

# Amphista: Bi-directional Multi-head Decoding for Accelerating LLM Inference

Zeping Li<sup>1\*</sup>, Xinlong Yang<sup>1,2\*</sup>, Ziheng Gao<sup>1</sup>, Ji Liu<sup>1</sup>, Guanchen Li<sup>1</sup>, Zhuang Liu<sup>1</sup>, Dong Li<sup>1</sup>, Jinzhang Peng<sup>1</sup>, Lu Tian<sup>1</sup>, Emad Barsoum<sup>1</sup>

<sup>1</sup> Advanced Micro Devices, Inc. <sup>2</sup> Peking University  
{zeping.li, emad.barsoum}@amd.com

## Abstract

Large Language Models (LLMs) inherently use autoregressive decoding, which lacks parallelism in inference and results in significantly slow inference speed. While methods such as Medusa constructs parallelized heads, they lack adequate information interaction across different prediction positions. To overcome this limitation, we introduce **Amphista**, an enhanced speculative decoding framework that builds upon Medusa. Specifically, Amphista models an *Auto-embedding Block* capable of parallel inference, incorporating bi-directional attention to enable interaction between different drafting heads. Additionally, Amphista integrates *Staged Adaptation Layers*, which ensure a seamless transition of semantic information from the target model’s autoregressive inference to the drafting heads’ non-autoregressive inference, effectively achieving paradigm shift and feature fusion. Experimental results on Vicuna models using MT-Bench and Spec-Bench demonstrate that Amphista achieves substantial acceleration while maintaining generation quality. On MT-Bench, Amphista delivers up to  $2.75\times$  speedup over vanilla autoregressive decoding and  $1.40\times$  over Medusa on Vicuna 33B in wall-clock time.

## 1 Introduction

Generative large language models (LLMs) have made remarkable advances in language processing by scaling the transformer decoder block, offering a potential pathway toward Artificial General Intelligence (AGI) (OpenAI, 2022; Touvron et al., 2023). However, the autoregressive nature of next-token prediction and the large parameter size of foundational models result in low inference efficiency, marked by high latency per token and low throughput per second during decoding.

In this context, acceleration during inference has become a burgeoning research area. Speculative de-

\*Equal contribution.

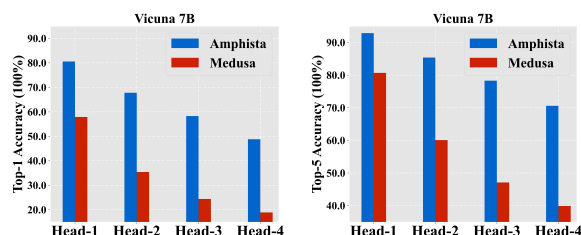


Figure 1: Top-1/5 accuracy for different heads of Medusa and Amphista. We perform testing with randomly sampled 5% ShareGPT conversation data. Amphista far outperforms Medusa in terms of head accuracy, especially for the latter two heads.

coding (Stern et al., 2018; Chen et al., 2023) uses a draft model for preliminary multi-step speculative inference and a target model to verify the speculative predictions, emerging as a very promising algorithmic strategy. Notably, by employing a rejection sampling strategy (Leviathan et al., 2023), the generation quality and accuracy of the speculate-and-verify framework are consistent with those of the target model, making speculative decoding a **lossless** acceleration framework. Medusa decoding (Cai et al., 2024) innovatively uses the target model’s last hidden states to implement a multi-heads inference framework. It is widely adopted for its efficient acceleration and simple structure.

Nonetheless, as illustrated in Figure 1, we find that the prediction accuracy of separately independent Medusa heads is relatively low, which progressively worsens and adversely impacts acceleration performance in downstream tasks. To mitigate these inaccuracies stemming from the absence of feature dependencies while maintaining parallel inference, we first introduce the Auto-embedding Block, which integrates a bi-directional self-attention mechanism (Vaswani et al., 2017). This structure not only allows earlier heads to attend to subsequent ones, but more importantly, enables backward heads to leverage information from preceding heads. This enhancement allows

drafting heads to better capture contextual information, thereby improving the acceptance rate of their predictions. Moreover, in the multi-step drafting framework, this non-autoregressive structure achieves lower drafting latency compared to an autoregressive approach.

Additionally, we identify a significant gap between the autoregressive target model and the non-autoregressive draft model in their prediction paradigms. To bridge this discrepancy and further enhance feature representations across different drafting heads, we propose the Staged Adaptation Layers. These layers serve as an intermediary module to facilitate feature integration and transformation between the target model and draft heads. Once adopted, semantically enriched features are passed through MLP activations and fed into the auto-embedding block. This enhances the bi-directional attention mechanism’s ability to fuse features across heads, ultimately boosting acceptance rates and reducing wall-clock time.

Lastly, to further align the draft model with the target model with minimal computational cost, we incorporate the sampled token from the target model’s latest prediction into the staged adaptation layers. This critically integrated information harmonizes Amphista with the target model, yielding a significant improvement in performance.

On MT-Bench, Amphista achieves up to  $2.75\times$  speedup over vanilla autoregressive decoding and  $1.40\times$  over Medusa on Vicuna 33B, as consistently high accuracy (see Figure 1). To summarize, our contributions are as follows:

- We present Amphista, a non-autoregressive and innovatively cost-efficient inference acceleration framework, built upon the foundational principles of Medusa decoding.
- We introduce the Auto-embedding Block, which enables bi-directional interaction among different heads by facilitating collaborative information exchange during the drafting phase. Additionally, the Staged Adaptation Layers are introduced to bridge the gap between autoregressive and non-autoregressive token prediction through a two-stage adaptation process. Finally, the integration of a sampled token from the target model further aligns the draft and target models with minimal computational overhead.
- We perform comprehensive evaluations on a diverse set of foundational models. The results

show that Amphista consistently outperforms Medusa in terms of both acceptance rate and speed-up, across various generation tasks.

## 2 Preliminaries

In this section, we introduce some preliminary background related to our work as follows:

**Speculative Decoding.** Speculative decoding has been successfully applied to LLM decoding algorithm recently (Leviathan et al., 2023; Chen et al., 2023). The core idea is to leverage a small, lower-quality model (draft model) together with a large, higher-quality model (target model) to accelerate token generation. Concretely, in each decoding step, the algorithm first uses the draft model to autoregressively generate a sequence of future tokens. These drafted tokens are then verified by the target model in a single forward pass. During the verification process, a certain strategy is applied to determine which tokens are accepted by the target model and which are rejected and discarded. Previous work (Leviathan et al., 2023) has theoretically and empirically demonstrated that the token output distribution of speculative decoding is consistent with the autoregressive generation of original target model, but with fewer decoding steps, thus enhancing generation efficiency.

**Medusa Decoding.** Medusa Decoding (Cai et al., 2024) represents an efficient speculative decoding algorithm based on the draft-and-verify principle, inheriting principles from blockwise parallel decoding (Stern et al., 2018). Specifically, Medusa integrates independent MLP layers, called drafting heads, with the target model to form a unified architecture. In each decoding step, the target model’s `lm_head` samples the next token, while the next-*i* MLP heads predict tokens at subsequent positions. These drafted tokens are then verified by the target model’s forward pass to decide their acceptance. By leveraging lightweight MLP layers, Medusa strikes an effective balance between computational efficiency and prediction accuracy, leading to substantial acceleration. Hydra (Ankner et al., 2024), which is a subsequent state-of-the-art optimization based on Medusa, transforms the independent MLP heads into sequentially dependent MLP heads, further enhancing the predictive accuracy.

