



EAGLE-3: Scaling up Inference Acceleration of Large Language Models via Training-Time Test

Yuhui Li^{1,3}, Fangyun Wei², Chao Zhang¹, Hongyang Zhang^{3,4}

¹Peking University ²Microsoft Research

³University of Waterloo ⁴Vector Institute

yuhui.li@stu.pku.edu.cn, fawe@microsoft.com

c.zhang@pku.edu.cn, hongyang.zhang@uwaterloo.ca

Abstract

The sequential nature of modern LLMs makes them expensive and slow, and speculative sampling has proven to be an effective solution to this problem. **Methods like EAGLE perform autoregression at the feature level, reusing top-layer features from the target model to achieve better results than vanilla speculative sampling.** A growing trend in the LLM community is scaling up training data to improve model intelligence without increasing inference costs. **However, we observe that scaling up data provides limited improvements for EAGLE. We identify that this limitation arises from EAGLE’s feature prediction constraints.** In this paper, we introduce EAGLE-3, which **abandons feature prediction in favor of direct token prediction and replaces reliance on top-layer features with multi-layer feature fusion via a technique named training-time test.** These improvements significantly enhance performance and enable the draft model to fully benefit from scaling up training data. Our experiments include both chat models and reasoning models, evaluated on five tasks. **The results show that EAGLE-3 achieves a speedup ratio up to 6.5x, with about 1.4x improvement over EAGLE-2.** In the SGLang framework, EAGLE-3 achieves a 1.38x throughput improvement at a batch size of 64. The code is available at <https://github.com/SafeAILab/EAGLE>.

1 Introduction

Modern Large Language Models (LLMs) are being applied to more domains, with their improved capabilities driven by scaling model parameters—some LLMs now exceed hundreds of billions of parameters. In autoregressive generation, each token requires accessing all model parameters, making LLM inference slow and costly.

Recently, test-time scaling up has gained significant attention. Models like ChatGPT o1 and

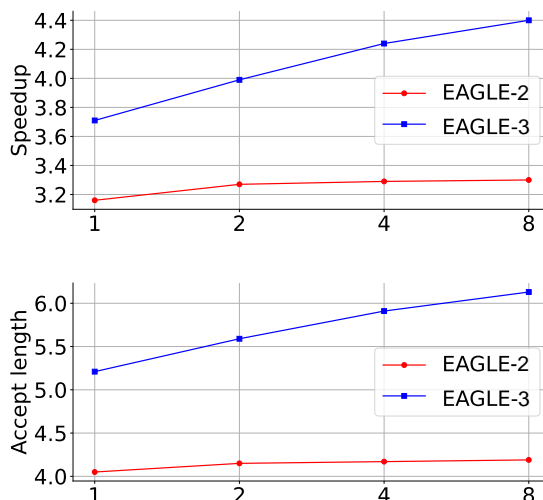


Figure 1: Scaling law evaluated on the MT-bench using LLaMA-Instruct 3.1 8B as the target model, with the x-axis representing the data scale relative to ShareGPT. The new architectural designs in EAGLE-3 enable an increasing scaling curve, which was never observed in the previous works.

DeepSeek-R1 (Guo et al., 2025) engage in deliberate reasoning before responding, pushing the boundaries of LLM capabilities at the cost of longer inference time. However, these models often require lengthy reasoning processes, making them extremely costly, while the increased response time severely impacts user satisfaction. These reasoning models significantly increase the proportion of inference costs in the overall LLM pipeline, driving researchers to explore cheaper and faster inference optimization methods.

Speculative sampling methods can reduce LLM latency by partially parallelizing the generation process. These methods rapidly generate draft tokens and then verify them in parallel. This allows multiple tokens to be produced in a single forward pass, significantly reducing inference latency. In vanilla speculative sampling, the draft model is a separate, smaller LLM, typically a lower-parameter version from the same series as the target model. This

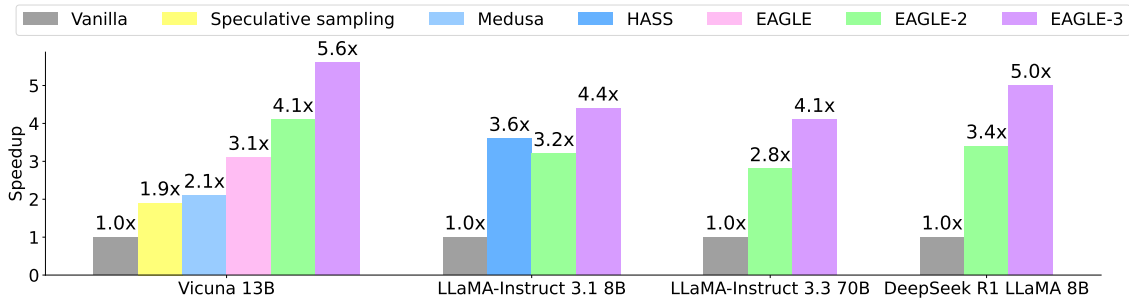


Figure 2: Speedup ratios of different methods at temperature=0. For the standard speculative sampling, Vicuna-13B uses Vicuna-68M as the draft model. In Table 1, we present comparisons with additional methods, but this figure only showcases a subset. Chat model’s evaluation dataset is MT-bench, and the reasoning model’s evaluation dataset is GSM8K. DeepSeek R1 LLaMA 8B refers to DeepSeek-R1-Distill-LLaMA 8B.

draft model operates independently of the target model. Unlike the vanilla speculative sampling, EAGLE (Li et al., 2024c) reuses the top-layer features of the target model (the features before the LM head). It trains the draft model to autoregressively predict the next feature and then uses the target model’s LM head to obtain the draft token. By leveraging the rich information from the target model, EAGLE achieves significantly better acceleration compared to vanilla speculative sampling. Subsequent methods such as HASS (Zhang et al., 2024) and Falcon (Gao et al., 2024) also adopt the approach of predicting the next feature using the current feature sequence.

