

Towards Fast Multilingual LLM Inference: Speculative Decoding and Specialized Drafters

Anonymous ACL submission

Abstract

Large language models (LLMs) have revolutionized natural language processing and broadened their applicability across diverse commercial applications. However, the deployment of these models is constrained by high inference time in multilingual settings. To mitigate this challenge, this paper explores a training recipe of an assistant model in speculative decoding, which are leveraged to draft and then its future tokens are verified by the target LLM. We show that language-specific draft models, optimized through a targeted *pretrain-and-finetune* strategy, substantially brings a speedup of inference time compared to the previous methods. We validate these models across various languages in inference time, out-of-domain speedup, and GPT-4o evaluation.

1 Introduction

Large language models (LLMs) such as Gemini (Team et al., 2023), GPT (Achiam et al., 2023), and Llama (Touvron et al., 2023a) have remarkable success across various natural language processing tasks. Their deployment in commercial settings has expanded to include applications such as coding assistance, writing support, conversational interfaces, and tools for search (Reid et al., 2024). Despite their potential, the practical deployment of these models is often limited by prohibitively high inference time, particularly in multilingual contexts (Ahia et al., 2023). For example, character-level and byte-level models exhibit encoding length discrepancies exceeding fourfold for certain language pairs, resulting in significant disparities in cost and inference time available to different language communities (Petrov et al., 2024). These challenges present substantial huddles to scalable and cost-efficient applications of LLMs.

Speculative decoding, utilizing assistant models, has emerged as a promising strategy to improve LLM inference efficiency (Leviathan et al.,

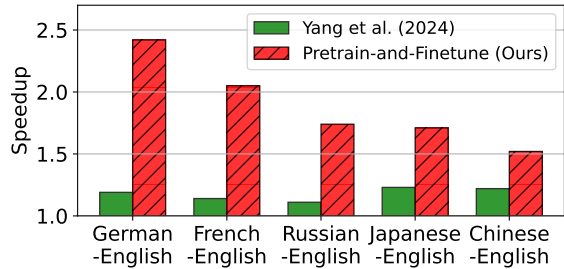


Figure 1: Speedup ratio¹ relative to the standard autoregressive greedy decoding on various multilingual datasets. Target model is Vicuna 7B v1.3 and the drafter is Vicuna 68M. Speculative greedy sampling is implemented with the drafter by Yang et al. (2024) (green) and our specialized drafter (*pretrain-and-finetune*) (red).

2023; Chen et al., 2023; Xia et al., 2024), inspired by speculative execution (Burton, 1985). This method drafts potential future tokens by leveraging a smaller model for the initial predictions. In parallel, these tokens are verified by the target LLM, ensuring only outputs aligned with the target LLM’s predictions are accepted. Recent efforts are focused on aligning these initial predictions with the target LLM’s outputs (Liu et al., 2023; Zhou et al., 2023). This involves advancing the training methods and modifying the architectural design of drafters (Miao et al., 2024; Li et al., 2024).

Although speculative decoding has garnered considerable hype recently, the adaptation of this approach to multilingual scenarios common in real-world applications remains largely unexplored. Prevailing methods (Cai et al., 2024; Li et al., 2024; Yang et al., 2024) use small drafters simply trained on datasets such as ShareGPT (ShareGPT, 2023) which is often used for instruction tuning of LLMs to learn a pattern of target LLM’s language modeling. However, our investigations reveal that such approach are insufficient for multilingual translation (Figure 1). This research also raises concerns

¹Evaluated on a single RTX3090 GPU with a batch size 1.

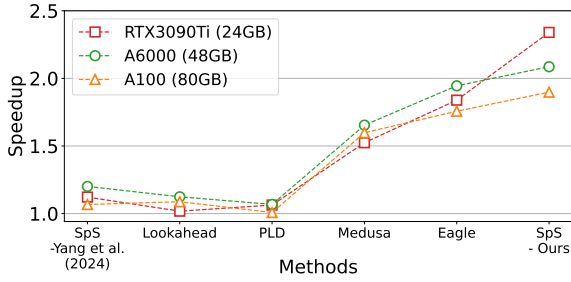


Figure 2: Speedup comparison of various speculative decoding methods on WMT16 De-En dataset (Bojar et al., 2016) with greedy settings ($T = 0.0$) across various hardware. Target model is Vicuna-7B.

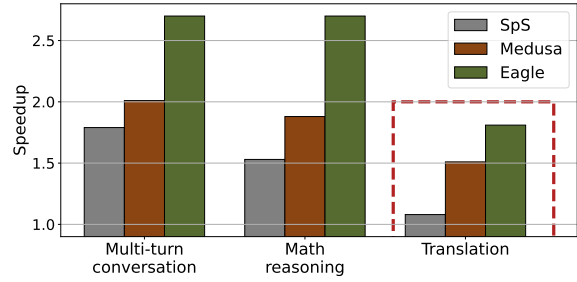


Figure 3: Speedup² comparison across categories containing multi-turn conversation (MT-Bench) (Zheng et al., 2024), math reasoning (GSM8K) (Cobbe et al., 2021), and translation (WMT16 De-En). Target model is Vicuna-7B with speculative greedy sampling.

066 regarding the capacity of such small drafters with
 067 simple tuning to comprehend the nuances of all
 068 target languages, thus questioning the feasibility of
 069 employing such models for universal speculative
 070 decoding. This paper aims to shed light upon the
 071 behaviors of drafters in speculative decoding within
 072 multilingual tasks and to explore their efficacy. Our
 073 contributions are as follows:

- 074 • We demonstrate that the strategy of *pretrain-*
 075 *and-finetune* significantly improves the alignment
 076 of drafter models, achieving the highest
 077 speedup ratio among the baselines (Figure 2).
- 078 • We find that the speedup ratio increases as the
 079 number of tokens specific to the target task
 080 used in training increases. This speedup is
 081 logarithmically proportional to the scale of
 082 token count in drafter training.
- 083 • In multilingual translation, we observe that
 084 input languages consistent with the training
 085 set result in notable speedup, whereas outputs
 086 aligned with the training domain do not necessarily
 087 lead to improved performance. Additionally, our
 088 results are corroborated by GPT-4o judgment scores
 089 and qualitative analyses.

090 2 Method

091 2.1 Preliminaries: speculative decoding

092 Speculative decoding employs a draft-verify-accept
 093 paradigm for fast inference. This method leverages
 094 a simpler assistant model (M_p) to predict easy tokens,
 095 thereby addressing memory bandwidth constraints in
 096 LLM inference (Shazeer, 2019):

- 097 1. **Draft:** An assistant model M_p , which is less
 098 computationally intensive than the target LLM
 099 M_q , drafts the future tokens $\{x_{t_1}, \dots, x_{t_K}\}$
 100 based on the input sequence x_1, \dots, x_t .