**Tree Attention.** Tree attention (Miao et al., 2024; Cai et al., 2024) enables parallel computation of attention scores for multiple draft candidates. Medusa uses a tree causal mask, allowing each

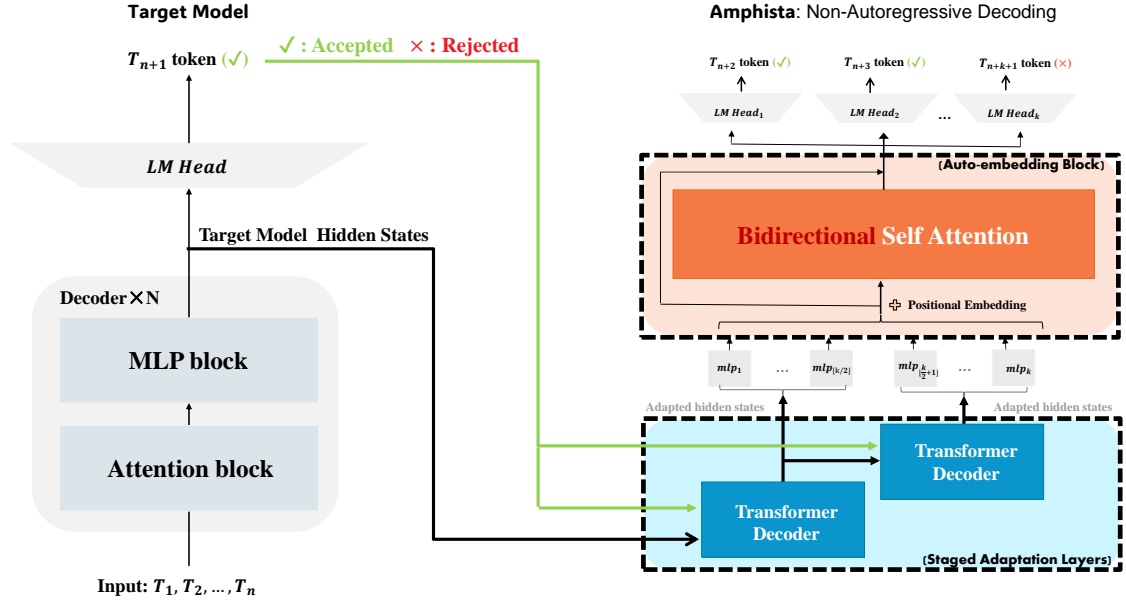


Figure 2: **The Framework of Amphista Decoding.** Our method improves Medusa in two folds: (1) We introduce staged adaptation layers, consisting of a group of causal Transformer Decoder layers built upon the target model, to adapt the target model’s hidden states and the sampled token in two stages. This module ensures that the adapted features contain richer contextual information, supporting multiple-token predictions rather than focusing solely on the immediate next-token prediction. (2) We introduce an auto-embedding block, which is a bi-directional Transformer Encoder module with positional encoding. This block allows each head to attend to others, fostering cooperative predictions and thereby enhancing the speculative accuracy during the drafting stage.

node to attend only to its ancestors, efficiently processing multiple candidate sequences simultaneously (see Appendix A.1 for details).

### 3 Amphista

The overview of Amphista is shown in Figure 2. Building its pipeline upon target model, Amphista contains two main modules: (1) Staged Adaptation Layers. They are causal Transformer Decoder layers that adapt the target model’s hidden states and sampled token embedding in two stages, each focusing on different drafting positions. This adaptation process results in hidden states that are enhanced with position-aware contextual information, improving overall prediction accuracy, especially for the latter steps. (2) Auto-embedding Block. It is a Transformer Encoder module that conducts *bi-directional* self-attention computations among the representations of different draft heads, allowing each head can be attended by the others. This facilitates collaborative prediction among these heads, thereby improving overall prediction accuracy.

#### 3.1 Staged Adaptation Layers

Figure 2 demonstrates the relevant details of our staged adaptation layers. Although target model’s

hidden states contain semantically rich information, there are still differences in the representation requirements between the target model and the draft heads. Specifically, the hidden states of the target model are trained only for predicting the next token, while draft heads need more contextual and position-aware hidden states to perform multi-step speculation. To address this problem, Medusa-2 applies LoRA (Hu et al., 2021) for joint training of the target model and draft heads, which may compromise the generality on downstream tasks. Hydra employs a single prefix layer for all positions, lacking targeted adaptation for different positions. We propose an effective adaptation method by incorporating two adaptation layers to transform and adapt the strong semantic information from the target model in stages. Specifically, given the hidden states  $h_t$  at position  $t$  from the target model’s final layer and the embedding of the token  $e_{t+1}$  sampled from  $h_t$ , we use the two adaptation layers to transform them in stages as below:

$$\begin{aligned} h_t^1 &= SAL^1(fc^1([h_t; e_{t+1}]), kv_{1:t-1}^1), \\ h_t^2 &= SAL^2(fc^2([h_t^1; e_{t+1}]), kv_{1:t-1}^2). \end{aligned} \quad (1)$$

$SAL^1$  stands for the Stage-one Adaptation Layer that adapts target model hidden states and sampled

token embedding, while  $SAL^2$  stands for the Stage-two Adaptation Layer that adapts  $SAL^1$ 's output hidden states as well as the sampled token embedding. The function  $fc^1$  and  $fc^2$  are fully connected layers employed to transform features derived from the concatenation of hidden states and token embeddings. The terms  $kv_{1:t-1}^1$  and  $kv_{1:t-1}^2$  represent the key-value caches for each adaptation layer. Subsequently, adapted hidden states  $h_t^1$  and  $h_t^2$  are fed into the first and second halves of the drafting heads respectively, ensuring that each adaptation layer focuses on adapting target model's semantic representations in specific future locations.

### 3.2 Auto-embedding Block

Figure 2 shows the detailed design of our Auto-embedding Block. Given a set of  $K$  drafting MLP heads,  $MLP_k$  head is tasked with predicting the token in the  $(t+k+1)$ -th position. Upon obtaining adapted hidden states  $h_t^1$  and  $h_t^2$  from the first and second staged adaptation layers, we first utilize the MLP layers to project them into more position-aware and semantically rich hidden states:

$$\begin{aligned} h'_k &= MLP_k(h_t^1), & k &= 1, 2, \dots, \lfloor K/2 \rfloor, \\ h'_k &= MLP_k(h_t^2), & k &= \lfloor K/2 \rfloor + 1, \dots, K, \end{aligned} \quad (2)$$

where  $MLP_i \in \mathbb{R}^{d \times d}$ , and  $d$  is the dimension of the target model hidden states. We then concatenate these  $K$  hidden states along the seq\_len dimension:

$$H' = \text{concat}([h'_1, h'_2, h'_3, \dots, h'_K]), \quad (3)$$

where  $H' \in \mathbb{R}^{K \times d}$ . In order to further enhance the relative positional information among different heads, we introduce additional positional encodings. Specifically, we introduce a learnable positional embedding  $PE \in \mathbb{R}^{K \times d}$ , and the position-encoded hidden states  $H_p$  are expressed as:

$$H_p = H' + PE. \quad (4)$$

Finally, we employ an effective and efficient bi-directional self-attention module to enable mutual awareness among the drafting heads and use additional learnable lm-head to sample the top- $k$  draft tokens in each position:

$$\text{attn}_o = \text{Self-Attention}(H_p), \quad (5)$$

$$d\_logits_k = \text{LM-Head}_k(\text{attn}_o[k]), \quad k = 1, \dots, K. \quad (6)$$

In the end, these draft tokens are organized into a draft tree and then verified by the LLM through tree attention. Unlike the independent heads in Medusa and the sequentially dependent heads in Hydra, Amphista adopts bi-directionally dependent heads. This approach enhances overall prediction accuracy while maintaining a non-autoregressive mechanism, potentially reducing the substantial computation overhead associated with sequential calculations (i.e., autoregressive manner).