Recent LLMs have increasingly relied on larger training datasets to achieve better performance. For example, LLaMA series models with sizes of 7B (8B) have used 1T, 2T, and 15T tokens of training data for LLaMA 1 (Touvron et al., 2023a), LLaMA 2 (Touvron et al., 2023b), and LLaMA 3 (Dubey et al., 2024), respectively, resulting in significant improvements across various metrics while keeping the model architecture and inference cost largely unchanged. Similarly, we aim to improve the acceptance rate and acceleration ratio of EAGLE by increasing its training data. Unfortunately, we observe that the gains from additional training data for EAGLE are limited. We analyze the reasons behind this phenomenon. As shown in the upper part of Figure 3, EAGLE performs autoregressive prediction at the feature level, predicting the next feature and then feeding the feature into the LM head of the target model to obtain the token distribution. EAGLE’s loss function consists of two components: the feature prediction loss l_{fea} and the token prediction loss l_{token} . Thanks to the feature prediction loss, the draft model trained only at Step

1 can adapt to Step 2 and acquire multi-step prediction capabilities. However, with token prediction as the ultimate goal, feature prediction can be seen as an additional constraint, which limits the expressiveness of the draft model and makes it difficult to benefit from increased data. After removing the feature constraint and expanding the training data (the middle part of Figure 3), as shown in Figure 4, the acceptance rate $0-\alpha$ of the first draft token improves significantly. However, the output of the draft model in Step 1, denoted as \hat{a}_{t+1} , is far away from the ground-truth f_{t+1} , causing the input sequence $f_1, f_2, \dots, f_t, \hat{a}_{t+1}$ in Step 2 to deviate significantly from the training distribution, resulting in a very low acceptance rate $1-\alpha$ for the second draft token, as shown in Figure 4. We can address this issue by incorporating Step 1 into the training process (the bottom of Figure 3). Using this method, the benefits of increasing training data become more pronounced. We name this technique as training-time test.

EAGLE and speculative sampling methods such as Medusa (Cai et al., 2024) reuse the top-layer features of the target model, specifically the features immediately before the LM head. For an LM head with a full-rank weight matrix, the top-layer features corresponding to the logits of the next token are unique, ensuring that the information contained in these features aligns directly with the logits of the next token. However, predicting the next-next token based solely on top-layer features—which are inherently limited to the next token—poses a significant challenge. Fortunately, the training-time test technique described above enables the use of features from intermediate layers instead of relying solely on the top layer, as the feature prediction loss l_{fea} has been removed during training.

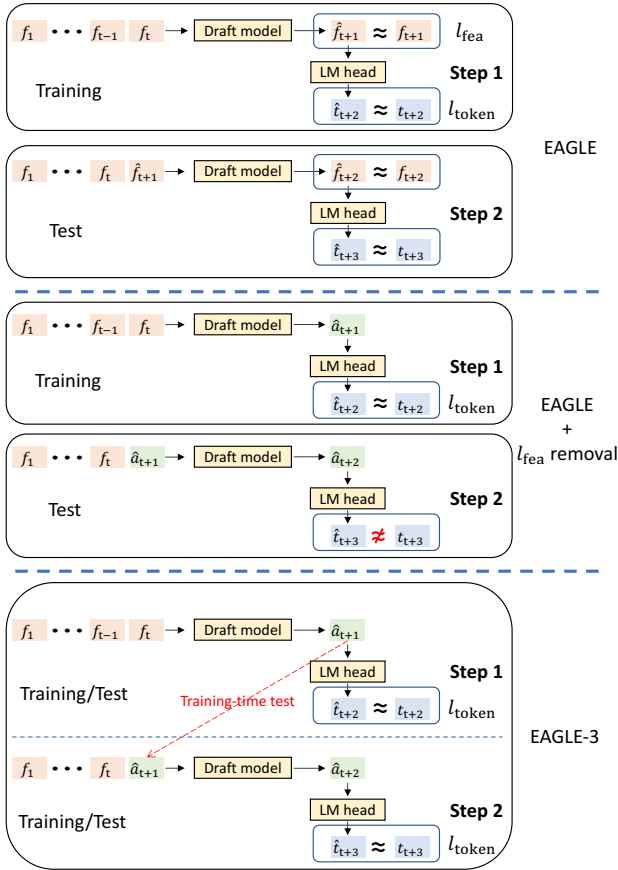


Figure 3: Illustration of **training-time test** (the bottom part) and its comparison with other draft methods (the upper and middle parts). f denotes the feature, t denotes the token, and a represents the unconstrained vectors. We use the hat to denote the predictions from models. All the methods shown in the figure use the token sequence from the previous time step, but for simplicity, this is not depicted in the figure. The input to EAGLE-3 is not actually f , but it is not shown in this figure. We will provide a detailed explanation in the following section.

To summarize, this paper introduces EAGLE-3, an enhanced version of EAGLE that achieves a significant speedup. EAGLE-3 is parallelized and fully compatible with the drafting tree technique from EAGLE-2 (Li et al., 2024b). Our key contributions include:

- **A novel training-time test architecture for the draft model:** We remove the feature prediction constraint and directly predict tokens while simulating multi-step generation during training. This direct token prediction provides complete flexibility in the draft model’s input. Instead of reusing only the top-layer features, we integrate and leverage low-, mid-, and high-level features from the target model, capturing rich semantic information from different lay-

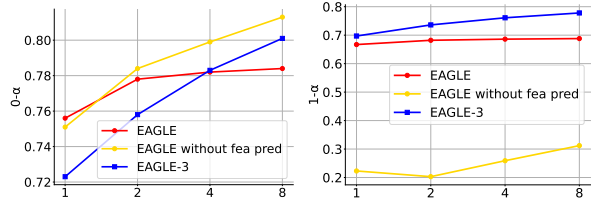


Figure 4: Comparison of acceptance rates across different methods, with the x-axis representing the data scale relative to ShareGPT.

ers.

- **Discovery of a scaling law for inference acceleration in large language models:** With the new architecture, we observe that increasing the amount of training data for the draft model leads to a proportional increase in the speedup ratio of EAGLE-3. This scaling behavior was not observed in the original EAGLE architecture, as shown in Figure 1

- **Improved inference acceleration:** EAGLE-3, trained with approximately 8x more data than EAGLE, achieves a 1.4x latency speedup over EAGLE-2 at batch size 1. Speculative sampling is often thought to reduce throughput at large batch sizes. However, in SGLang (Zheng et al., 2024), a production-grade framework, EAGLE-3 improves throughput by 40% at a batch size of 64. We expect larger data size would lead to further improved speedup ratio.