- 101 2. **Verify:** The target LLM M_q assesses
 102 each token x_{t_i} regarding whether it is
 103 aligned with its own predictions: $p_i =$
 104 $M_p(x_{t_i}|x_1, \dots, x_t, x_{t_1}, \dots, x_{t_{i-1}})$, $q_i =$
 105 $M_q(x_{t_i}|x_1, \dots, x_t, x_{t_1}, \dots, x_{t_{i-1}})$.

- 106 3. **Accept:** Tokens meeting the validation criteria
 107 (e.g., rejection sampling) aligned with
 108 M_q 's outputs are retained. Tokens failing verification
 109 are either discarded or corrected, and
 110 the draft-verify cycle is repeated.

111 In this paper, the verification process employs
 112 rejection sampling (Leviathan et al., 2023; Li et al.,
 113 2024) when the temperature parameter is above
 114 zero to accept only tokens that closely match M_q 's
 115 predictions. For greedy decoding with a temperature
 116 of zero, tokens are accepted if they are identical
 117 to M_q 's predictions.

118 2.2 Motivation

119 Our evaluation of various speculative models, including
 120 SpS (Chen et al., 2023), Medusa (Cai et al., 2024),
 121 Eagle (Li et al., 2024), as shown in Figure 3,
 122 demonstrates that speedup ratios significantly differ
 123 by task domain. While these models excel in English
 124 tasks such as multi-turn conversations and mathematical
 125 reasoning, where they achieve notable speed improvements,
 126 they underperform in translation tasks (red dotted box
 127 in Figure 3). This result confirms that the effectiveness
 128 of these models is not universal but may be highly
 129 language-specific. The consistent underperformance in
 130 translation tasks highlights a key weakness and drives
 131 our study towards developing specialized drafters.
 132

²Evaluated on a single RTX3090 GPU with a batch size 1.

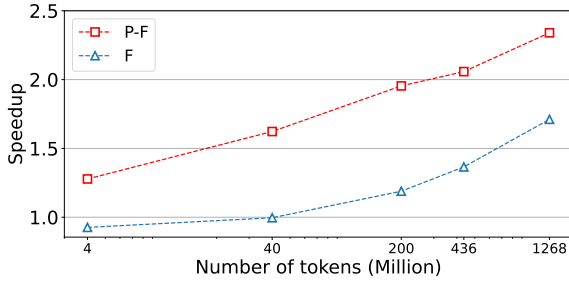


Figure 4: Speedup with speculative greedy sampling on the WMT16 De-En dataset as the training token for finetune (F) count varies, displayed on a logarithmic x-axis. ‘P-F’ represents our strategy and ‘F’ involves training solely on De-En without pretrain step (P).

2.3 Training specialized assistant models

At the core of our approach is the recognition that smaller models, due to their inherent limited capacity, struggle to capture the diverse token distributions across languages. To address this challenge, we present specialized drafter models tailored to each language. Our strategy consists of:

1. **Pretrain (P):** Assistant models are pretrained on a part of C4 (Raffel et al., 2019) and ShareGPT dataset (ShareGPT, 2023) for language modeling.
2. **Finetune (F):** The models are finetuned on the target lingual task with instructions to refine their responses to non-English inputs.

While the practices of pretraining and finetuning are well-established paradigms in language model training, applying these steps to drafter models represents a novel adaptation within the field. Traditionally, assistant models have been trained from scratch with little strategic rationale or clear justification for dataset selection.

Figure 4 shows that the *pretrain-and-finetune* strategy significantly boosts the speedup ratio as the number of training tokens increases. Our ‘P-F’ approach outperforms models that are only finetuned (F), and even surpasses the speedup rates by Yang et al. (2024), which stood at 1.12.

Dataset with self-distillation The training dataset for our assistant models is generated through self-distillation from the target LLM, ensuring alignment with its outputs (Kim and Rush, 2016; Zhou et al., 2023; Cai et al., 2024). To capture the full range of the target’s output variability, we generate multiple responses at a range of temperatures—{0.0, 0.3, 0.7, 1.0}.

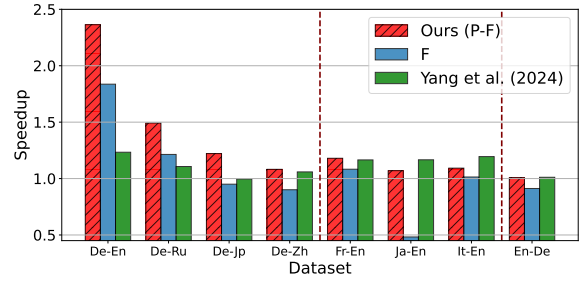


Figure 5: Speedup with speculative greedy sampling on various out-of-domain dataset as the drafters for ‘Ours (P-F)’ and ‘F’ are trained on WMT16 De-En dataset.

3 Experiment

3.1 Experimental setup

Models We utilize Vicuna 7B (Chiang et al., 2023), Gemma-Instruct 7B (Team et al., 2024), and Llama2-chat (Touvron et al., 2023b) as target LLMs. The drafter models employed include Vicuna 68M (Yang et al., 2024), a custom Gemma 250M drafter and Llama 68M (Miao et al., 2024). Training configurations are outlined in Appendix F.

Number of drafts For speculative sampling (SpS), we use a single draft candidate (Chen et al., 2023). In contrast, Medusa and Eagle models are evaluated using multiple drafts via tree-attention mechanism by following their original settings.

Training and evaluation Training datasets for each target model are generated via self-distillation and comprise five datasets: German (De)→English (En), French (Fr)→En, Russian (Ru)→En, Japanese (Ja)→En and Chinese (Zh)→En, each with 4 million (M) conversations (~ 1.3 billion (B) tokens) sourced from WMT14 Fr-En (Bojar et al., 2014), WMT16 De-En, and Ru-En (Bojar et al., 2016), and JParaCrawl-v3.0 (Morishita et al., 2022). Evaluations are conducted using a single NVIDIA 3090 GPU, under both greedy settings ($T=0.0$) and for diversity at $T=1.0$ with three different seeds. The details are in Appendix F.

3.2 Main result

Overall Table 1 shows that our specialized drafter (*pretrain-and-finetune*) for targeted languages demonstrates superior performance across multiple translation tasks, recording the highest speedup in both deterministic ($T=0.0$) and diverse ($T=1.0$) settings. At $T=0.0$, our model outperforms all competitors with an average speedup ratio of

Table 1: Speedup comparison of different methods for Vicuna 7B v1.3. Results are provided for $T=0.0$ and $T=1.0$ across various translation tasks. For our approach, each drafter is finetuned with the corresponding dataset.