### 3.3 Training Objective

Our loss function integrates two components to achieve a dual objective. First, we employ a Cross-Entropy (CE) loss between the logits of Amphista and those of the target model to align their output token distributions. Second, we utilize a language modeling (LM) loss to quantify the discrepancy between Amphista's outputs and the ground truth tokens. This approach enables Amphista not only to emulate the target model but also to assimilate predictive capabilities from the real corpus.

$$\mathcal{L}_{\text{Amphista}} = \lambda_1 \mathcal{L}_{\text{alignment}} + \lambda_2 \mathcal{L}_{\text{lm}}, \quad (7)$$

$$\mathcal{L}_{\text{alignment}} = \text{CE}(d\_logits_i, logits_{T_{t+1+i}}), \quad (8)$$

$$\mathcal{L}_{\text{lm}} = \text{CE}(d\_logits_i, gt_{t+1+i}). \quad (9)$$

Note that  $d\_logits_i$  and  $logits_{T_{t+1+i}}$  are the logits from Amphista and the target model for token  $T_{t+1+i}$ , while  $gt_{t+1+i}$  represent the ground truth labels of token  $T_{t+1+i}$ . The terms  $\lambda_1$  and  $\lambda_2$  are weighting factors for the two objectives.

## 4 Experiments

### 4.1 Experimental Settings

**Models and Baselines.** Following (Cai et al., 2024; Li et al., 2024; Ankner et al., 2024), we use Vicuna family of models (Zheng et al., 2024) as our target model. Specifically, we implement our method on Vicuna 7, 13, and 33B models with four drafting heads. As for compared baseline methods, we choose original Speculative Decoding, Lookahead (Fu et al., 2024), Medusa (Cai et al., 2024) and Hydra (Ankner et al., 2024) for comparison.

**Training and Datasets.** For the training stage, again following (Cai et al., 2024; Ankner et al., 2024), we use ShareGPT<sup>1</sup> dataset to fine-tune our proposed module while keeping target model frozen. Training is conducted using HuggingFace

<sup>1</sup>ShareGPT. 2023. [https://huggingface.co/datasets/Aeala/ShareGPT\\_Vicuna\\_unfiltered](https://huggingface.co/datasets/Aeala/ShareGPT_Vicuna_unfiltered)

Table 1: The speed-up comparison on MT-Bench and Spec-Bench between different methods under **greedy sampling** setting (Temperature = 0). We regard the speed-up of vanilla autoregressive decoding as  $1.00\times$ .

Model Size	Method	MT-Bench	Spec-Bench					Avg
			Translation	Summarization	QA	Math	RAG	
7B	Vanilla	$1.00\times$	$1.00\times$	$1.00\times$	$1.00\times$	$1.00\times$	$1.00\times$	$1.00\times$
	Spec-decoding	$1.62\times$	$1.11\times$	$1.66\times$	$1.46\times$	$1.45\times$	$1.61\times$	$1.45\times$
	Lookahead	$1.44\times$	$1.15\times$	$1.26\times$	$1.25\times$	$1.56\times$	$1.13\times$	$1.27\times$
	Medusa	$1.87\times$	$1.42\times$	$1.42\times$	$1.50\times$	$1.74\times$	$1.39\times$	$1.50\times$
	Hydra++	$2.37\times$	$1.92\times$	$1.80\times$	<b><math>1.94\times</math></b>	$2.43\times$	$2.04\times$	$2.03\times$
	<b>Amphista (ours)</b>	<b><math>2.44\times</math></b>	<b><math>1.96\times</math></b>	<b><math>2.11\times</math></b>	<b><math>1.94\times</math></b>	<b><math>2.45\times</math></b>	<b><math>2.20\times</math></b>	<b><math>2.13\times</math></b>
13B	Vanilla	$1.00\times$	$1.00\times$	$1.00\times$	$1.00\times$	$1.00\times$	$1.00\times$	$1.00\times$
	Spec-decoding	$1.66\times$	$1.17\times$	$1.75\times$	$1.44\times$	$1.59\times$	$1.73\times$	$1.53\times$
	Lookahead	$1.34\times$	$1.08\times$	$1.23\times$	$1.15\times$	$1.51\times$	$1.15\times$	$1.22\times$
	Medusa	$1.85\times$	$1.55\times$	$1.55\times$	$1.53\times$	$1.88\times$	$1.51\times$	$1.60\times$
	Hydra++	$2.34\times$	$1.75\times$	$1.85\times$	$1.85\times$	$2.31\times$	$1.86\times$	$1.92\times$
	<b>Amphista (ours)</b>	<b><math>2.49\times</math></b>	<b><math>1.88\times</math></b>	<b><math>2.14\times</math></b>	<b><math>1.88\times</math></b>	<b><math>2.41\times</math></b>	<b><math>2.04\times</math></b>	<b><math>2.07\times</math></b>
33B	Vanilla	$1.00\times$	$1.00\times$	$1.00\times$	$1.00\times$	$1.00\times$	$1.00\times$	$1.00\times$
	Spec-decoding	$1.73\times$	$1.28\times$	$1.76\times$	$1.54\times$	$1.71\times$	$1.69\times$	$1.60\times$
	Lookahead	$1.32\times$	$1.09\times$	$1.21\times$	$1.16\times$	$1.55\times$	$1.16\times$	$1.24\times$
	Medusa	$1.97\times$	$1.72\times$	$1.62\times$	$1.66\times$	$2.06\times$	$1.61\times$	$1.73\times$
	Hydra++	$2.54\times$	$1.93\times$	$2.10\times$	$2.04\times$	$2.63\times$	$2.17\times$	$2.17\times$
	<b>Amphista (ours)</b>	<b><math>2.75\times</math></b>	<b><math>2.11\times</math></b>	<b><math>2.49\times</math></b>	<b><math>2.12\times</math></b>	<b><math>2.83\times</math></b>	<b><math>2.44\times</math></b>	<b><math>2.40\times</math></b>

Trainer, which we employ with AdamW optimizer ( $\beta_1=0.9$ ,  $\beta_2=0.999$ ) and a cosine learning rate schedule with warmup strategy, the initial learning rate is set to  $1e-3$  and we train 4 epochs. At the evaluation stage, we use MT-Bench (Zheng et al., 2024) and Spec-Bench (Xia et al., 2024) as our benchmark. MT-Bench is an open source multi-turn conversation benchmark. Spec-Bench is a well-acknowledged and comprehensive benchmark designed for assessing speculative decoding methods across diverse application scenarios.

**Metrics.** Following previous speculative decoding work, we choose tokens/s and tokens/step as our main metrics. Tokens/step measures the average token length accepted per forward pass of the target LLM. Tokens/s represents the overall throughput of the acceleration algorithm, which is influenced by both the prediction accuracy of the speculator and the drafting latency of the speculator.

## 4.2 Evaluation of Amphista

Amphista builds on Medusa to support parallel decoding, distinctly diverging from auto-regression drafting methods. Thus, the representative work of parallel drafting (i.e., Lookahead), and the state-of-the-art work based on Medusa (i.e., Hydra), have been chosen as a competitive baseline method for comparison. Specifically, Hydra’s best-performing model (i.e., Hydra++) is used for fair evaluation and vicuna-68m (Yang et al., 2024) is used as

draft model for the vanilla speculative decoding method. We conduct all the experiments on A100 40G GPUs, and all the experimental settings are kept the same for fair comparison.