2 Preliminaries

2.1 Speculative Sampling

Speculative sampling (Leviathan et al., 2023; Chen et al., 2023; Sun et al., 2024c,b) is a lossless LLM acceleration technique that alternates between drafting and verification, where drafting is performed at low cost and verification is parallelized, corresponding to the generation of drafts and the verification process, respectively. We use t_i to denote the i -th token and $T_{a:b}$ to represent the token sequence t_a, t_{a+1}, \dots, t_b . When $T_{1:j}$ is used as the prefix, the two stages of speculative sampling are as follows.

In the drafting stage, speculative sampling utilizes a draft model (a smaller version from the same series as the target model) to autoregressively generate k tokens to form the draft. $\hat{T}_{j+1:j+k}$, while also recording the probability \hat{p} for each token.

In the verification stage, speculative sampling invokes the target model to evaluate the draft $\hat{T}_{j+1:j+k}$ and records its probability p . Speculative sampling then determines the acceptance of draft tokens sequentially, from front to back. For token \hat{t}_{j+i} , the probability of acceptance is given by $\min(1, p_{j+i}(\hat{t}_{j+i})/\hat{p}_{j+i}(\hat{t}_{j+i}))$. If the token is accepted, the process moves to the next token. Otherwise, a token is sampled from the distribution $\text{norm}(\max(0, p_{j+i} - \hat{p}_{j+i}))$ to replace \hat{t}_{j+i} , and the remaining tokens in the draft are discarded. Appendix A.1 of (Leviathan et al., 2023) proves that speculative sampling is consistent with the distribution of vanilla autoregressive decoding.

2.2 EAGLE and EAGLE-2

The draft model with limited capacity struggles to precisely approximate the large-scale target model. EAGLE leverages the top-layer features of the target model as additional information and performs autoregression at the feature level, simplifying the drafting process. EAGLE performs autoregression at the feature level and then uses the LM head of the target model to obtain the draft token. Due to the sampling results at the token layer being hidden, feature-level autoregression introduces uncertainty. EAGLE addresses this issue by feeding the token sequence from the previous time step, i.e., the sampling results, into the draft model. Unlike the chain-like drafts of Vanilla speculative sampling, EAGLE generates multiple draft tokens at the same position, resulting in a tree-like draft. In the verification stage, EAGLE uses tree attention to parallelize the verification of the draft tree. Interestingly, EAGLE inspired the *multi-token prediction* technique used in the pre-training of DeepSeek-v3 (Liu et al., 2024a), which in turn inspired new architectural designs in EAGLE-3.

EAGLE (Li et al., 2024c) and Medusa (Cai et al., 2024), among others, use tree-shaped drafts, where the structure of the draft tree is predefined, static, and context-independent. The difficulty of drafting is closely related to the context, and a static draft tree can lead to resource wastage. EAGLE-2 (Li et al., 2024b) approximates the acceptance rate using the confidence of the draft model and dynamically generates the draft tree based on this, performing pruning of the draft tree at the end of the drafting stage. EAGLE-3 also adopts the context-aware dynamic draft tree proposed in EAGLE-2.

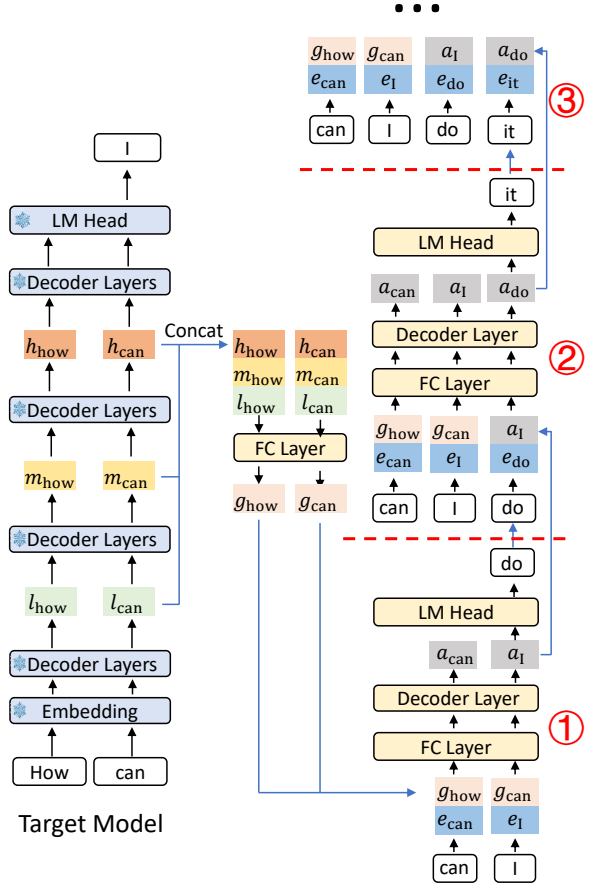


Figure 5: Diagram of the EAGLE-3 inference pipeline, illustrating the three steps of the draft model. l , m , and h represent the low, middle, and high-level features of the target model, respectively. e denotes the embedding.

3 EAGLE-3

In this section, we provide a detailed description of the implementation of EAGLE-3.

3.1 Inference Pipeline

Consistent with other speculative sampling methods, EAGLE-3 alternates between the drafting and verification stages. The difference between EAGLE-3 and EAGLE lies in the drafting stage, which we introduce with an example, as shown in Figure 5. Consider the prefix “How can”. During the prefill phase or the previous verification stage, the target model performs a forward pass to generate the next token, “I”. We record the low, middle, and high-level feature sequences from the target model’s forward pass, denoted as l , m , and h , respectively. We concatenate the k -dimensional vectors l , m , and h to form a $3k$ -dimensional vector, then pass it through a fully connected (FC) layer to reduce it to k -dimensions, obtaining a feature g that integrates information from different layers. Here, k refers to the hidden size of the target model.

Our goal is to generate a draft token sequence with the prefix “How can I”. By inputting only g_{how} and g_{can} , the draft model cannot access the random sampling process. Therefore, similar to EAGLE (Li et al., 2024c), we introduce the embedding e_I of the sampled token “I”. The concatenated vector is then passed through an FC layer to reduce its dimensionality to k , and subsequently inputted into a single layer decoder, producing the output a . Finally, we input a_I into the LM head and sample to obtain the draft token “do”.