Method	T=0.0					Avg	T=1.0					Avg
	De→En	Fr→En	Ru→En	Ja→En	Zh→En		De→En	Fr→En	Ru→En	Ja→En	Zh→En	
Sps - Yang et al. (2024)	1.19±0.06	1.14±0.05	1.11±0.04	1.23±0.03	1.22±0.00	1.18±0.04	1.07±0.03	1.06±0.02	1.04±0.01	1.15±0.02	1.11±0.02	1.09±0.02
Lookahead (Fu et al., 2024)	1.03±0.01	1.01±0.02	0.98±0.01	1.00±0.01	0.96±0.00	1.00±0.01	1.03±0.03	1.04±0.03	0.99±0.00	0.98±0.05	0.98±0.00	1.01±0.02
PLD (Saxena, 2023)	1.13±0.06	1.05±0.04	1.03±0.00	1.09±0.05	0.99±0.07	1.06±0.05	-	-	-	-	-	-
Medusa (Cai et al., 2024)	1.58±0.05	1.57±0.01	1.52±0.01	1.55±0.01	1.43±0.00	1.53±0.02	1.61±0.03	1.69±0.01	1.62±0.00	1.72±0.01	1.60±0.01	1.65±0.01
Eagle (Li et al., 2024)	1.90±0.05	1.88±0.00	1.67±0.05	1.88±0.01	1.75±0.01	1.81±0.02	1.57±0.00	1.61±0.01	1.45±0.02	1.63±0.01	1.51±0.03	1.55±0.01
Sps - <i>pretrain-and-finetune</i> (Ours)	2.42±0.02	2.05±0.04	1.74±0.02	1.71±0.01	1.52±0.01	1.89±0.02	1.99±0.01	1.86±0.03	1.58±0.00	1.67±0.01	1.44±0.00	1.71±0.01

Table 2: Examples of speculative decoding on WMT16 De-En dataset. Black indicates standard decoded output and magenta indicates accepted draft tokens.

Input
Translate German to English: So wie er gestartet ist , wird es nicht lange dauern , bis er auf der „ Pferd des Jahres “ -Schau ist – und ich bin mir sicher , dass er gut abschneiden wird.
SpS with a drafter by Yang et al. (2024)
As he started, it won't take long until he's on the "Horse of the Year" show, and I'm sure he'll do well.
Eagle (Li et al., 2024)
As he started, it won't take long until he's on the "Horse of the Year" show, and I'm sure he'll do well.
SpS with our specialized drafter (<i>pretrain-and-finetune</i>)
As he started, it won't take long until he's on the "Horse of the Year" show, and I'm sure he'll do well.

1.89. Similarly, at $T=1.0$, it maintains robust performance with an overall speedup ratio of 1.71.

Speedup on out-of-domain translation tasks

As Figure 5 shows, our analysis reveals variability applying the drafter, trained on the WMT16 De-En dataset, across diverse translation pairs. Speedups are consistently higher when translating from German to other languages, highlighting the importance of input domain consistency with the training data. Conversely, translations involving non-German languages with English and English-German pairings show limited gains. This result emphasizes that effective speculation relies more on matching the translation task's input domain with the training data than on the output domain.

Qualitative analysis on responses Table 2 provides examples of speculative inference on the WMT16 De-En dataset. Both Eagle and our method incorporate a significant number of accepted tokens from drafts. However, our model achieves this with $\sim 75\%$ fewer parameters, leading to reduced latency and faster inference time (Table 1). Speculation typically takes place at critical junctions of the sentence such as transitions between clauses and phrases.

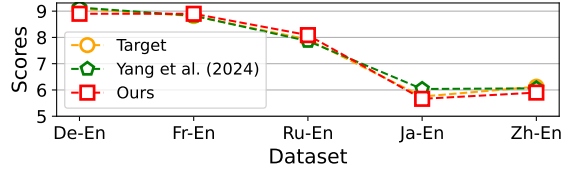


Figure 6: GPT-4o judgment scores following the Zheng et al. (2024) on various multilingual translation dataset. The score is evaluated random sampling with $T=1.0$.

Table 3: Ablations with speedup as the training stages continue on WMT19 Zh→En.

Target LLM - Drafter	P	+ F
Gemma-Instruct 7B - Gemma 250M	0.92±0.01	1.04±0.02
Llama2-chat 7B - Llama 68M	1.47±0.00	1.95±0.01

GPT-4o judgment analysis Figure 6 depicts the GPT-4o judgment scores (Zheng et al., 2024) generated using a temperature of 1.0. Our drafter closely matches the target Vicuna LLM across multiple datasets. The setup and further results are in Appendix F and Appendix G.

Ablation study Table 3 presents the ablation results, illustrating the progressive impact of the *pretrain-and-finetune* approach on the performance of Gemma and Llama2-chat models.

4 Conclusion

This paper has demonstrated that the *pretrain-and-finetune* strategy for training drafters significantly enhances speedup ratio relative to standard autoregressive decoding in multilingual translation tasks. This gain grows logarithmically with the increase in the number of training tokens. Supported by qualitative analysis, out-of-domain analysis, and GPT-4o evaluation, this strategy substantially outperforms the state-of-the-art methods in various language pairs. Our study uncovers approaches to maximize the benefits from drafter models, thereby setting a new benchmark in this area.

251 Limitations

252 Despite the improvement, our approach, requir-
253 ing separate drafters for each language, introduces
254 complexities in deployment, especially in multilin-
255 gual settings. For instance, in environments where
256 multiple languages are frequently interchanged,
257 such as multinational corporations or global cus-
258 tomer service platforms, the lack of an automated
259 drafter selection system could hinder operational
260 efficiency. Currently, our study focuses on inde-
261 pendent drafters; however, examining systems that
262 utilize interdependent models, similar to Eagle and
263 Medusa, might offer insights into more interest-
264 ing strategies. Additionally, while our findings
265 are promising for translation tasks, expanding this
266 methodology to other multilingual applications,
267 like real-time multilingual generation or summa-
268 rization, is essential to understand its broader ap-
269 plicability and uncover additional constraints.

270 This work primarily presents no direct ethical
271 concerns. Further discussions are detailed in Ap-
272 pendix B and Appendix H.

273 References

274 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama
275 Ahmad, Ilge Akkaya, Florencia Leoni Aleman,
276 Diogo Almeida, Janko Altenschmidt, Sam Altman,
277 Shyamal Anadkat, et al. 2023. Gpt-4 technical report.
278 *arXiv preprint arXiv:2303.08774*.

279 Orevaoghene Ahia, Sachin Kumar, Hila Gonen, Jungo
280 Kasai, David R Mortensen, Noah A Smith, and Yulia
281 Tsvetkov. 2023. Do all languages cost the same? tok-
282 enization in the era of commercial language models.
283 *arXiv preprint arXiv:2305.13707*.