Table 1 and Table 2 present the speed-up metrics compared on MT-Bench and Spec-Bench under greedy and random sampling settings (see A.2 for more experiment results). Overall, Amphista demonstrates significant performance superiority over Medusa and surpasses Hydra’s best results by a considerable margin across a variety of generation tasks, and also greatly exceeding the speed-up achieved by vanilla speculative decoding. In detail, Amphista achieves a  $2.44\times$  -  $2.75\times$  speed-up on MT-Bench and  $2.13\times$  -  $2.40\times$  speed-up on Spec-Bench under greedy decoding setting. Similarly, under random sampling setting, Amphista achieves a  $2.37\times$  -  $2.85\times$  speed-up and  $1.99\times$  -  $2.43\times$  speed-up on MT-Bench and Spec-Bench with different target model sizes. These robust results demonstrate that enhancing non-autoregressive drafting can surpass autoregressive drafting in terms of speed-up, highlighting the efficiency of our Amphista architecture. During the drafting stage, all computations in non-autoregressive modeling (i.e., Amphista) can be processed in parallel, better leveraging the parallel computing capabilities of modern GPU accelerators. This leads to a more optimal trade-off between drafting acceptance rate and drafting latency.

Table 2: The speed-up comparison on MT-Bench and Spec-bench between different methods under **random sampling** setting (Temperature = 0.7). We regard the speed-up of vanilla autoregressive decoding as 1.00 $\times$ .

Model Size	Method	MT-Bench	Spec-Bench					Avg
			Translation	Summarization	QA	Math	RAG	
7B	Vanilla	1.00 $\times$	1.00 $\times$	1.00 $\times$	1.00 $\times$	1.00 $\times$	1.00 $\times$	1.00 $\times$
	Spec-decoding	1.39 $\times$	1.02 $\times$	1.41 $\times$	1.24 $\times$	1.32 $\times$	1.43 $\times$	1.28 $\times$
	Lookahead	1.28 $\times$	1.05 $\times$	1.21 $\times$	1.12 $\times$	1.25 $\times$	1.14 $\times$	1.16 $\times$
	Medusa	1.86 $\times$	1.51 $\times$	1.47 $\times$	1.57 $\times$	1.89 $\times$	1.43 $\times$	1.57 $\times$
	Hydra++	2.35 $\times$	<b>1.81</b> $\times$	1.81 $\times$	<b>1.97</b> $\times$	2.41 $\times$	1.74 $\times$	1.95 $\times$
	<b>Amphista (ours)</b>	<b>2.37</b> $\times$	<b>1.81</b> $\times$	<b>1.92</b> $\times$	1.96 $\times$	<b>2.43</b> $\times$	<b>1.79</b> $\times$	<b>1.99</b> $\times$
13B	Vanilla	1.00 $\times$	1.00 $\times$	1.00 $\times$	1.00 $\times$	1.00 $\times$	1.00 $\times$	1.00 $\times$
	Spec-decoding	1.52 $\times$	1.08 $\times$	1.57 $\times$	1.33 $\times$	1.42 $\times$	1.46 $\times$	1.37 $\times$
	Lookahead	1.30 $\times$	1.07 $\times$	1.19 $\times$	1.15 $\times$	1.38 $\times$	1.14 $\times$	1.19 $\times$
	Medusa	2.01 $\times$	1.65 $\times$	1.62 $\times$	1.71 $\times$	2.01 $\times$	1.57 $\times$	1.71 $\times$
	Hydra++	2.57 $\times$	1.90 $\times$	1.99 $\times$	2.12 $\times$	2.56 $\times$	2.04 $\times$	2.12 $\times$
	<b>Amphista (ours)</b>	<b>2.65</b> $\times$	<b>1.93</b> $\times$	<b>2.16</b> $\times$	<b>2.17</b> $\times$	<b>2.64</b> $\times$	<b>2.15</b> $\times$	<b>2.22</b> $\times$
33B	Vanilla	1.00 $\times$	1.00 $\times$	1.00 $\times$	1.00 $\times$	1.00 $\times$	1.00 $\times$	1.00 $\times$
	Spec-decoding	1.58 $\times$	1.21 $\times$	1.62 $\times$	1.48 $\times$	1.59 $\times$	1.54 $\times$	1.48 $\times$
	Lookahead	1.29 $\times$	1.04 $\times$	1.18 $\times$	1.15 $\times$	1.52 $\times$	1.14 $\times$	1.21 $\times$
	Medusa	2.06 $\times$	1.71 $\times$	1.79 $\times$	1.76 $\times$	2.10 $\times$	1.79 $\times$	1.83 $\times$
	Hydra++	2.74 $\times$	2.01 $\times$	2.24 $\times$	2.24 $\times$	2.82 $\times$	2.26 $\times$	2.31 $\times$
	<b>Amphista (ours)</b>	<b>2.85</b> $\times$	<b>2.05</b> $\times$	<b>2.51</b> $\times$	<b>2.29</b> $\times$	<b>2.90</b> $\times$	<b>2.39</b> $\times$	<b>2.43</b> $\times$

Moreover, Amphista exhibits a discernible upward trend in speed-up when employed on larger target models. This can be attributed to Amphista’s cost-efficient non-autoregressive modeling and effective transformation of semantic information from the target model. Amphista allows for appropriate increases in accepted token length without introducing excessive additional inference costs. For more exploration on the performance potential of Amphista, please refer to A.2.4. For more exploration on the parameter complexity optimization, please refer to A.2.5.

Last but not least, we further provide the actual throughput of different methods on MT-Bench with a batch size of 1. As depicted in Figure 3, Amphista achieves an actual throughput of approximately 120 tokens/s with a 7B target model and about 80 tokens/s with a 13B target model under both temperature settings. This performance surpasses that of Medusa and Hydra, underscoring Amphista’s advantages in practical deployment.

### 4.3 Generation Quality of Amphista

We perform evaluation on XSUM (Narayan et al., 2018) and CNN/DM (See et al., 2017) to validate the generation quality of our Amphista (more results can be found in appendix A.2.1). From the ROUGE-1/2/L scores (Lin, 2004) in Table 3, we can find that Amphista can reserve the output distribution quality while achieving 2.10 $\times$ -2.31 $\times$  speed-

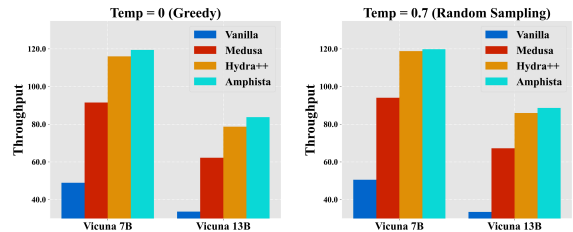


Figure 3: Throughput (tokens/s) on MT-Bench with different target model sizes and temperatures.

Table 3: Results on CNN/DM and XSUM with different temperatures, AR means Auto-Regressive decoding.

Benchmark	Temp	Method	ROUGE-1	ROUGE-2	ROUGE-L	Speed-up
CNN/DM	0.0	AR	18.74	8.44	12.59	1.00 $\times$
		Amphista	18.70	8.44	12.59	2.15 $\times$
	0.7	AR	17.92	7.65	11.93	1.00 $\times$
		Amphista	17.91	7.65	11.92	2.31 $\times$
XSUM	0.0	AR	17.32	5.05	12.16	1.00 $\times$
		Amphista	17.30	5.05	12.15	2.25 $\times$
	0.7	AR	15.99	4.44	11.42	1.00 $\times$
		Amphista	15.96	4.43	11.40	2.10 $\times$

up compared with vanilla auto-regressive decoding.

### 4.4 Multi-Batching Exploration

In this section, we evaluate the speed-up of Amphista in multi-batch scenarios (batch size > 1). For varying sentence lengths within a batch, we use padding to align them and always track the position of the last valid token for each sentence. The experimental results, presented in Table 6, are based on randomly sampled prompts from MT-Bench to generate various batch sizes. Generally, as batch

Table 4: Ablation experiments of different model variants on MT-Bench and Spec-Bench, with the target model being Vicuna 7B and the evaluation metric being **speed-up**. Medusa can be considered as Amphista w/o any added modules, and Hydra can be seen as Medusa w/ sequential dependency heads.