In Step 1, with the prefix “How can”, we reuse g_{how} and g_{can} from the target model. In Step 2, the prefix becomes “How can I”. Ideally, we would reuse g_{how} , g_{can} , and g_I from the target model. However, this is not possible because the token “I” has not yet been checked by the target model, and we cannot obtain g_I . Instead, we use the output a_I from the draft model in the previous step to replace g_I , and concatenate a_I with the embedding e_{do} of the sampled result “do” as the input to the draft model in Step 1. In Step 3, we similarly cannot obtain g_{do} , so we use a_{do} as a replacement, concatenating a_{do} with e_{it} as the input to the draft model. The same approach is followed for subsequent steps.

3.2 Draft Model Training

The input to the draft model in EAGLE is either, or at least approximately, the top-layer features f_1, f_2, \dots, f_t of the target model. In contrast, the input to the draft model in EAGLE-3 may include the features g_1, g_2, \dots, g_t from the target model, or it may include the output $a_{t+1}, a_{t+2}, \dots, a_{t+j}$ from the draft model. Therefore, we need to train the draft model to adapt to different inputs. During training, we perform test steps, where we generate a and feed it back into the draft model for further training.

The core of the draft model in EAGLE-3 is a Transformer decoder layer. Aside from the self-attention operation, no other components interact with the context, so no further modifications are required during training or testing. The only component that requires slight modification is the self-attention, which we will describe in detail below.

Although the actual input consists of features, for clarity, we describe the process using tokens as input. As shown in Figure 6, the original training data is a sequence of length 3, “How can I”, with a normal sequential dependency in the context. Therefore, the attention mask is a standard lower triangular matrix. The outputs at the three positions

are “are”, “we”, and “do”, which have a tree-like contextual relationship with “how”, “can”, and “I”. As a result, when the input “are”, “we”, and “do” is fed into Step 2, the attention mask needs to be adjusted accordingly, as shown in the top-right corner of Figure 6. All attention masks are diagonal, except when the original training data is used as the key. Using matrix multiplication in this case would result in significant computational waste, so we can use vector dot products to calculate the attention score only for the corresponding positions.

HASS (Zhang et al., 2024) and EAGLE-3 both make similar modifications to the attention mechanism to simulate the testing process during training, but this is not the main focus of EAGLE-3. The motivations, methods, and outcomes of the two approaches are distinctly different. The motivation behind HASS is to mitigate the error accumulation caused by inaccurate feature predictions in EAGLE. HASS still performs feature prediction, includes a feature prediction loss l_{fea} , and the input to the draft model must be the top-layer features. In contrast, the motivation behind EAGLE-3 is to remove unnecessary constraints to enhance the model’s expressive power. EAGLE-3 no longer requires the draft model’s output to fit the top-layer features of the target model, thus avoiding error accumulation. After removing feature prediction, the input to EAGLE-3 is completely free, and it is replaced by a fusion of features from different layers of semantic information. The removal of the feature prediction loss also enables us to discover a new scaling law for inference acceleration which was never found before. Figure 2 also shows the speedup of EAGLE-3 and HASS, with EAGLE-3 demonstrating significantly better performance.

4 Experiments

Models. We conduct experiments with state-of-the-art open-source chat and reasoning models, including Vicuna 13B (Chiang et al., 2023), LLaMA-Instruct 3.1 8B, LLaMA-Instruct 3.3 70B (Dubey et al., 2024), and DeepSeek-R1-Distill-LLaMA 8B (DeepSeek-AI et al., 2025). Due to the GPU constraint, we are unable to test EAGLE-3 on the 405B and 671B models.

Tasks. Following EAGLE (Li et al., 2024c) and Spec-Bench (Xia et al., 2024), we evaluate on five common tasks, using the same weights for all tasks without fine-tuning on the respective tasks. For multi-turn conversation, code genera-

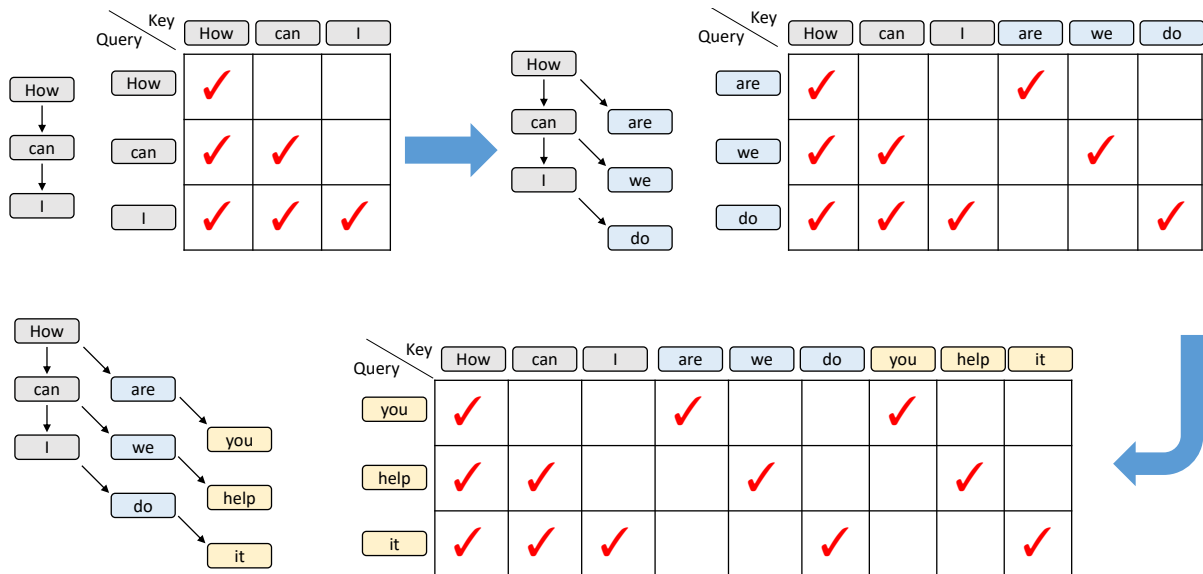


Figure 6: Diagram of the attention causal masks during training-time test. It sequentially shows a native training step (the first step) and two simulated training steps (the second and third steps). The arrows between tokens represent contextual relationships. The gray tokens represent the training data while the blue and yellow tokens represent the first- and second-round predictions by the draft model, respectively.

tion, mathematical reasoning, instruction following, and summarization, we chose the MT-bench (Zheng et al., 2023), HumanEval (Chen et al., 2021), GSM8K (Cobbe et al., 2021), Alpaca (Taori et al., 2023), and CNN/Daily Mail (Nallapati et al., 2016) datasets, respectively.