284 Sangmin Bae, Jongwoo Ko, Hwanjun Song, and Se-
285 Young Yun. 2023. Fast and robust early-exiting
286 framework for autoregressive language models with
287 synchronized parallel decoding. *arXiv preprint*
288 *arXiv:2310.05424*.

289 Nikhil Bhendawade, Irina Belousova, Qichen Fu, Henry
290 Mason, Mohammad Rastegari, and Mahyar Na-
291 jibi. 2024. Speculative streaming: Fast llm infer-
292 ence without auxiliary models. *arXiv preprint*
293 *arXiv:2402.11131*.

294 Ondřej Bojar, Christian Buck, Christian Federmann,
295 Barry Haddow, Philipp Koehn, Johannes Leveling,
296 Christof Monz, Pavel Pecina, Matt Post, Herve Saint-
297 Amand, et al. 2014. Findings of the 2014 workshop
298 on statistical machine translation. In *Proceedings of*
299 *the ninth workshop on statistical machine translation*,
300 pages 12–58.

Ondrej Bojar, Rajen Chatterjee, Christian Federmann,
Yvette Graham, Barry Haddow, Matthias Huck, An-
tonio Jimeno Yepes, Philipp Koehn, Varvara Lo-
gacheva, Christof Monz, et al. 2016. Findings of
the 2016 conference on machine translation (wmt16).
In *First conference on machine translation*, pages
131–198. Association for Computational Linguistics.

F Warren Burton. 1985. Speculative computation, par-
allelism, and functional programming. *IEEE Trans-*
actions on Computers, 100(12):1190–1193.

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. *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. *arXiv preprint*
arXiv:2302.01318.

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.

Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian,
Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias
Plappert, Jerry Tworek, Jacob Hilton, Reiichiro
Nakano, et al. 2021. Training verifiers to solve math
word problems. *arXiv preprint arXiv:2110.14168*.

Mostafa Elhoushi, Akshat Shrivastava, Diana Liskovich,
Basil Hosmer, Bram Wasti, Liangzhen Lai, Anas
Mahmoud, Bilge Acun, Saurabh Agarwal, Ahmed
Roman, et al. 2024. Layer skip: Enabling early
exit inference and self-speculative decoding. *arXiv*
preprint arXiv:2404.16710.

Yimin Fan, Yaobo Liang, Alexandre Muzio, Hany Has-
san, Houqiang Li, Ming Zhou, and Nan Duan. 2021.
Discovering representation sprachbund for multilin-
gual pre-training. *arXiv preprint arXiv:2109.00271*.

Yichao Fu, Peter Bailis, Ion Stoica, and Hao Zhang.
2024. Break the sequential dependency of llm infer-
ence using lookahead decoding. *arXiv preprint*
arXiv:2402.02057.

Fabian Gloeckle, Badr Youbi Idrissi, Baptiste Rozière,
David Lopez-Paz, and Gabriel Synnaeve. 2024. Bet-
ter & faster large language models via multi-token
prediction. *arXiv preprint arXiv:2404.19737*.

Yuxian Gu, Li Dong, Furu Wei, and Minlie Huang. 2023.
Minillm: Knowledge distillation of large language
models. In *The Twelfth International Conference on*
Learning Representations.

Dan Hendrycks and Kevin Gimpel. 2016. Gaus-
sian error linear units (gelus). *arXiv preprint*
arXiv:1606.08415.

357	Taehyeon Kim, Ananda Theertha Suresh, Kishore Papineni, Michael Riley, Sanjiv Kumar, and Adrian Benton. 2024. Exploring and improving drafts in blockwise parallel decoding . <i>Preprint</i> , arXiv:2404.09221.	412
358		413
359		414
360		415
361	Yoon Kim and Alexander M Rush. 2016. Sequence-level knowledge distillation. <i>arXiv preprint arXiv:1606.07947</i> .	416
362		417
363		418
364	Jongwoo Ko, Sungnyun Kim, Tianyi Chen, and Se-Young Yun. 2024. Distillm: Towards streamlined distillation for large language models. <i>arXiv preprint arXiv:2402.03898</i> .	419
365		420
366		421
367		422
368	Yaniv Leviathan, Matan Kalman, and Yossi Matias. 2023. Fast inference from transformers via speculative decoding. In <i>International Conference on Machine Learning</i> , pages 19274–19286. PMLR.	423
369		424
370		425
371		426
372	Yuhui Li, Fangyun Wei, Chao Zhang, and Hongyang Zhang. 2024. Eagle: Speculative sampling requires rethinking feature uncertainty. <i>arXiv preprint arXiv:2401.15077</i> .	427
373		428
374		429
375		430
376	Zhuohan Li, Eric Wallace, Sheng Shen, Kevin Lin, Kurt Keutzer, Dan Klein, and Joey Gonzalez. 2020. Train big, then compress: Rethinking model size for efficient training and inference of transformers. In <i>International Conference on machine learning</i> , pages 5958–5968. PMLR.	431
377		432
378		433
379		434
380		435
381		436
382	Xiaoxuan Liu, Lanxiang Hu, Peter Bailis, Ion Stoica, Zhijie Deng, Alvin Cheung, and Hao Zhang. 2023. Online speculative decoding. <i>arXiv preprint arXiv:2310.07177</i> .	437
383		438
384		439
385		440
386	Ilya Loshchilov and Frank Hutter. 2017. Decoupled weight decay regularization. <i>arXiv preprint arXiv:1711.05101</i> .	441
387		442
388		443
389	Xupeng Miao, Gabriele Oliaro, Zhihao Zhang, Xinhao Cheng, Zeyu Wang, Zhengxin Zhang, Rae Ying Yee Wong, Alan Zhu, Lijie Yang, Xiaoxiang Shi, et al. 2024. Specinfer: Accelerating large language model serving with tree-based speculative inference and verification. In <i>Proceedings of the 29th ACM International Conference on Architectural Support for Programming Languages and Operating Systems, Volume 3</i> , pages 932–949.	444
390		445
391		446
392		447
393		448
394		449
395		450
396		451
397		452
398	Makoto Morishita, Katsuki Chousa, Jun Suzuki, and Masaaki Nagata. 2022. Jparacrawl v3. 0: A large-scale english-japanese parallel corpus. <i>arXiv preprint arXiv:2202.12607</i> .	453
399		454
400		455
401		456
402	OpenAI. 2024. Hello GPT-4o . Accessed: Insert the current date.	457
403		458
404	Jiayi Pan. 2023. Tiny-vicuna 1b. https://huggingface.co/Jiayi-Pan/Tiny-Vicuna-1B .	459
405		460
406	David A Patterson. 2004. Latency lags bandwidth. <i>Communications of the ACM</i> , 47(10):71–75.	461
407		462
408	Aleksandar Petrov, Emanuele La Malfa, Philip Torr, and Adel Bibi. 2024. Language model tokenizers introduce unfairness between languages. <i>Advances in Neural Information Processing Systems</i> , 36.	463
409		464
410		465
411		466
	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2019. Exploring the limits of transfer learning with a unified text-to-text transformer . <i>arXiv e-prints</i> .	467
	Machel Reid, Nikolay Savinov, Denis Teplyashin, Dmitry Lepikhin, Timothy Lillicrap, Jean-baptiste Alayrac, Radu Soricut, Angeliki Lazaridou, Orhan Firat, Julian Schrittwieser, et al. 2024. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. <i>arXiv preprint arXiv:2403.05530</i> .	
	Apoorv Saxena. 2023. Prompt lookup decoding .	
	Tal Schuster, Adam Fisch, Jai Gupta, Mostafa Dehghani, Dara Bahri, Vinh Tran, Yi Tay, and Donald Metzler. 2022. Confident adaptive language modeling. <i>Advances in Neural Information Processing Systems</i> , 35:17456–17472.	
	ShareGPT. 2023. Sharegpt: Vicuna unfiltered dataset. https://huggingface.co/datasets/Aeala/ShareGPT_Vicuna_unfiltered . Accessed: 2024.	
	Noam Shazeer. 2019. Fast transformer decoding: One write-head is all you need. <i>arXiv preprint arXiv:1911.02150</i> .	
	Mitchell Stern, Noam Shazeer, and Jakob Uszkoreit. 2018. Blockwise parallel decoding for deep autoregressive models. <i>Advances in Neural Information Processing Systems</i> , 31.	
	Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. 2023. Gemini: a family of highly capable multimodal models. <i>arXiv preprint arXiv:2312.11805</i> .	
	Gemma Team, Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya Pathak, Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, et al. 2024. Gemma: Open models based on gemini research and technology. <i>arXiv preprint arXiv:2403.08295</i> .	
	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023a. Llama: Open and efficient foundation language models. <i>arXiv preprint arXiv:2302.13971</i> .	
	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shrutu 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,	

