency between the prefix and suffix for each tokens group, and obtain the dependency set D by sliding window sampling. D can be defined in the following form:

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

where f is the frequency of dependency $< P_i, S_i >$ 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:

$$SP_{i} = \{SP_{i,j} | j \in (0, L_{p})\} \\ = \{P_{i}[j : L_{p}] | j \in (0, L_{p})\}, i \in [1, l],$$
(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:

276

277

278

279

281

283

284

287

293

294

295

296

297

299

300

301

303

304

$$node = \begin{cases} child, & if \quad t \in node.children \\ node_t, & if \quad t \notin node.children \end{cases},$$
(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 $T_s = \{t_1, t_2, ..., t_k\}$. 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

251

258

260

261

267

271

272

273



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 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 substract 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 gernerated tokens 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, refering 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:

$$d_{k} = < \operatorname{Path}(v_{r}, v_{k}), f_{k} > = < \{v_{r}, v_{1}, ..., v_{i}, ..., v_{k}\}, f_{k} >,$$
(5)
$$i \in [1, k], v_{i} \in \tau_{s},$$

where $Path(v_r, v_k)$ 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.

332 **Model Verification** Figure 2 shows an example of tree attention verifying the draft trie. For the 333 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 336

level. The specific form can be expressed as:

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

where p_i is the position id of t_i . 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 2dimention 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 taget 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 Spec-Bench(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-rerankerv2-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.

305

307

308

337

339

341

342

343

344

345

346

347

350

351

352

353

354

355

356

357

359

360

361

362

363

364

365

366

367

369

370

371

372

373

374

375

376

377

378

379

380

381

Model	Method	Spec-	Bench RAG	TriviaQA	Hagrid	Mean Speedup
	Vanilla		1.00 × (1.00)	1.00 × (1.00)	1.00 × (1.00)	1 00m
	vaiinia	1.00×(1.00)	$1.00 \times (1.00)$	$1.00 \times (1.00)$	$1.00 \times (1.00)$	1.00X
	SpS	$1.69 \times (2.44)$	$1.59 \times (2.30)$	$1.74 \times (2.49)$	$1.40 \times (2.46)$	1.61x
	Medusa	$1.48 \times (2.01)$	$1.45 \times (2.08)$	$1.45 \times (2.03)$	$1.56 \times (2.17)$	1.49x
Vicupa 7B	SPACE	$1.69 \times (2.26)$	$1.47 \times (1.91)$	$1.57 \times (2.26)$	$1.26 \times (2.05)$	1.50x
viculia-7D	Hydra	1.86 × (2.70)	<u>1.88×</u> (<u>2.90</u>)	$1.86 \times (2.84)$	1.52×(<u>2.98</u>)	$1.78 \times$
	Lookahead	$1.29 \times (1.54)$	$1.19 \times (1.48)$	$1.27 \times (1.46)$	$0.95 \times (1.47)$	1.18x
	REST	$1.13 \times (1.65)$	$1.32 \times (1.89)$	$1.31 \times (1.71)$	$1.33 \times (1.82)$	1.27x
	Ours	$1.75 \times (2.39)$	3.48 × (5.19)	1.92 × (2.36)	1.94× (3.07)	2.27×
	Vanilla	1.00×(1.00)	$1.00 \times (1.00)$	1.00×(1.00)	1.00×(1.00)	1.00x
	SpS	$1.25 \times (1.54)$	$1.47 \times (1.91)$	1.35×(<u>1.86</u>)	$1.38 \times (1.62)$	1.36x
Llama2-7B-Chat	Lookahead	1.44 × (<u>1.59</u>)	$1.40 \times (1.63)$	$1.53 \times (1.71)$	1.38×(<u>1.97</u>)	$1.44 \times$
	REST	$1.03 \times (1.54)$	1.14×(<u>1.91</u>)	$1.22 \times (1.68)$	$1.42 \times (1.47)$	1.20x
	Ours	<u>1.28×</u> (1.76)	3.62 × (5.00)	$\textbf{1.89} \times \textbf{(2.88)}$	$1.61 \times (2.34)$	2.10 ×
	Vanilla	1.00×(1.00)	$1.00 \times (1.00)$	1.00×(1.00)	1.00×(1.00)	1.00x
Llomo 2 9D Instant	Lookahead	1.25× (1.60)	<u>1.18×</u> (1.51)	$1.58 \times (1.54)$	$1.38 \times (1.73)$	1.50 imes
Liama5-8B-Instruct	REST	$0.93 \times (1.54)$	$1.14 \times (1.91)$	$1.13 \times (1.61)$	$1.02 \times (1.69)$	1.05x
	Ours	$1.06 \times (1.42)$	$1.77 \times (2.11)$	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 <u>underline</u>.

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).