Method Variants	MT-Bench	Spec-Bench					Avg
		Translation	Summary	QA	Math	RAG	
Medusa	1.86×	1.51×	1.47×	1.57×	1.89×	1.43×	1.57×
Hydra++	2.37×	1.92×	1.80×	1.94×	2.43×	2.04×	2.03×
Amphista w/o Auto-embedding	2.30×	1.82×	2.00×	1.81×	2.25×	1.99×	1.97×
Amphista w/o Position-Encoding	2.42×	1.96×	2.08×	1.92×	2.42×	2.18×	2.11×
Amphista w/o Staged-Adaptation	2.14×	1.85×	1.75×	1.78×	2.10×	1.91×	1.88×
Amphista w/ One-Adaptation-Layer	2.31×	1.90×	1.99×	1.83×	2.35×	2.14×	2.04×
Amphista w/o Sampled-Token	2.25×	1.88×	1.80×	1.81×	2.26×	2.01×	1.95×
<b>Amphista (ours)</b>	2.44×	1.96×	2.11×	1.94×	2.45×	2.20×	<b>2.13×</b>

Table 5: Ablation experiments of different model variants on MT-Bench and Spec-Bench, with the target model being Vicuna 7B and evaluation metric being **average accepted length**. Medusa can be considered as Amphista w/o any added modules, and Hydra can be seen as Medusa w/ sequential dependency heads.

Method Variants	MT-Bench	Spec-Bench					Avg
		Translation	Summary	QA	Math	RAG	
Medusa	2.52	2.12	2.01	2.05	2.48	2.09	2.15
Hydra++	3.58	2.80	2.70	2.91	3.61	2.90	<b>2.98</b>
Amphista w/o Auto-embedding	3.16	2.41	2.66	2.40	3.11	2.49	2.60
Amphista w/o Position-Encoding	3.47	2.61	2.90	2.78	3.47	2.91	2.93
Amphista w/o Staged-Adaptation	2.91	2.42	2.24	2.30	2.85	2.38	2.43
Amphista w/ One-Adaptation-Layer	3.36	2.49	2.68	2.71	3.37	2.75	2.80
Amphista w/o Sampled-Token	3.11	2.43	2.48	2.45	3.15	2.55	2.61
<b>Amphista (ours)</b>	3.50	2.62	3.01	2.80	3.50	2.96	<b>2.98</b>

size increases, the GPU’s idle computational resources gradually decrease, resulting in a reduced speed-up. Additionally, despite the additional computational overhead from different multi-batching strategies, Amphista consistently achieves around 2× speed-up using the simplest padding method, demonstrating its acceleration advantage in multi-batch settings.

Table 6: Speed-up evaluation of Amphista on MT-Bench with batch size > 1.

Batch Size	1	2	4	6	8
Speed-up	2.32×	2.30×	2.23×	2.11×	2.06×

#### 4.5 Ablation Study

Diverging from other approaches based on speculative sampling and Medusa, Amphista’s main insight lies in adapting transformation through Staged Adaptation Layers and enhancing integration via the non-autoregressive Auto-embedding Block. These approaches strengthen semantic information derived from the target model. In this section, we conduct comprehensive ablation ex-

periments based on the vicuna 7B model to validate the effectiveness of each proposed module in our Amphista. Specifically, we conduct five model variants as follows: (1) **Amphista w/o Auto-embedding** which means removing the Auto-embedding Block. (2) **Amphista w/o Position-Encoding** which means removing the additional position embedding matrix in Auto-embedding Block. (3) **Amphista w/o Staged-Adaptation** which means removing staged adaptation layers. (4) **Amphista w/ One-Adaptation-Layer** which means using only one adaptation layer for all the drafting heads. (5) **Amphista w/o Sampled-Token** which means removing sampled token during adaptation process. The experimental results are presented in Table 4, 5. From these comparative results, some observations can be found as follows:

- **Amphista w/o Auto-embedding** exhibits an approximate 5%-8% decrease in speed-up performance and about a 10%-12% reduction in average accepted length. This highlights the effectiveness of the Auto-embedding Block in mitigating inaccuracies deriving from the independent speculation of Medusa heads, and demonstrating the

efficiency of non-autoregressive drafting computations. Additionally, **Amphista w/o Position-Encoding** exhibits a slight performance decline, with an approximate 2% decrease in inference speed-up, suggesting that position encoding provides additional benefits.

- **Amphista w/o Staged-Adaptation** leads to a more significant decline in speed-up (14%) and average accepted length (16%). This emphasizes the importance of bridging the feature gap between the target model and drafting heads, and further underscores the critical role of the staged adaptation layer in enhancing the auto-embedding block. Additionally, it is noteworthy that **Amphista w/ One-Adaptation-Layer** utilizes only a single adaptation layer for all drafting positions. In contrast to staged adaptation, this approach poses greater challenges to the adaptation process, resulting in some performance degradation, thereby validating the rationale behind our staged adaptation design.
- **Amphista w/o Sampled-Token** also causes an approximate 8% performance decline. Unlike previous works (e.g., Hydra), we do not use the sampled token directly for the next step of prediction. Instead, we adapt it along with the target model’s hidden states. This not only indicates that the sampled token, in addition to target model hidden states, contains important semantic information, but also demonstrates the effectiveness of our staged adaptation approach.
- Thanks to the autoregressive characteristics and the substantial number of parameters in the MLP layers, Hydra exhibits great performance in average token length. However, the computational overhead of auto-regressive methods is huge, resulting in significant reductions when translated into final speed-up. In contrast, Amphista achieves a comparable average token length to Hydra, and due to the parallelism and efficiency of its non-autoregressive computations, it ultimately attains a more favorable overall trade-off.

## 5 Related Work

Increasing techniques have been proposed to enhance the inference speed of large language models (LLMs), covering aspects of system hardware, model architecture, and decoding algorithms. A significant branch of these techniques is **Model Compression**, which includes methods such as

model quantization (Yao et al., 2023; Dettmers et al., 2024; Liu et al., 2023a; Ma et al., 2024), pruning (Belcak and Wattenhofer, 2023; Liu et al., 2023b; Zhong et al., 2024), and distillation (Zhou et al., 2024; Sun et al., 2024; Touvron et al., 2021). Additionally, techniques like kv-cache (Ge et al., 2023; Kwon et al., 2023), flash-attention (Dao et al., 2022), and early exiting (Bae et al., 2023; Elhoushi et al., 2024; Liu et al., 2024a) have also significantly reduced inference overhead. Another important line is **Speculative Decoding**, which our work is based on. It can be broadly categorized into two types. The first treats the target model and draft model separately and independently, involving the use of a small language model (Kim et al., 2024; Leviathan et al., 2023; Liu et al., 2024b; Monea et al., 2023; Chen et al., 2024; Du et al., 2024), external database, or n-grams pool (He et al., 2024; Fu et al., 2024; Kou et al., 2024; Ou et al., 2024) to generate candidate token sequences or token trees (Miao et al., 2024), which the LLM then verifies. The second type views the draft model as a dependent approximation of the target model, deriving the draft model directly from the target model or building additional modules on top of the target model for drafting (Stern et al., 2018; Zhang et al., 2023, 2024; Li et al., 2024; Cai et al., 2024; Kim et al.; Xiao et al., 2024; Ankner et al., 2024). Unlike these approaches, we propose a novel method using an auto-embedding block combined with staged adaptation layers to further enhance acceleration.