468 Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Di-
469 ana Liskovich, Yinghai Lu, Yuning Mao, Xavier Mar-
470 tinet, Todor Mihaylov, Pushkar Mishra, Igor Moly-
471 bog, Yixin Nie, Andrew Poulton, Jeremy Reizen-
472 stein, Rashi Rungta, Kalyan Saladi, Alan Schelten,
473 Ruan Silva, Eric Michael Smith, Ranjan Subrama-
474 nian, Xiaoqing Ellen Tan, Binh Tang, Ross Tay-
475 lor, Adina Williams, Jian Xiang Kuan, Puxin Xu,
476 Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan,
477 Melanie Kambadur, Sharan Narang, Aurelien Ro-
478 driguez, Robert Stojnic, Sergey Edunov, and Thomas
479 Scialom. 2023b. [Llama 2: Open foundation and
480 fine-tuned chat models](#). *Preprint*, arXiv:2307.09288.

481 Neeraj Varshney, Agneet Chatterjee, Mihir Parmar, and
482 Chitta Baral. 2023. Accelerating llama inference by
483 enabling intermediate layer decoding via instruction
484 tuning with lite. *arXiv e-prints*, pages arXiv–2310.

485 Heming Xia, Zhe Yang, Qingxiu Dong, Peiyi Wang,
486 Yongqi Li, Tao Ge, Tianyu Liu, Wenjie Li, and
487 Zhifang Sui. 2024. Unlocking efficiency in large
488 language model inference: A comprehensive sur-
489 vey of speculative decoding. *arXiv preprint*
490 *arXiv:2401.07851*.

491 Guangxuan Xiao, Ji Lin, Mickael Seznec, Hao Wu,
492 Julien Demouth, and Song Han. 2023. Smoothquant:
493 Accurate and efficient post-training quantization for
494 large language models. In *International Conference
495 on Machine Learning*, pages 38087–38099. PMLR.

496 Nan Yang, Tao Ge, Liang Wang, Binxing Jiao, Daxin
497 Jiang, Linjun Yang, Rangan Majumder, and Furu
498 Wei. 2023. Inference with reference: Lossless ac-
499 celeration of large language models. *arXiv preprint*
500 *arXiv:2304.04487*.

501 Sen Yang, Shujian Huang, Xinyu Dai, and Jiajun Chen.
502 2024. Multi-candidate speculative decoding. *arXiv*
503 *preprint arXiv:2401.06706*.

504 Jun Zhang, Jue Wang, Huan Li, Lidan Shou, Ke Chen,
505 Gang Chen, and Sharad Mehrotra. 2023. Draft
506 & verify: Lossless large language model accelera-
507 tion via self-speculative decoding. *arXiv preprint*
508 *arXiv:2309.08168*.

509 Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan
510 Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin,
511 Zhuohan Li, Dacheng Li, Eric Xing, et al. 2024.
512 Judging llm-as-a-judge with mt-bench and chatbot
513 arena. *Advances in Neural Information Processing
514 Systems*, 36.

515 Yongchao Zhou, Kaifeng Lyu, Ankit Singh Rawat,
516 Aditya Krishna Menon, Afshin Rostamizadeh, San-
517 jiv Kumar, Jean-François Kagy, and Rishabh Agar-
518 wal. 2023. Distillspec: Improving speculative de-
519 coding via knowledge distillation. *arXiv preprint*
520 *arXiv:2310.08461*.

A Overview of appendix	521
This appendix provides supplementary material that expands on the main contents. Each section is designed to complement the research presented:	522
	523
	524
• Appendix B: Broader impact - Examines the wider implications of our findings on speculative decoding.	525
	526
	527
• Appendix C: Future work - Outlines possible directions for future research, building upon the current study’s findings to explore new avenues and improvements.	528
	529
	530
	531
• Appendix D: Related works - Provides a comprehensive review of literature and previous research that relate to the speculative decoding techniques discussed in the paper.	532
	533
	534
	535
• Appendix E: Algorithm - Details the algorithms used in the speculative decoding processes, providing pseudocode and explanations to support reproducibility.	536
	537
	538
	539
• Appendix F: Implementation details - Offers an in-depth look at the practical implementation of the speculative decoding methods, including baselines, self-distillation, training, and GPT-4o evaluation.	540
	541
	542
	543
	544
• Appendix G: Additional experimental results - Presents extra experimental data and analyses that were not included in the main sections due to space constraints.	545
	546
	547
	548
• Appendix H: Discussions - Engages in discussions on results, such as foundational beliefs that underpin our research approach, the number of drafts used, and drafter size.	549
	550
	551
	552
Each appendix is intended to provide clarity and additional context to the research.	553
	554
B Broader impact	555
Implementing language-specific drafters significantly enhances the speed of large language models tailored to diverse linguistic environments. For instance, a system could leverage heuristic analysis of input prompt token distributions to automatically select an optimal drafter, streamlining processing efficiency. Moreover, if a user interface allows individuals to choose their preferred language, the system can instantly apply the corresponding drafter, thereby accelerating response times considerably.	556
	557
	558
	559
	560
	561
	562
	563
	564
	565