Metrics. EAGLE-3 does not modify the target model’s weights and uses strict speculative sampling acceptance conditions, ensuring no loss in performance. Therefore, we do not evaluate generation quality. Instead, we use the following metrics to assess the acceleration performance:

- **Speedup Ratio:** The actual test speedup ratio relative to vanilla autoregressive decoding.
- **Average Acceptance Length τ :** The average number of tokens generated per drafting-verification cycle, which corresponds to the number of tokens accepted from the draft.
- **Acceptance Rate $n-\alpha$:** The proportion of draft tokens accepted, which directly reflects the draft model’s approximation to the target model. Following EAGLE’s setup, we use a chain-like draft rather than a tree-like draft when testing acceptance rates. EAGLE suffers from error accumulation, meaning that the input to the draft model may be its own estimates rather than the exact values from the target model. Therefore, EAGLE uses $n-\alpha$ to represent the acceptance

rate when the input contains n estimated features, under the condition that the previous estimated tokens are all accepted by the target model. In other words, the acceptance rate for inputs $f_1, f_2, \dots, f_i, \hat{f}_{i+1}, \dots, \hat{f}_{i+n}$, where f is the exact value and \hat{f} is the draft model’s estimate. Similarly, we use $n-\alpha$ to represent the acceptance rate in EAGLE-3 when the input contains n self-predicted values a , i.e., the acceptance rate for inputs $g_1, g_2, \dots, g_i, a_{i+1}, \dots, a_{i+n}$, where g is the fused feature from the target model.

Implementation. We use the AdamW optimizer, with beta values (β_1, β_2) set to (0.9, 0.95) and implemented gradient clipping of 0.5. The learning rate is set to $5e-5$. We use ShareGPT and UltraChat-200K (Ding et al., 2023) as training data, containing approximately 68K and 464K data entries, respectively. We call the target model to generate responses rather than using a fixed dataset. For the reasoning model DeepSeek-R1-Distill-LLaMA 8B, we also used the OpenThoughts-114k-math dataset for training.

Comparison. We use vanilla autoregressive decoding as the baseline, which serves as the benchmark for speedup ratios (1.00x). We compare EAGLE-3 with recent lossless speculative sampling methods, including standard speculative sampling (Leviathan et al., 2023; Chen et al., 2023; Gante, 2023), PLD (Saxena, 2023), Medusa (Cai

Table 1: Speedup ratios and average acceptance lengths τ of different methods. V represents Vicuna, L31 represents LLaMA-Instruct 3.1, L33 represents LLaMA-Instruct 3.3, and DSL represents DeepSeek-R1-Distill-LLaMA. SpS denotes standard speculative sampling, with its draft model being Vicuna-68M. Methods like Medusa relax acceptance conditions under non-greedy settings, which do not guarantee lossless acceleration. Therefore, we do not compare EAGLE-3 with these methods when temperature=1.

Model	Method	MT-bench		HumanEval		GSM8K		Alpaca		CNN/DM		Mean	
		Speedup	τ	Speedup	τ	Speedup	τ	Speedup	τ	Speedup	τ	Speedup	τ
Temperature=0													
V 13B	SpS	1.93x	2.27	2.23x	2.57	1.77x	2.01	1.76x	2.03	1.93x	2.33	1.92x	2.24
	PLD	1.58x	1.63	1.85x	1.93	1.68x	1.73	1.16x	1.19	2.42x	2.50	1.74x	1.80
	Medusa	2.07x	2.59	2.50x	2.78	2.23x	2.64	2.08x	2.45	1.71x	2.09	2.12x	2.51
	Lookahead	1.65x	1.69	1.71x	1.75	1.81x	1.90	1.46x	1.51	1.46x	1.50	1.62x	1.67
	Hydra	2.88x	3.65	3.28x	3.87	2.93x	3.66	2.86x	3.53	2.05x	2.81	2.80x	3.50
	EAGLE	3.07x	3.98	3.58x	4.39	3.08x	3.97	3.03x	3.95	2.49x	3.52	3.05x	3.96
	EAGLE-2	4.26x	4.83	4.96x	5.41	4.22x	4.79	4.25x	4.89	3.40x	4.21	4.22x	4.83
	EAGLE-3	5.58x	6.65	6.47x	7.54	5.32x	6.29	5.16x	6.17	5.01x	6.47	5.51x	6.62
L31 8B	EAGLE-2	3.16x	4.05	3.66x	4.71	3.39x	4.24	3.28x	4.12	2.65x	3.45	3.23x	4.11
	EAGLE-3	4.40x	6.13	4.85x	6.74	4.48x	6.23	4.82x	6.70	3.65x	5.34	4.44x	6.23
L33 70B	EAGLE-2	2.83x	3.67	3.12x	4.09	2.83x	3.69	3.03x	3.92	2.44x	3.55	2.85x	3.78
	EAGLE-3	4.11x	5.63	4.79x	6.52	4.34x	6.15	4.30x	6.09	3.27x	5.02	4.12x	5.88
DSL 8B	EAGLE-2	2.92x	3.80	3.42x	4.29	3.40x	4.40	3.01x	3.80	3.53x	3.33	3.26x	3.92
	EAGLE-3	4.05x	5.58	4.59x	6.38	5.01x	6.93	3.65x	5.37	3.52x	4.92	4.16x	5.84
Temperature=1													
V 13B	SpS	1.62x	1.84	1.72x	1.97	1.46x	1.73	1.52x	1.78	1.66x	1.89	1.60x	1.84
	EAGLE	2.32x	3.20	2.65x	3.63	2.57x	3.60	2.45x	3.57	2.23x	3.26	2.44x	3.45
	EAGLE-2	3.80x	4.40	4.22x	4.89	3.77x	4.41	3.78x	4.37	3.25x	3.97	3.76x	4.41
	EAGLE-3	4.57x	5.42	5.15x	6.22	4.71x	5.58	4.49x	5.39	4.33x	5.72	4.65x	5.67
L31 8B	EAGLE-2	2.44x	3.16	3.39x	4.39	2.86x	3.74	2.83x	3.65	2.44x	3.14	2.80x	3.62
	EAGLE-3	3.07x	4.24	4.13x	5.82	3.32x	4.59	3.90x	5.56	2.99x	4.39	3.45x	4.92
L33 70B	EAGLE-2	2.73x	3.51	2.89x	3.81	2.52x	3.36	2.77x	3.73	2.32x	3.27	2.65x	3.54
	EAGLE-3	3.96x	5.45	4.36x	6.16	4.17x	5.95	4.14x	5.87	3.11x	4.88	3.95x	5.66
DSL 8B	EAGLE-2	2.69x	3.41	3.01x	3.82	3.16x	4.05	2.64x	3.29	2.35x	3.13	2.77x	3.54
	EAGLE-3	3.20x	4.49	3.77x	5.28	4.38x	6.10	3.16x	4.30	3.08x	4.27	3.52x	4.89

et al., 2024), Lookahead (Fu et al., 2024), Hydra (Ankner et al., 2024), HASS (Zhang et al., 2024), EAGLE (Li et al., 2024c), and EAGLE-2 (Li et al., 2024b).