## 6 Conclusion

We propose Amphista, an efficient non-autoregressive speculative decoding framework that accelerates inference through parallel decoding and improves alignment between target and draft models via feature adaptation. Amphista integrates two core components: the Auto-embedding Block, leveraging bi-directional self-attention for collaborative speculation among drafting heads, and Staged Adaptation Layers, transforming target model semantics for multi-step predictions. Additionally, Amphista exploits sampled tokens to further optimize alignment. Extensive experiments confirm the superiority of Amphista, showcasing the promise of non-autoregressive methods in speculative decoding.



## Limitations

While we have found and adhered to using bi-directional self-attention for non-autoregressive modeling as an efficient inference structure, we have not yet fully explored the optimal structure of the Auto-embedding Block module. Specifically, this includes experimenting with different intermediate sizes (i.e., the hidden dimensions used in self-attention computations) and increasing the number of self-attention layers within the auto-embedding block to enhance its modeling depth (see A.2.4). Both of these structural optimizations could potentially improve Amphista’s acceleration performance within the current framework. Additionally, this work primarily focuses on scenarios where the batch size is equal to one, with limited optimization for larger batch sizes. We leave these areas as our future work and also hope that researchers interested in non-autoregressive inference acceleration will build upon this foundation.

## Acknowledgement

We acknowledge the helpful discussions from Kolorin Yan, Fuwei Yang, Ethan Yang, Xiandong Zhao, Mahdi Kamani, and Vikram Appia during the writing process of this work.

## References

- Zachary Ankner, Rishab Parthasarathy, Aniruddha Nrusimha, Christopher Rinard, Jonathan Ragan-Kelley, and William Brandon. 2024. [Hydra: Sequentially-dependent draft heads for medusa decoding](#). *Preprint*, arXiv:2402.05109.
- Sangmin Bae, Jongwoo Ko, Hwanjun Song, and Se-Young Yun. 2023. [Fast and robust early-exiting framework for autoregressive language models with synchronized parallel decoding](#). pages 5910–5924, Singapore.
- Peter Belcak and Roger Wattenhofer. 2023. [Exponentially faster language modelling](#). *arXiv preprint arXiv:2311.10770*.
- 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](#). *Preprint*, arXiv:2401.10774.
- Charlie Chen, Sebastian Borgeaud, Geoffrey Irving, Jean-Baptiste Lespiau, Laurent Sifre, and John Jumper. 2023. [Accelerating large language model decoding with speculative sampling](#). *Preprint*, arXiv:2302.01318.
- Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian, Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fotios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebgen Guss, Alex Nichol, Alex Paino, Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders, Christopher Hesse, Andrew N. Carr, Jan Leike, Josh Achiam, Vedant Misra, Evan Morikawa, Alec Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. 2021. [Evaluating large language models trained on code](#).
- Zhuoming Chen, Avner May, Ruslan Svirschevski, Yuhsun Huang, Max Ryabinin, Zhihao Jia, and Beidi Chen. 2024. [Sequoia: Scalable, robust, and hardware-aware speculative decoding](#). *Preprint*, arXiv:2402.12374.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. [Training verifiers to solve math word problems](#). *arXiv preprint arXiv:2110.14168*.
- Tri Dao, Daniel Y. Fu, Stefano Ermon, Atri Rudra, and Christopher Ré. 2022. [Flashattention: Fast and memory-efficient exact attention with io-awareness](#). *Preprint*, arXiv:2205.14135.
- Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. 2024. [Qlora: Efficient finetuning of quantized llms](#). *Advances in Neural Information Processing Systems*, 36.
- Cunxiao Du, Jing Jiang, Xu Yuanchen, Jiawei Wu, Sicheng Yu, Yongqi Li, Shenggui Li, Kai Xu, Liqiang Nie, Zhaopeng Tu, and Yang You. 2024. [Glide with a cape: A low-hassle method to accelerate speculative decoding](#). *Preprint*, arXiv:2402.02082.
- Mostafa Elhoushi, Akshat Shrivastava, Diana Liskovich, Basil Hosmer, Bram Wasti, Liangzhen Lai, Anas Mahmoud, Bilge Acun, Saurabh Agarwal, Ahmed Roman, Ahmed A Aly, Beidi Chen, and Carole-Jean Wu. 2024. [Layerskip: Enabling early exit inference and self-speculative decoding](#). *Preprint*, arXiv:2404.16710.
- Yichao Fu, Peter Bailis, Ion Stoica, and Hao Zhang. 2024. [Break the sequential dependency of llm inference using lookahead decoding](#). *Preprint*, arXiv:2402.02057.
- Suyu Ge, Yunan Zhang, Liyuan Liu, Minjia Zhang, Jiawei Han, and Jianfeng Gao. 2023. [Model tells you](#)

- what to discard: Adaptive kv cache compression for llms. *arXiv preprint arXiv:2310.01801*.
- Zhenyu He, Zexuan Zhong, Tianle Cai, Jason D. Lee, and Di He. 2024. [Rest: Retrieval-based speculative decoding](#). *Preprint*, arXiv:2311.08252.
- Edward J Hu, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al. 2021. Lora: Low-rank adaptation of large language models. In *International Conference on Learning Representations*.
- Sehoon Kim, Karttikeya Mangalam, Suhong Moon, Jitendra Malik, Michael W Mahoney, Amir Gholami, and Kurt Keutzer. 2024. Speculative decoding with big little decoder. *Advances in Neural Information Processing Systems*, 36.
- Taehyeon Kim, Ananda Theertha Suresh, Kishore A Papineni, Michael Riley, Sanjiv Kumar, and Adrian Benton. Exploring and improving drafts in blockwise parallel decoding. In *Workshop on Efficient Systems for Foundation Models II@ ICML2024*.
- Siqi Kou, Lanxiang Hu, Zhezhi He, Zhijie Deng, and Hao Zhang. 2024. [Cllms: Consistency large language models](#). *Preprint*, arXiv:2403.00835.
- 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 *Proceedings of the 29th 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.
- Yuhui Li, Fangyun Wei, Chao Zhang, and Hongyang Zhang. 2024. [Eagle: Speculative sampling requires rethinking feature uncertainty](#). *Preprint*, arXiv:2401.15077.
- Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*, pages 74–81.
- Fangcheng Liu, Yehui Tang, Zhenhua Liu, Yunsheng Ni, Kai Han, and Yunhe Wang. 2024a. [Kangaroo: Lossless self-speculative decoding via double early exiting](#). *Preprint*, arXiv:2404.18911.
- Xiaoxuan Liu, Lanxiang Hu, Peter Bailis, Alvin Cheung, Zhijie Deng, Ion Stoica, and Hao Zhang. 2024b. [Online speculative decoding](#). *Preprint*, arXiv:2310.07177.
- Zechun Liu, Barlas Oguz, Changsheng Zhao, Ernie Chang, Pierre Stock, Yashar Mehdad, Yangyang Shi, Raghuraman Krishnamoorthi, and Vikas Chandra. 2023a. [Llm-qat: Data-free quantization aware training for large language models](#). *Preprint*, arXiv:2305.17888.
- Zichang Liu, Jue Wang, Tri Dao, Tianyi Zhou, Binhang Yuan, Zhao Song, Anshumali Shrivastava, Ce Zhang, Yuandong Tian, Christopher Re, et al. 2023b. Deja vu: Contextual sparsity for efficient llms at inference time. In *International Conference on Machine Learning*, pages 22137–22176. PMLR.
- Shuming Ma, Hongyu Wang, Lingxiao Ma, Lei Wang, Wenhui Wang, Shaohan Huang, Li Dong, Ruiping Wang, Jilong Xue, and Furu Wei. 2024. [The era of 1-bit llms: All large language models are in 1.58 bits](#). *Preprint*, arXiv:2402.17764.
- 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 *Proceedings of the 29th ACM International Conference on Architectural Support for Programming Languages and Operating Systems, Volume 3*, pages 932–949.
- Giovanni Monea, Armand Joulin, and Edouard Grave. 2023. Pass: Parallel speculative sampling. *arXiv preprint arXiv:2311.13581*.
- Shashi Narayan, Shay B. Cohen, and Mirella Lapata. 2018. Don’t give me the details, just the summary! topic-aware convolutional neural networks for extreme summarization. *ArXiv*, abs/1808.08745.
- OpenAI. 2022. Chatgpt: Chatgpt: Optimizing language models for dialogue.
- Jie Ou, Yueming Chen, and Wenhong Tian. 2024. Lossless acceleration of large language model via adaptive n-gram parallel decoding. *arXiv preprint arXiv:2404.08698*.
- Abigail See, Peter J. Liu, and Christopher D. Manning. 2017. [Get to the point: Summarization with pointer-generator networks](#). In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1073–1083, Vancouver, Canada. Association for Computational Linguistics.
- Mitchell Stern, Noam Shazeer, and Jakob Uszkoreit. 2018. Blockwise parallel decoding for deep autoregressive models. *Advances in Neural Information Processing Systems*, 31.
- Ziteng Sun, Ananda Theertha Suresh, Jae Hun Ro, Ahmad Beirami, Himanshu Jain, and Felix Yu. 2024. Spectr: Fast speculative decoding via optimal transport. *Advances in Neural Information Processing Systems*, 36.
- Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou. 2021. Training data-efficient image transformers & distillation through attention. In *International conference on machine learning*, pages 10347–10357. PMLR.

- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems*, 30.
- 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](#). *Preprint*, arXiv:2401.07851.
- Bin Xiao, Chunan Shi, Xiaonan Nie, Fan Yang, Xiangwei Deng, Lei Su, Weipeng Chen, and Bin Cui. 2024. Clover: Regressive lightweight speculative decoding with sequential knowledge. *arXiv preprint arXiv:2405.00263*.
- Sen Yang, Shujian Huang, Xinyu Dai, and Jiajun Chen. 2024. [Multi-candidate speculative decoding](#). *Preprint*, arXiv:2401.06706.
- Zhewei Yao, Cheng Li, Xiaoxia Wu, Stephen Youn, and Yuxiong He. 2023. A comprehensive study on post-training quantization for large language models. *arXiv preprint arXiv:2303.08302*.
- Aonan Zhang, Chong Wang, Yi Wang, Xuanyu Zhang, and Yunfei Cheng. 2024. [Recurrent drafter for fast speculative decoding in large language models](#). *Preprint*, arXiv:2403.09919.
- 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*.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. 2024. Judging llm-as-a-judge with mt-bench and chatbot arena. *Advances in Neural Information Processing Systems*, 36.
- Shuzhang Zhong, Zebin Yang, Meng Li, Ruihao Gong, Runsheng Wang, and Ru Huang. 2024. Propd: Dynamic token tree pruning and generation for llm parallel decoding. *arXiv preprint arXiv:2402.13485*.
- Yongchao Zhou, Kaifeng Lyu, Ankit Singh Rawat, Aditya Krishna Menon, Afshin Rostamizadeh, Sanjiv Kumar, Jean-François Kagy, and Rishabh Agarwal. 2024. [Distillspec: Improving speculative decoding via knowledge distillation](#). *Preprint*, arXiv:2310.08461.

## A Appendix

### A.1 Draft Tree

For a fully fair comparison, we adopt the same draft tree structure as Medusa and Hydra. As shown in Figure 4, this tree is a sparse structure with a depth of 4, representing four drafting heads, and includes a total of 64 nodes, including the root node (the token sampled in the final step of the target model). Each layer’s nodes represent the tokens obtained by top\_k sampling from the corresponding drafting head. The entire tree is constructed using an auxiliary dataset by maximizing the acceptance probability of the whole tree (Cai et al., 2024). Moreover, a specially designed tree mask is used to correctly compute attention scores while simultaneously handling multiple paths, as described in Figure 5.

However, in some cases, due to the lack of redundant computational power (such as in high-throughput inference service scenarios) or parallel accelerators, an excessive number of tree nodes may lead to significant computation overhead, thereby affecting the acceleration efficiency of the algorithm. Consequently, we configure varying numbers of draft tree nodes without changing the tree depth for more comprehensive comparison, and the experimental results are shown in Table 7. From these results we observe that as the number of tree nodes decreases, the width of the tree reduces, leading to a decrease in speed-up for all compared methods. However, the decline is slightly less pronounced for Amphista, owing to its higher head accuracy. Furthermore, across various tree node configurations, we consistently achieve optimal performance, demonstrating the advantages of our algorithm in practical deployment and low-resource scenarios.

Table 7: Speed-up comparison on MT-Bench with varying number of draft tree nodes.

Method	Node = 22	Node = 35	Node = 45	Node = 64
Medusa	1.71×	1.80×	1.87×	1.87×
Hydra++	2.17×	2.26×	2.28×	2.37×
Amphista	<b>2.29×</b>	<b>2.37×</b>	<b>2.42×</b>	<b>2.44×</b>

### A.2 Additional Experiments Results

#### A.2.1 Evaluation on XSUM and CNN/DM

We use XSUM (Narayan et al., 2018) and CNN/DM (See et al., 2017) for evaluating the generation quality of Amphista, the target model

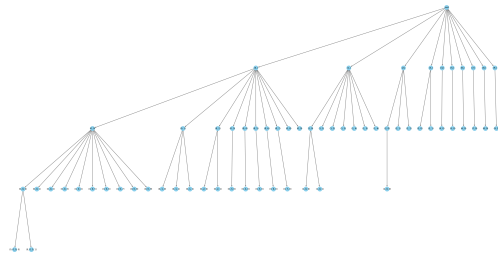


Figure 4: Draft tree used in Medusa, Hydra and our Amphista.

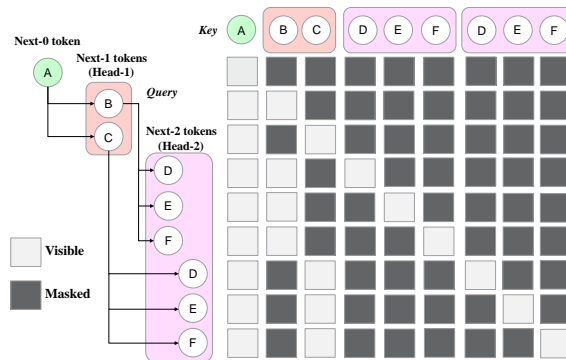


Figure 5: **An Illustration of Tree Attention.** Assuming Medusa has only 2 heads, where head-1 generates the top-2 tokens and head-2 generates the top-3 tokens, resulting in 6 candidate sequences (e.g., ABD). Additionally, a special tree mask is designed to ensure causal relationships among the top-k nodes of each head.

is vicuna 7B. Specifically, we perform zero-shot evaluation and the input prompt template is 'Article:+'+'Original Text'+'Summary:'. Additionally, for input prompts exceeding a length of 2048, we perform truncation to meet the target model’s input requirements.

Table 8: The speed-up metric comparison on Humaneval and GSM8K between different methods under greedy setting. The target model is vicuna 7B and 13B, and we regard the speed-up of vanilla auto-regressive decoding as 1.00×

Model Size	Benchmark	Vinilla AR	Medusa	Hydra++	Amphista
7B	Humaneval	1.00×	2.40×	2.76×	<b>3.02×</b>
	GSM8K	1.00×	1.87×	2.14×	<b>2.32×</b>
13B	Humaneval	1.00×	2.11×	2.75×	<b>3.00×</b>
	GSM8K	1.00×	1.98×	2.39×	<b>2.68×</b>

#### A.2.2 Code Generation and Math Reasoning

In this section, we provide more experimental results on code generation and math reasoning. we choose public Humaneval (Chen et al., 2021) and GSM8k (Cobbe et al., 2021) benchmark for evalua-

Table 9: The speed-up and average accepted length metric comparison with the target model being vicuna 7B. We regard the speed-up of vanilla auto-regressive decoding as 1.00 $\times$ .