566	Such advancements not only reduce computational	(FFN) heads directly into their architecture. These	615
567	load but also enrich the user experience by pro-	FFN heads facilitate the early drafting of token	616
568	viding rapid and culturally relevant responses in	sequences, enhancing throughput and reducing	617
569	multilingual contexts.	latency. Similarly, approaches such as the self-	618
570	C Future work	speculative model (Zhang et al., 2023) and El-	619
571	Future projects will explore broadening the scope	houshi et al. (2024) incorporate early exiting and	620
572	of our speculative decoding framework to cover	layer skipping strategies, allowing for a reduction	621
573	general multi-task environments, extending beyond	in computational load by prematurely terminating	622
574	multilingual translation to include varied domains	decoding processes or bypassing less impactful	623
575	such as legal and medical text processing. A sig-	neural layers. Another line of research explores the	624
576	nificant challenge lies in developing an efficient	blockwise parallel language models with multiple	625
577	method for selecting the appropriate drafter among	softmax heads pretrained from scratch presented by	626
578	multiple options when direct user input is unavail-	Stern et al. (2018) by either refining its drafts (Kim	627
579	able or when inputs consist of mixed languages.	et al., 2024) or scaling up the model size (Gloeckle	628
580	This issue becomes more complex as the ambigu-	et al., 2024).	629
581	ity of language indicators increases. To alleviate	D.2 Inference acceleration of LLM	630
582	this, designing an advanced router capable of intel-	As LLMs continue to evolve rapidly, enhancing	631
583	ligently assigning tasks to the most suitable drafter	their inference speed has become a focal area of	632
584	based on the nature of the input presents a promis-	research. Traditional techniques such as knowl-	633
585	ing direction. Training such a router involves lever-	edge distillation (Gu et al., 2023; Ko et al., 2024),	634
586	aging advanced techniques to understand and pre-	model compression (Li et al., 2020), and quantiza-	635
587	dict the optimal drafter based on contextual rep-	tion (Xiao et al., 2023) aim to optimize these mod-	636
588	resentations. This approach aims to improve the	els but often require extensive training adjustments	637
589	overall efficiency and accuracy of language model	or significant architectural modifications. More re-	638
590	applications across diverse and dynamically chang-	cent strategies have shifted towards applying early	639
591	ing content landscapes.	exiting mechanisms, particularly within series like	640
592	D Related works	T5 (Schuster et al., 2022; Bae et al., 2023) and	641
593	D.1 Speculative decoding	decoder-only architectures (Varshney et al., 2023),	642
594	Speculative decoding, advancing from blockwise	to streamline inference processes. Although early	643
595	parallel decoding introduced by Stern et al. (2018),	exiting can significantly hasten model responses by	644
596	adopts a draft-then-verify paradigm to enhance	truncating computational sequences, this method	645
597	LLM inference efficiency. This method addresses	typically introduces a trade-off with performance	646
598	latency issues in autoregressive decoding, which	degradation (Schuster et al., 2022).	647
599	stem from the extensive memory transfers required	E Algorithm: speculative sampling	648
600	for each token generation, leading to computational	By referring to Chen et al. (2023), Algorithm 1	649
601	underutilization (Xia et al., 2024; Patterson, 2004).	demonstrates the speculative sampling process. Ini-	650
602	To further advance this paradigm, Leviathan et al.	tiating with an initial prompt, an assistant model is	651
603	(2023) and Chen et al. (2023) introduced specu-	utilized to generate multiple prospective continua-	652
604	lative decoding and sampling, which includes the	tions at each step, which are concurrently verified	653
605	lossless acceleration of various sampling methods.	against the target LLM’s predictions.	654
606	These methods utilize smaller models from the	Each candidate token’s acceptance probability is	655
607	same series, such as T5-small, to accelerate infer-	calculated based on the target LLM’s relative con-	656
608	ence for larger counterparts like T5-XXL without	fidence compared to the assistant model’s sugges-	657
609	additional training.	tion (i.e., rejection sampling). If a value, randomly	658
610	Recent advancements in speculative decoding,	drawn from a uniform distribution, falls below this	659
611	exemplified by models like EAGLE (Li et al.,	threshold, the token is accepted and incorporated	660
612	2024) and Medusa (Cai et al., 2024), have sig-	into the ongoing sequence. If not, the algorithm	661
613	nificantly refined the efficiency of LLMs by in-	recalibrates, adjusting the speculative path by di-	662
614	tegrating lightweight feedforward neural network	rectly sampling from the differences in predictions,	663

Algorithm 1: Speculative sampling

input : Target LLM \mathcal{M}_q , a small assistant model \mathcal{M}_p , initial prompt sequence x_1, \dots, x_t and target sequence length T .

- 1: Initialize $t \leftarrow 1$
- 2: **while** $t < T$ **do**
- 3: **for** $k \leftarrow 1, \dots, K$ **do**
- 4: $x_{t_k} \sim \mathcal{M}_p(x|x_1, \dots, x_t, x_{t_1}, \dots, x_{t_{k-1}})$
- 5: **end for**
- 6: In parallel, compute $K + 1$ sets of logits drafts x_{t_1}, \dots, x_{t_K} with the target LLM \mathcal{M}_q :
 $\mathcal{M}_q(x|x_1, \dots, x_t), \mathcal{M}_q(x|x_1, \dots, x_t, x_{t_1}), \dots, \mathcal{M}_q(x|x_1, \dots, x_t, x_{t_1}, \dots, x_{t_K})$
- 7: **for** $j \leftarrow 1, \dots, K$ **do**
- 8: Sample $r \sim U[0, 1]$ from a uniform distribution
- 9: **if** $r < \min(1, \frac{\mathcal{M}_q(x|x_1, \dots, x_{t+j-1})}{\mathcal{M}_p(x|x_1, \dots, x_{t+j-1})})$ **then**
- 10: Set $x_{t+j} \leftarrow x_{t_j}$ and $t \leftarrow t + 1$
- 11: **else**
- 12: Sample $x_{t+j} \sim (\mathcal{M}_q(x|x_1, \dots, x_{t+j-1}) - \mathcal{M}_p(x|x_1, \dots, x_{t+j-1}))_+$ and exit for loop.
- 13: **end if**
- 14: **end for**
- 15: If all tokens x_{t+1}, \dots, x_{t+K} are accepted, sample extra token $x_{t+K+1} \sim \mathcal{M}_q(x|x_1, \dots, x_t, x_{t+K})$ and set $t \leftarrow t + 1$
- 16: **end while**

enhancing accuracy and contextual relevance.