4.1 Effectiveness

Figure 1 and Table 1 demonstrate the acceleration performance of EAGLE-3. On all tasks and target models, EAGLE-3 achieves the highest speedup ratio and average acceptance length. EAGLE-3 provides a speedup of approximately 3.0x-6.5x compared to vanilla autoregressive generation, with a 20%-40% improvement over EAGLE-2. Different tasks affect the draft model’s acceptance rate, so both the average acceptance length and speedup ratio are task-dependent. Due to the presence of many fixed templates in code generation tasks, generating drafts is the easiest, which is why EAGLE-3 performs best on HumanEval, achieving a speedup ratio of up to 6.5x and an average acceptance length

of up to 7.5. DeepSeek-R1-Distill-LLaMA 8B is an exception, with the highest speedup ratio on the mathematical reasoning dataset GSM8K. This may be because we trained the draft model of DeepSeek-R1-Distill-LLaMA 8B using the OpenThoughts-114k-math dataset.

Figure 7 shows the acceptance rates of EAGLE and EAGLE-3 on MT-bench with LLaMA-Instruct 3.1 8B as the target model. The acceptance rate of EAGLE-3 is significantly higher than that of EAGLE. As the input from the draft model itself increases, the acceptance rate of EAGLE drops significantly, whereas EAGLE-3’s acceptance rate remains almost unchanged, demonstrating the effectiveness of the Training-time test.

4.2 Ablation Study

The improvements of EAGLE-3 mainly come from two aspects: first, the removal of the feature regression constraint, and second, the improvement

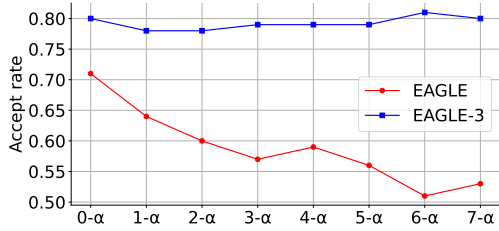


Figure 7: Acceptance rate of EAGLE and EAGLE-3 on MT-bench, with the target model being LLaMA-Instruct 3.1 8B. Hereby, $n-\alpha$ refers to the acceptance rate when the input contains n estimated features, under the condition that the previous estimated tokens are all accepted by the target model.

from reusing only the top-layer features to reusing a mix of low, middle, and high-level features. We conducted an ablation study on MT-bench with LLaMA-Instruct 3.1 8B as the target model. The results, shown in Table 2, indicate that both improvements in EAGLE-3 significantly enhance the acceptance length and speedup ratio, demonstrating the rationality of the EAGLE-3 design.

Table 2: Ablation study results with LLaMA-Instruct 3.1 8B as the target model. “Remove fea con” refers to the first improvement of EAGLE-3, which removes the feature prediction constraint. “Fused features” refers to the second improvement of EAGLE-3, where low, middle, and high-level feature fusion replaces the use of top-layer features.

Method	MT-bench		GSM8K	
	Speedup	τ	Speedup	τ
EAGLE-2	3.16x	4.05	3.39x	4.24
+ remove fea con	3.82x	5.37	3.77x	5.22
+ fused features (ours)	4.40x	6.13	4.48x	6.23

4.3 EAGLE-3 in SGLang

Speculative sampling algorithms like EAGLE-3 reduce memory accesses and lower latency during memory-bound decoding by leveraging redundant computational power. As batch sizes increase, this redundancy decreases, reducing the effectiveness of speculative sampling. Efficiency improvements are more challenging in highly optimized production-grade frameworks. The performance of EAGLE-3 for large batches on a single H100 GPU and LLaMA-Instruct 3.1 8B in the SGLang v0.4.4 environment (Zheng et al., 2024) was evaluated by the SGLang team, shown in Table 3. This part of the experiment did not use the tree structure, the chain length was set to 3, and the testing dataset was MT-Bench. EAGLE reduces throughput at a

batch size of 24, whereas EAGLE-3 still achieves a 38% throughput improvement at a batch size of 64.

Table 3: Throughput improvement under different batch sizes on H100 and LLaMA-Instruct 3.1 8B for the MT-Bench dataset, with SGLang without speculative sampling as the baseline (1.00x). The experiments were conducted by the SGLang team.

Batch size	2	4	8	16	24	32	48	56	64
EAGLE	1.40x	1.38x	1.23x	1.02x	0.93x	0.94x	0.88x	0.99x	0.99x
EAGLE-3	1.81x	1.82x	1.62x	1.48x	1.39x	1.32x	1.38x	1.34x	1.38x

The SGLang team also tested the throughput of EAGLE-3 at batch size = 1 on H100 when the target model is LLaMA-Instruct 3.1 8B and the testing dataset is MT-bench. The results are shown in Table 4.

Table 4: Throughput at batch size = 1 on H100 when the target model is LLaMA-Instruct 3.1 8B and the testing dataset is MT-bench. The experiments were conducted by the SGLang team.

Method	Throughput (bs=1)
SGLang (w/o speculative, 1x H100)	158.34 tokens/s
SGLang + EAGLE-2 (1x H100)	244.10 tokens/s
SGLang + EAGLE-3 (1x H100)	373.25 tokens/s

4.4 EAGLE-3 in vLLM

We also conducted a study on the impact of EAGLE-3 on throughput for large batch sizes based on vLLM (Kwon et al., 2023), a widely used production-grade framework, and the results on RTX3090 and LLaMA-Instruct 3.1 8B are shown in Table 5. EAGLE shows the maximum throughput improvement at a batch size of 24, while EAGLE-3 shows this at 56. This part of the experiment did not use the tree structure, the maximum chain length was set to 2, and the testing dataset was MT-Bench.