4.1.4 Hyperparameters

In the experiment, there are two hypermeter that need to be tuned: matched prefix L_p and gramlength n. So we conduct the experiment to test the efficiency. The details can be seen in table2. 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-v2m3 (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* 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). 411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

432

433

434

435

436

437

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 $(2.27 \times \text{ on Vicuna-7B}, 2.10 \times \text{ on Llama2-7B-Chat} \text{ and } 1.56 \times \text{ on Llama3-8B-Instruct})$. 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.

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 $3.48 \times$, $3.62 \times$ and $1.77 \times$. 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

400

401

402

403

404

405

406

407

408

409

410

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 disc ourage would-be a ttackers , but the use of two identical vans has raised e y eb rows . The vans were seen driving separately to Clinton's a ppointed location before leaving together in a seven-car motorcade. The Secret Service has decl ined to comment on the sec urity arrangements for dignitaries. The use of two Scooby vans has been observed in other instances , including when Presi dent 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 hel icopter, which usually travel s with two deco ys .

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.

438 second place on all models.

439

440

441

442

443 444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472 473

474

475

476

477

On TriviaQA dataset, the speedup rate of our 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.

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 datset, which outperforms other approaches by $0.19x-0.99 \times$. Compared to the basic method REST, we have more speedup on all models. This fully demonstrates the advantages brought by n-gram trie.

Summary Task. The performance of our method on Spec-Bench Summary dataset is not the best. Global draft-getting method will result in wrongly draft clipping and unnecessarily draft-searching. In summarization tasks, input articles can be segmented into discrete text blocks, with most generated outputs demonstrating primary dependency on individual text units. The proposed methodology employs a global draft selection mechanism that may inadvertently incorporate non-essential drafts, potentially introducing redundant verification overhead. However, in RAG and multi-document QA scenarios, generated content exhibits stronger reliance on comprehensive document analysis, necessitating preservation and rigorous evaluation of multiple drafts. The experimental validation confirms that our method optimizes context utilization while maintaining computational efficiency through adaptive draft management. But our method still ranks second $(1.75 \times \text{ on})$ Vicuna-7B, $1.28 \times$ on Llama2-7B-Chat and $1.06 \times$ on Llama3-8B-Instruct) in terms of speedup and outperform REST.

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

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. 478

479

480

481

482

483

484

485

486

487

4.4 Hyperparameter Analysis

In this section, we will will conduct experiments 488 on the hyperparameters n and L_p in our method 489 to get the best hyperparameter configuration. We 490 use the RAG task of the Spec-bench (Xia et al., 491 2024) on Llama2-7B-Chat to test the performance 492 of N-Gram-Trie. We try the value between 8-16 493 for n-gram length n and 2-5 for maximum prefix 494 length L_p . The performance of the N-Gram-Trie 495 with different hyperparameter is shown in the Table 496 2. We can find that when n is small, the speedup ef-497 fect will gradually deteriorate with the increase of 498 L_p . We think this is because the excessively long 499 L_p limits the length of the suffix, which in turn 500 reduces the acceleration ability. When n is large, 501 the token generation speed first speeds up and then slows down as L_p increases. This is mainly be-503



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



Figure 5: Time distribution in searching operation. The red opponent indicates time for searching in tree, the blue opponent means suffix processing and the cyan opponent means the attention tree making.

cause 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 to 3.

4.5 Time Analysis

504

506

508

509

510

511

512

513The time expended on Trie search versus Tree con-514struction during each step of draft retrieval from the515Trie is illustrated in the figure 5. It can be seen that516draft processing and draft tree making cost much517more time than Trie searching.

4.6 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. 518

519

520

521

522

523

524

525

526

527

528

529

530

5 Conclusion

In this paper, we propose N-Gram Trie Spec-531 ulative Decoding, a novel approach to acceler-532 ate in-context inference for large language mod-533 els. By constructing an n-gram trie from the con-534 text through prefix and suffix dependencies, our 535 method efficiently generates speculative decoding drafts, leveraging the overlap between context and 537 model output. Extensive experiments on sum-538 marization, RAG, and context QA tasks demon-539 strate significant speedups—2.27x on Vicuna-7B, 540 2.10x on Llama2-7B-Chat, and 1.56x on Llama3-541 8B-Instruct—without compromising output quality. 542 This work addresses a critical limitation of ICL, 543 providing an effective and scalable solution for 544 real-world LLM deployment. 545

6 Limitations

546

564

567

569

570

571

572

574

576

577

578

579

580

582

586

588

589

591

595

This approach presents several limitations. First, 547 while the trie-based generation and search mech-548 anism offers efficiency advantages, its current im-549 plementation has suboptimal aspects. A key issue arises when multiple suffix candidates share 551 identical frequency scores, which may lead to the 552 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 re-558 trieval 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. 563