F Implementation details

F.1 Baselines

Following the Spec-Bench settings (Xia et al., 2024), we have selected 5 speculative decoding methods, all open-source and rigorously tested for reliability. Each method represents a unique approach to improving LLM inference speeds:

1. **SpS** (Chen et al., 2023): SpS employs a smaller LM from the same model series as the drafter. In the verification, this method corrects the last token with residual probability if the token is rejected.
2. **Medusa** (Cai et al., 2024) and **Eagle** (Li et al., 2024): Both methods enhance the target LLM by integrating additional lightweight FFN heads. These heads are designed to efficiently draft potential token sequences depending on the penultimate representations from the target LLM.
3. **Lookahead** (Fu et al., 2024): This method appends multiple special tokens to the end of the input prompt. These tokens are used for parallel drafting, with the resultant drafts transformed into n-gram candidates for efficient prediction.

4. **PLD** (Saxena, 2023): Serving as the practical code implementation of Yang et al. (2023), PLD selects text spans directly from the input to serve as drafts, optimizing the relevance and accuracy of the initial predictions.

F.2 Self-distillation

We follow the self-distillation pipeline as described by Cai et al. (2024). Initially, a public dataset, such as WMT 16 De-En, is selected as the training dataset. The target model’s responses are then generated using the OpenAI API server, with input prompts derived directly from the training dataset.

Install prerequisites For software dependencies, CUDA 12.1 and PyTorch 2.1.2 are required. To start the server, install the necessary dependencies:

```
vllm==0.4.0, openai==0.28.0
```

Use of vLLM We utilize the vLLM library for self-distillation, executing the following command:

```
python -m
vllm.entrypoints.openai.api_server
--model lmsys/vicuna-7b-v1.3
--port 8000 --max-model-len 2048
```

Input prompt For instance, when self-distillation the WMT14 Fr-En dataset using the

Table 4: Custom Gemma 250M model configuration.

Configuration	Value
Activation function	GeLU (Hendrycks and Gimpel, 2016)
Hidden size	768
Intermediate size	6144
Number of attention heads	16
Number of hidden layers	2
Number of key-value heads	2
RMS epsilon	1e-06
Vocabulary size	256000

Vicuna7b v1.3 model, the input prompt consists of a system prompt and a user prompt. In the user prompt, we prepend "Translate French to English: ".

A chat between a curious user and an artificial intelligence assistant. The assistant gives helpful, detailed, and polite answers to the user’s questions. USER: Translate French to English: Madame la Présidente, c’est une motion de procédure. ASSISTANT:

F.3 Details on training setup

For the shared settings across all training drafters, we employ the Fastchat³ framework. We utilize a cosine learning rate scheduler with a warmup ratio of 0.03 and the AdamW (Loshchilov and Hutter, 2017) optimizer. The drafter is trained using the ‘P-F’ strategy (ours) for 3 epochs, and using the ‘F’ strategy (without the pretraining step ‘P’) for 5 epochs to ensure sufficient learning. The model’s maximum length is set to 2048 tokens. The training is conducted using 4 GPUs with a batch size of 2 per GPU.

For finetuning the Vicuna 68M drafter (Yang et al., 2024), the learning rate is set to 2e-5. Similarly, for finetuning the Llama 68M model (Miao et al., 2024), the learning rate is set to 3e-5.

As a drafter for Gemma-Instruct 7B model, we newly design a Gemma 250M model as a drafter (Table 4). We use the same training recipe with Vicuna 68M and Llama 68M.

F.4 Details on GPT-4o evaluation

We follow LLM-as-a-Judge framework (Zheng et al., 2024) to evaluate the model’s answers. The GPT-4o model is utilized as a judge, which has greater performance on both English and non-English than GPT-4 Turbo (OpenAI, 2024). For

³<https://github.com/lm-sys/FastChat/tree/main>

Single answer grading, used prompt is followed:

[System]
You are a helpful assistant. Please act as an impartial judge and evaluate the quality of the response provided by an AI assistant to the user question displayed below. Your evaluation should consider factors such as the helpfulness, relevance, accuracy, depth, creativity, and level of detail of the response. Begin your evaluation by providing a short explanation. Be as objective as possible. After providing your explanation, you must rate the response on a scale of 1 to 10 by strictly following this format: "[[rating]]", for example: "Rating: [[5]]".

[Question]
{question}

[The Start of Assistant’s Answer]
{answer}
[The End of Assistant’s Answer]

The detail implementation of LLM-as-a-judge is in the following GitHub repository⁴.

G Additional experimental results

G.1 Out-of-domain speedup

Building on the findings discussed in the main body, this subsection further explores the speedup variations achieved by employing a drafter trained on each dataset across a range of translation tasks. Figure 8 depicts the speedup results using speculative greedy sampling for drafters trained on different datasets: Ru-En, Ja-En, and Zh-En.

Most observations align with those discussed in Section 3. Notably, drafters trained on the Ja-En (Figure 8 (b)) and Zh-En (Figure 8 (c)) datasets consistently outperform Yang et al. (2024)’s drafter, even on out-of-domain tasks. We hypothesize these into two folds. Firstly, this suggests that certain intrinsic properties of the Japanese and Chinese languages may improve the efficacy of speculative decoding when applied to unrelated language pairs, possibly due to specific syntactic or lexical features that are effectively captured during training. In another scenario, the target LLM does not work well on those tasks, and thus drafters are easier to catch

⁴https://github.com/lm-sys/FastChat/tree/main/fastchat/llm_judge

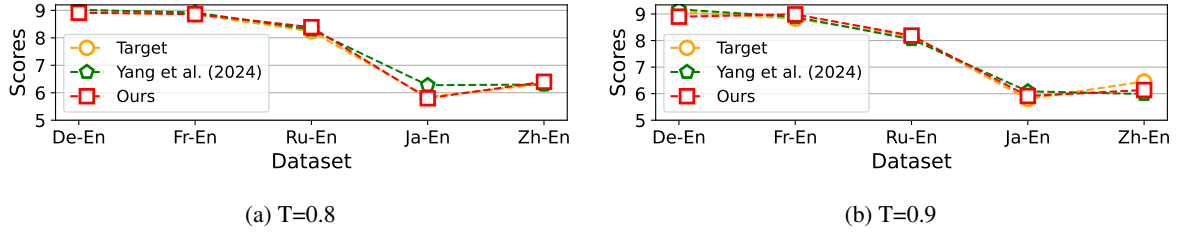


Figure 7: GPT-4o evaluation scores following the Zheng et al. (2024) on various multilingual translation dataset. Each figure denotes the score of random sampling with different temperature on the output whose target LLM is Vicuna 7B v1.3.