Table 5: Throughput improvement under different batch sizes on A100 and LLaMA-Instruct 3.1 8B for the MT-Bench dataset, with vLLM without speculative sampling as the baseline (1.00x).

Batch size	2	4	8	16	24	32	48	56
EAGLE	1.30x	1.25x	1.21x	1.10x	1.03x	0.93x	0.82x	0.71x
EAGLE-3	1.75x	1.68x	1.58x	1.49x	1.42x	1.36x	1.21x	1.01x

5 Related Work

Many methods have been used to accelerate inference in LLMs, such as quantization (Hubara et al., 2018; Shen et al., 2020; Kim et al., 2021;

Zadeh et al., 2020; Zafrir et al., 2019) and distillation (Hinton et al., 2015). These methods generally have trade-offs, where there is a need to balance model performance with acceleration benefits.

Speculative sampling uses the target model for verification to ensure lossless acceleration. Early speculative decoding methods (Stern et al., 2018; Sun et al., 2021) accelerated generation in greedy settings, while Leviathan et al. (2023); Chen et al. (2023) introduced speculative sampling to extend the draft verification framework to non-greedy generation. Many subsequent works have improved upon speculative sampling. EAGLE (Li et al., 2024c), EAGLE-2 (Li et al., 2024b), Medusa (Cai et al., 2024), and Hydra (Ankner et al., 2024) reused the features of the target model. HASS (Zhang et al., 2024) simulates a multi-step draft process during training to mitigate the issues of training-inference inconsistency and error accumulation in EAGLE. GLIDE and CAPE (Du et al., 2024) reuse the target model’s KV cache, while methods (Hooper et al., 2023; Yang et al., 2023; Monea et al., 2023; Li et al., 2024a; Yi et al., 2024; Liu et al., 2024b; Sun et al., 2024a; Elhoushi et al., 2024; Svirschevski et al., 2024) like Draft & Verify (Zhang et al., 2023) use layer skipping or early exits to reuse parts of the target model’s parameters.

6 Conclusion

In this paper, we introduce EAGLE-3. Building upon EAGLE, EAGLE-3 incorporates two key improvements. First, it removes the feature prediction constraint, instead directly predicting draft tokens through a Training-time test. Second, it replaces the use of the target model’s top-layer features with a fusion of the target model’s lower, middle, and upper-layer features to obtain richer information. With these improvements, EAGLE-3 continues to benefit from the augmentation of training data, achieving a maximum speedup of 6.5x.

Acknowledgement

We would like to thank James Liu, Ke Bao, Yineng Zhang, Lianmin Zheng, Ying Sheng, and many others in the SGLang team for merging and evaluating EAGLE-3 in the SGLang environment.

References

Zachary Ankner, Rishab Parthasarathy, Aniruddha Nrusimha, Christopher Rinard, Jonathan Ragan-Kelley, and William Brandon. 2024. Hydra:

Sequentially-dependent draft heads for Medusa decoding. *arXiv preprint arXiv:2402.05109*.

Tianle Cai, Yuhong Li, Zhengyang Geng, Hongwu Peng, Jason D. Lee, Deming Chen, and Tri Dao. 2024. Medusa: Simple LLM inference acceleration framework with multiple decoding heads. *arXiv preprint arXiv: 2401.10774*.

Charlie Chen, Sebastian Borgeaud, Geoffrey Irving, Jean-Baptiste Lespiau, Laurent Sifre, and John Jumper. 2023. Accelerating large language model decoding with speculative sampling. *arXiv preprint arXiv:2302.01318*.

Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. 2021. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*.

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

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*.

Daya Guo DeepSeek-AI, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*.

Ning Ding, Yulin Chen, Bokai Xu, Yujia Qin, Zhi Zheng, Shengding Hu, Zhiyuan Liu, Maosong Sun, and Bowen Zhou. 2023. Enhancing chat language models by scaling high-quality instructional conversations. *arXiv preprint arXiv:2305.14233*.

Cunxiao Du, Jing Jiang, Xu Yuanchen, Jiawei Wu, Sicheng Yu, Yongqi Li, Shenggui Li, Kai Xu, Liqiang Nie, Zhaopeng Tu, et al. 2024. Glide with a cape: A low-hassle method to accelerate speculative decoding. *arXiv preprint arXiv:2402.02082*.

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The Llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.

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*.