References

- 565 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 Ilm 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 596 Johnson, Gergely Szilvasy, Pierre-Emmanuel Mazaré, 597 Maria Lomeli, Lucas Hosseini, and Hervé Jégou. 2025. The faiss library. *Preprint*, arXiv:2401.08281. 599 Yichao Fu, Peter Bailis, Ion Stoica, and Hao Zhang. 600 2023. Breaking the sequential dependency of llm 601 inference using lookahead decoding. 602 Yichao Fu, Peter Bailis, Ion Stoica, and Hao Zhang. 603 2024. Break the sequential dependency of llm inference using lookahead decoding. arXiv preprint 605 arXiv:2402.02057. 606 Yuxian Gu, Li Dong, Furu Wei, and Minlie Huang. 607 2023. Pre-training to learn in context. Preprint, 608 arXiv:2305.09137. 609 Zhenyu He, Zexuan Zhong, Tianle Cai, Jason D Lee, 610 and Di He. 2023. Rest: Retrieval-based speculative 611 decoding. Preprint, arXiv:2311.08252. 612 Mandar Joshi, Eunsol Choi, Daniel Weld, and Luke 613 Zettlemoyer. 2017. triviaqa: A Large Scale Distantly 614 Supervised Challenge Dataset for Reading Compre-615 hension. arXiv e-prints, arXiv:1705.03551. 616 Ehsan Kamalloo, Aref Jafari, Xinyu Zhang, Nandan 617 618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

- Thakur, and Jimmy Lin. 2023. HAGRID: A humanllm collaborative dataset for generative informationseeking 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

651

- 662 663 664 665
- 666 667

66

6

671

- 672 673 674
- 675 676 677
- 678 679
- 6
- 6
- 684 685

686 687 688

69 69 69

694

695 696

6

6

700

7

70

707 708 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, Çaglar Gülçehre, and Bing Xiang. 2016.
 Abstractive text summarization using sequence-tosequence 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 finetuned 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 Zhifang 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. 709

710

711

712

713

714

715

716

718

719

720

721

722

723

724

725

726

727

728

729

730

731

732

733

734

735

737

739

740

741

742

743

744

745

746

- 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.

A More experiments about parameters

We conduct more experiments in summary and RAG tasks. The experiment results can be seen in table??:

There is a downgrade of the experiment result because the experiments are running with other processes using GPU. These results show that if the model and the dataset change, n and L_p also need to adjust for better speculative performance. What's more, the strategy of choosing the parameters may vary in GPU conditions. Also, we found that in Llama3, there is no significant degrade in the performance as the data changes. The choice of the n and L_p may need further discussion because now we haven't found the best way of selecting n and L_p . The best choices may be different due to many factors, which needs more experiments to check.

Llama2 Summarization					Llama2 RAG					
n/L_p	2	3	4	5		n/L_p	2	3	4	4
8	17.4	14.7	17.7	17.5		8	27.5	41.1	37.6	24
9	13.7	18.1	16.6	18.0		9	45.2	42.6	31.2	35
10	18.3	15.5	17.9	14.7		10	45.1	43.7	36.6	39
11	17.8	18.2	16.0	16.8		11	50.9	43.8	44.0	36
12	17.7	17.9	18.6	18.4		12	52.3	49.7	45.5	40
13	14.5	17.8	18.2	18.0		13	38.7	38.7	50.1	34
14	15.1	18.1	17.8	17.8		14	54.3	48.2	51.0	49
15	18.0	17.7	13.5	18.1		15	55.3	38.9	48.3	39

Table 3: Llama2 Performance

Vicuna Summarization						Vicuna RAG					
n/L_p	2	3	4	5		n/L_p	2	3	4	5	
8	18.7	23.1	22.5	18.8		8	35.7	42.8	29.5	25.2	
9	25.7	22.1	23.1	17.7		9	50.3	47.6	33.2	29.9	
10	25.5	25.9	25.6	23.1		10	38.2	50.5	36.1	44.0	
11	24.3	25.6	25.8	23.8		11	53.2	48.9	37.9	35.9	
12	25.1	25.1	26.0	25.9		12	54.8	52.1	50.8	50.8	
13	24.7	24.5	24.5	25.5		13	55.1	54.0	52.8	50.9	
14	18.5	24.5	25.1	24.9		14	55.3	41.8	54.5	53.8	
15	25.1	25.1	23.9	24.2		15	59.6	54.1	41.2	54.6	

Table 4: Vicuna Performance

Llama3 Summarization					Llama3 RAG						
n/L_p	2	3	4	5		n/L_p	2	3	4	5	
8	13.6	13.6	13.9	13.7		8	20.9	21.0	19.6	19.2	
9	13.7	13.8	13.6	13.6		9	20.6	20.7	19.8	19.8	
10	13.6	13.5	13.5	13.5		10	20.0	20.2	20.6	19.9	
11	13.7	13.5	13.4	13.5		11	20.0	20.1	20.3	20.5	
12	13.5	13.5	13.5	13.5		12	19.9	20.0	20.5	20.3	
13	13.3	13.5	13.4	13.7		13	20.4	20.5	20.0	19.1	
14	13.3	13.2	13.5	13.5		14	20.1	20.1	20.3	20.2	
15	13.6	13.3	13.4	13.4		15	20.2	20.0	20.0	20.6	

Table 5: Llama3 Performance