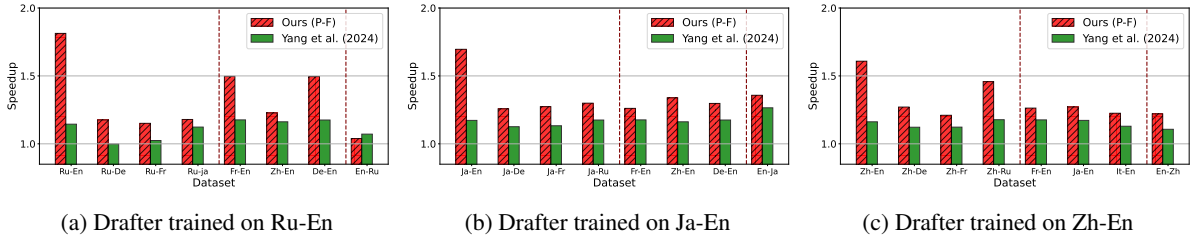


Figure 8: Speedup with speculative greedy sampling with the same settings in Figure 5.

766 the target token distribution. More precisely, for
 767 instance, in Zh-Ru task, Vicuna 7B should translate
 768 the Chinese to Russian, but to English, and thus
 769 the speedup seems to happen for us due to English
 770 generation.

771 In the case of the Ru-En (Figure 8 (a)) trained
 772 drafter, translations from Russian to other lan-
 773 guages generally surpass Yang et al. (2024)’s re-
 774 sults. Interestingly, translations from French to
 775 English and German to English exhibit unexpect-
 776 edly high speedups. This could hint at underlying
 777 linguistic similarities or shared grammatical struc-
 778 tures between Russian, French, and German that
 779 the Ru-En drafter is particularly adept at handling,
 780 thereby facilitating more efficient speculative dec-
 781 oding. While Fan et al. (2021) demonstrates that
 782 Russian belongs to another cluster from En / Fr /
 783 De, perhaps our results provide a different perspec-
 784 tive in lens of speculative decoding.

785 G.2 GPT-4o judgments

786 Figure 7 show additional GPT-4o evaluation scores
 787 for various multilingual translation datasets. The
 788 graphs display the comparative performance across
 789 different language pairs under two sampling con-
 790 ditions, at temperatures $T=0.8$ and $T=0.9$, respec-
 791 tively. Each data point reflects the quality of trans-
 792 lations produced by the target model (orange cir-
 793 cle), SpS with the instruction tuned model using
 794 ShareGPT (Yang et al., 2024) (green pentagon),
 795 and SpS with our specialized drafter (*pretrain-and-*

finetune) (red square). For the red points, each
 796 drafter is trained with the corresponding dataset.
 797 For instance, when the red point specify De-En, it
 798 indicates that the drafter has been fine-tuned with
 799 the De-En dataset.
 800

801 The results demonstrate negligible differences
 802 in quality among the three methods, underscoring
 803 the efficacy of speculative decoding in delivering
 804 translations with lossless quality. Both tempera-
 805 ture settings show that our speculative decoding
 806 strategy closely matches the performance of the
 807 established target model across various language
 808 pairs. This consistent performance across different
 809 settings and language pairs illustrates that specu-
 810 lative decoding effectively maintains high-quality
 811 outputs without compromising accuracy due to in-
 812 creased randomness in sampling.

813 H Discussion

814 H.1 Why is pretrain-and-finetune better in 815 small-size LM drafter?

816 Drafting in speculative decoding has been treated
 817 akin to n-gram prediction (Bhendawade et al.,
 818 2024), often relying on straightforward pretrain-
 819 ing using datasets designed to replicate target LLM
 820 behaviors, such as the ShareGPT dataset (Yang
 821 et al., 2024). This approach posits that generat-
 822 ing a limited sequence of future tokens suffices for
 823 speculative inference.

824 Contrary to this belief, our empirical result

Table 5: Speedup comparison of speculative greedy sampling across different drafter sizes on WMT16 De-En dataset.

Drafter	Vicuna 68M (Yang et al., 2024)	Vicuna 68M (<i>pretrain-and-finetune</i> ; Ours)	Tiny-Vicuna 1B (Pan, 2023)
Speedup	1.19	2.42	0.75
Mean of accepted tokens	1.47	3.03	3.06

presents a different narrative. Figure 5 illustrates that even in seemingly straightforward translation tasks, such as from German to English, outcomes are not as effective. This suggests that drafting requires a broader array of language modeling capabilities to manage complex linguistic structures and context variations effectively.

Drafters, therefore, benefit significantly from a robust *pretrain-and-finetune* approach, where they are first exposed to a wide array of linguistic contexts and then finely tuned to specific tasks. This training regimen transforms them into compact, yet comprehensive, language models capable of handling diverse and challenging speculative decoding scenarios with better alignment.

H.2 Number of drafts

This study primarily explores the speculative decoding process utilizing a single draft. In contrast, advanced baseline methods such as EAGLE and Medusa deploy multiple drafts, leveraging tree-attention mechanisms to enrich draft selection. This technique allows for a broader exploration of multiple draft candidates at each decoding step, potentially increasing the rate and quality of accepted drafts.

Adapting our approach to incorporate multiple drafts with tree-attention could significantly enhance performance, suggesting an untapped potential in our method. Experimenting with this expanded setup could lead to notable improvements in the speculative sampling’s effectiveness, particularly in increasing the mean number of high-quality tokens accepted per sequence. This prospect opens a critical path for future research, where deeper explorations could elevate the capabilities of our specialized drafters.

H.3 Is scaling up drafter size better for SpS?

Evaluating the efficacy of increasing drafter size reveals nuanced insights into speculative decoding performance. Table 5 compares three versions of drafters: the Vicuna 68M by Yang et al. (2024), our *pretrain-and-finetune* Vicuna 68M, and Tiny-Vicuna 1B (Pan, 2023)—a larger model with 1B

parameters that has been instruction-tuned.

Despite Tiny-Vicuna 1B’s substantial parameter count, it achieves a lower speedup of 0.75 compared to 2.34 by our optimized Vicuna 68M. Both models show similar mean accepted tokens, suggesting that increasing size does not proportionally enhance computational efficiency. This is due to speculative decoding’s reliance on minimizing memory bottlenecks to exploit parallel computation effectively. Larger models like Tiny-Vicuna 1B exacerbate these bottlenecks, diminishing the potential speed gains from increased parallelism.

Conversely, our *pretrain-and-finetune* Vicuna 68M demonstrates that strategic training and optimization of a smaller model can achieve high efficiency and speed, highlighting the importance of model configuration over mere size increase. This balance between model size and computational dynamics is crucial for optimizing speculative decoding, suggesting that enhancing model capabilities through targeted training may be more effective than scaling size.