- Yichao Fu, Peter Bailis, Ion Stoica, and Hao Zhang. 2024. Break the sequential dependency of LLM inference using lookahead decoding. *arXiv preprint arXiv:2402.02057*.
- Joao Gante. 2023. Assisted generation: a new direction toward low-latency text generation. <https://huggingface.co/blog/assisted-generation>.
- Xiangxiang Gao, Weisheng Xie, Yiwei Xiang, and Feng Ji. 2024. Falcon: Faster and parallel inference of large language models through enhanced semi-autoregressive drafting and custom-designed decoding tree. *arXiv preprint arXiv:2412.12639*.
- Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. 2025. DeepSeek-R1: Incentivizing reasoning capability in LLMs via reinforcement learning. *arXiv preprint arXiv:2501.12948*.
- Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. 2015. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*.
- Coleman Hooper, Sehoon Kim, Hiva Mohammadzadeh, Hasan Genc, Kurt Keutzer, Amir Gholami, and Sophia Shao. 2023. Speed: Speculative pipelined execution for efficient decoding. *arXiv preprint arXiv:2310.12072*.
- Itay Hubara, Matthieu Courbariaux, Daniel Soudry, Ran El-Yaniv, and Yoshua Bengio. 2018. Quantized neural networks: Training neural networks with low precision weights and activations. *journal of machine learning research*, 18(187):1–30.
- Sehoon Kim, Amir Gholami, Zhewei Yao, Michael W Mahoney, and Kurt Keutzer. 2021. I-bert: Integer-only bert quantization. In *International conference on machine learning*, pages 5506–5518. PMLR.
- Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph Gonzalez, Hao Zhang, and Ion Stoica. 2023. Efficient memory management for large language model serving with pagedattention. In *Symposium on Operating Systems Principles*, pages 611–626.
- 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.
- Minghan Li, Xilun Chen, Ari Holtzman, Beidi Chen, Jimmy Lin, Wen-tau Yih, and Xi Victoria Lin. 2024a. Nearest neighbor speculative decoding for llm generation and attribution. *arXiv preprint arXiv:2405.19325*.
- Yuhui Li, Fangyun Wei, Chao Zhang, and Hongyang Zhang. 2024b. EAGLE-2: Faster inference of language models with dynamic draft trees. In *Conference on Empirical Methods in Natural Language Processing*.
- Yuhui Li, Fangyun Wei, Chao Zhang, and Hongyang Zhang. 2024c. EAGLE: Speculative sampling requires rethinking feature uncertainty. In *International Conference on Machine Learning*.
- Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. 2024a. DeepSeek-v3 technical report. *arXiv preprint arXiv:2412.19437*.
- Fangcheng Liu, Yehui Tang, Zhenhua Liu, Yunsheng Ni, Kai Han, and Yunhe Wang. 2024b. Kangaroo: Lossless self-speculative decoding via double early exiting. *arXiv preprint arXiv:2404.18911*.
- Giovanni Monea, Armand Joulin, and Edouard Grave. 2023. Pass: Parallel speculative sampling. *arXiv preprint arXiv:2311.13581*.
- Ramesh Nallapati, Bowen Zhou, Caglar Gulcehre, Bing Xiang, et al. 2016. Abstractive text summarization using sequence-to-sequence rnns and beyond. *arXiv preprint arXiv:1602.06023*.
- Apoorv Saxena. 2023. [Prompt lookup decoding](#).
- Sheng Shen, Zhen Dong, Jiayu Ye, Linjian Ma, Zhewei Yao, Amir Gholami, Michael W Mahoney, and Kurt Keutzer. 2020. Q-bert: Hessian based ultra low precision quantization of bert. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 8815–8821.
- Mitchell Stern, Noam Shazeer, and Jakob Uszkoreit. 2018. Blockwise parallel decoding for deep autoregressive models. *Advances in Neural Information Processing Systems*, 31.
- Hanshi Sun, Zhuoming Chen, Xinyu Yang, Yuandong Tian, and Beidi Chen. 2024a. Triforce: Lossless acceleration of long sequence generation with hierarchical speculative decoding. *arXiv preprint arXiv:2404.11912*.
- Xin Sun, Tao Ge, Furu Wei, and Houfeng Wang. 2021. Instantaneous grammatical error correction with shallow aggressive decoding. *arXiv preprint arXiv:2106.04970*.
- Ziteng Sun, Jae Hun Ro, Ahmad Beirami, and Ananda Theertha Suresh. 2024b. Optimal block-level draft verification for accelerating speculative decoding. *arXiv preprint arXiv:2403.10444*.
- Ziteng Sun, Ananda Theertha Suresh, Jae Hun Ro, Ahmad Beirami, Himanshu Jain, and Felix Yu. 2024c. SpecTr: Fast speculative decoding via optimal transport. *Advances in Neural Information Processing Systems*, 36.
- Ruslan Svirschevski, Avner May, Zhuoming Chen, Beidi Chen, Zhihao Jia, and Max Ryabinin. 2024. Specexec: Massively parallel speculative decoding for interactive llm inference on consumer devices. *arXiv preprint arXiv:2406.02532*.

- Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B Hashimoto. 2023. Alpaca: A strong, replicable instruction-following model. *Stanford Center for Research on Foundation Models*. <https://crfm.stanford.edu/2023/03/13/alpaca.html>, 3(6):7.
- 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. *arXiv preprint arXiv:2302.13971*.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajwal Bhargava, Shruti Bhosale, et al. 2023b. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.
- 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. *arXiv preprint arXiv:2401.07851*.
- Seongjun Yang, Gibbeum Lee, Jaewoong Cho, Dimitris Papailiopoulos, and Kangwook Lee. 2023. Predictive pipelined decoding: A compute-latency trade-off for exact llm decoding. *arXiv preprint arXiv:2307.05908*.
- 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*.
- Ali Hadi Zadeh, Isak Edo, Omar Mohamed Awad, and Andreas Moshovos. 2020. Gobo: Quantizing attention-based nlp models for low latency and energy efficient inference. In *2020 53rd Annual IEEE/ACM International Symposium on Microarchitecture (MICRO)*, pages 811–824. IEEE.
- Ofir Zafrir, Guy Boudoukh, Peter Izsak, and Moshe Wasserblat. 2019. Q8bert: Quantized 8bit bert. In *2019 Fifth Workshop on Energy Efficient Machine Learning and Cognitive Computing-NeurIPS Edition (EMC2-NIPS)*, pages 36–39. IEEE.
- Jun Zhang, Jue Wang, Huan Li, Lidan Shou, Ke Chen, Gang Chen, and Sharad Mehrotra. 2023. Draft & verify: Lossless large language model acceleration via self-speculative decoding. *arXiv preprint arXiv:2309.08168*.
- Lefan Zhang, Xiaodan Wang, Yanhua Huang, and Ruiwen Xu. 2024. Learning harmonized representations for speculative sampling. *arXiv preprint arXiv:2408.15766*.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. 2023. Judging LLM-as-a-judge with MT-Bench and chatbot Arena. *arXiv preprint arXiv:2306.05685*.
- Lianmin Zheng, Liangsheng Yin, Zhiqiang Xie, Chuyue Livia Sun, Jeff Huang, Cody Hao Yu, Shiyi Cao, Christos Kozyrakis, Ion Stoica, Joseph E Gonzalez, et al. 2024. SGLang: Efficient execution of structured language model programs. *Advances in Neural Information Processing Systems*, 37:62557–62583.

A Implementation Details

Vanilla: We use models from the Huggingface.transformers library with the PyTorch backend and pre-allocated KV cache. Other methods also use these models as their base.

(Standard) Speculative Sampling: We use the assisted generation feature from the HuggingFace Transformers library.

PLD, Lookahead, Medusa, and Hydra: We use the default settings and the officially released weights.

EAGLE: Vicuna and LLaMA2-Chat draft models use the officially released weights, while LLaMA3-Instruct is trained using the ShareGPT dataset (consistent with Medusa and Hydra).

EAGLE-2: For the 7B (8B), 13B, and 70B original LLMs, we set the total number of draft tokens to 60, 50, and 48, respectively, with a draft tree depth of 6, and select 10 nodes during the expansion phase.

EAGLE-3: EAGLE-3’s draft model achieves a significantly higher acceptance rate, allowing us to increase the draft tree depth from 6 to 8 while keeping the number of nodes the same as in EAGLE-2.