Metric	Method	MT-Bench	Spec-Bench					Avg
			Translation	Summarization	QA	Math	RAG	
Speed-up	Vanilla	1.00 $\times$	1.00 $\times$	1.00 $\times$	1.00 $\times$	1.00 $\times$	1.00 $\times$	1.00 $\times$
	Hydra++	2.37 $\times$	1.92 $\times$	1.80 $\times$	1.94 $\times$	2.43 $\times$	2.04 $\times$	2.03 $\times$
	EAGLE	2.58 $\times$	1.94 $\times$	2.21 $\times$	2.02 $\times$	2.57 $\times$	2.30 $\times$	2.21 $\times$
	Amphista	2.44 $\times$	1.96 $\times$	2.11 $\times$	1.94 $\times$	2.45 $\times$	2.20 $\times$	2.13 $\times$
	Amphista- $\alpha$	<b>2.63<math>\times</math></b>	<b>2.09<math>\times</math></b>	<b>2.23<math>\times</math></b>	<b>2.06<math>\times</math></b>	<b>2.61<math>\times</math></b>	<b>2.34<math>\times</math></b>	<b>2.27<math>\times</math></b>
Average Accepted Length	Vanilla	1.00	1.00	1.00	1.00	1.00	1.00	1.00
	Hydra++	3.58	2.80	2.70	2.91	3.61	2.90	2.98
	EAGLE	3.84	2.92	3.32	3.14	3.93	3.31	<b>3.32</b>
	Amphista	3.50	2.62	3.01	2.80	3.50	2.96	2.98
	Amphista- $\alpha$	3.58	2.70	3.14	2.90	3.62	3.08	3.09

tion, and the target model is vicuna 7B and vicuna 13B. According to the results in Table 8, we can observe that due to the universal template and notation of code generation and mathematical reasoning, almost all compared methods achieve a higher speed-up. Furthermore, Amphista algorithm consistently attains optimal performance, demonstrating the superiority of our approach.

### A.2.3 Exploring The Optimal Number of Adaptation Stages.

In this section, we analyze the impact of varying the number of adaptation stages on the final performance of Amphista (see 10). The staged adaptation approach was introduced to facilitate the progressive transformation of features, and our ablation study has demonstrated that employing two adaptation stages outperforms a single-stage design. This indicates that using multiple adaptation layers to adapt to different positions is more effective than relying on a single layer for all positions.

Although it is theoretically possible to further increase the number of adaptation stages, our experiments reveal two critical limitations. First, as the number of layers increases, the parameter count increases significantly, leading to increased model complexity. Second, the additional adaptation layers provide diminishing returns, with only minimal improvements in mean accept length, while introducing higher inference time overhead and reducing the overall speed-up metric. Given these trade-offs, we identified the two-layer approach as the optimal balance, effectively enhancing performance without unnecessary computational burden.

Table 10: Experimental results for different numbers of staged adaptation layers.

Stages	1	2	3	4
Accept Length	3.36	3.50	3.46	3.51
Speed-up	2.31 $\times$	2.44 $\times$	2.23 $\times$	2.24 $\times$

### A.2.4 Exploring The Potential of Amphista

In this section, we conduct a preliminary exploration of Amphista’s scaling ability to demonstrate its potential for performance enhancement. By leveraging the efficiency of non-autoregressive modeling, we increase the number of auto-embedding blocks, which are essential modules within Amphista, while maintaining parallel inference. This approach yields remarkable results, detailed in Table 9. Specifically, we employ two layers of self-attention in the auto-embedding module, renaming our method as Amphista- $\alpha$ . This adjustment leads to an average accepted length increase of approximately 0.1-0.2 tokens and a notable 5%-8% improvement in overall speed-up, highlighting Amphista’s performance growth potential. We anticipate this to be a highly promising and potent attribute of Amphista.

Table 11: Experiment results of LoRA-like lm heads optimization. Note that we consider the speed-up of full rank lm head as 1.00 $\times$ .

Benchmark	rank=4096 (full)	rank=64	rank=128	rank=256	rank=512
MT-bench	1.00 $\times$	0.98 $\times$	1.01 $\times$	<b>1.02<math>\times</math></b>	0.99 $\times$
Spec-bench	1.00 $\times$	1.00 $\times$	1.00 $\times$	<b>1.01<math>\times</math></b>	0.98 $\times$

### A.2.5 Parameter Complexity Optimization of Amphista

In this part, we propose LoRA-like drafting lm heads to further optimize the original learnable lm

Input Prompt: Write a code to calculate Fibonacci numbers.

```

Here is a Python code to calculate Fibonacci numbers:
python
def fibonacci(n):
    if n <= 1:
        return n
    else:
        return fibonacci(n-1) + fibonacci(n-2)

for i in range(10):
    print(fibonacci(i))
...
This code defines a function 'fibonacci' that takes an integer 'n' as input and returns the 'n'-th Fibonacci number. The function uses recursion to calculate the Fibonacci numbers, starting with the first two numbers (1 and 1) and then using the previous two numbers to calculate the next two number.

The code then uses a 'for' loop to print the first 10 Fibonacci numbers.</>

```

Response

Figure 6: Case Study on Code Generation. Tokens in **red** means those generated by Amphista and tokens in **black** means those generated by target model itself.

Input Prompt: What good habits can make the body healthier?

```

There are many good habits that can help keep the body healthy. Some of these include:
1. Eating a balanced diet that includes a variety of fruits, vegetables, whole grains, and lean proteins.
2. Getting regular exercise, such as walking, jogging, or lifting weights.
3. Getting enough sleep, ideally 7-9 hours per night.
4. Practicing good hygiene, such as washing your hands regularly and brushing your teeth twice a day.
5. Avoiding smoking and excessive alcohol consumption.
6. Managing stress through activities such as meditation, yoga, or deep breathing exercises.
7. Getting regular check-ups and screenings, such as blood pressure checks and cancer screenings.
8. Practicing safe sex and avoiding risky behaviors that can spread sexually transmitted infections.
9. Staying hydrated by drinking plenty of water throughout the day.
10. Practicing mindfulness and being aware of your body's needs and signals.</>

```

Response

Figure 7: Case Study on Text Generation. Tokens in **red** means those generated by Amphista and tokens in **black** means those generated by target model itself.

heads of Amphista, which significantly reduces the parameter amount and complexity. Specifically, we use two low-rank matrices to replace the original lm head matrix. The experimental results are shown in Table 11, we choose Vicuna 7B as target model, so the parameter count of lm head is 4096 \* 32000. With the increase of rank, we can reduce the number of learnable parameters by up to 45% while maintaining almost the same performance, which greatly reduces the complexity of model parameters and reflects the advantages and potential of Amphista in practical deployment.

### A.3 Case Study

Here we show some real case studies (see Figure 6, 7) on Amphista inference, the target model is Vicuna 7B. Note that we do not apply any special processing to the tokenizer's output, preserving the original results. Tokens highlighted in **red** represent those generated by Amphista during each step of decoding. Tokens in **black** indicate those generated by target model. From these practical examples, we can observe that in the vast majority of cases, Amphista generates at least two tokens per decoding step. This generally results in a stable at least 2x speed-up, demonstrating the efficiency of our algorithm. Additionally, Amphista's output is consistent with the target model's auto-regressive decoding output, ensuring the generation quality of Amphista.