

# TABED: Test-Time Adaptive Ensemble Drafting for Robust Speculative Decoding in LVLMs

Minjae Lee<sup>1\*</sup> Wonjun Kang<sup>1\*</sup> Byeongkeun Ahn<sup>1</sup> Christian Classen<sup>2</sup> Kevin Galim<sup>1</sup>  
Seunghyuk Oh<sup>1</sup> Minghao Yan<sup>2</sup> Hyung Il Koo<sup>1</sup> Kangwook Lee<sup>2,3</sup>

<sup>1</sup> FuriosaAI <sup>2</sup> UW-Madison <sup>3</sup> KRAFTON

{minjae.lee, kangwj1995}@furiosa.ai, kangwook.lee@wisc.edu

## Abstract

Speculative decoding (SD) has proven effective for accelerating LLM inference by quickly generating draft tokens and verifying them in parallel. However, SD remains largely unexplored for Large Vision-Language Models (LVLMs), which extend LLMs to process both image and text prompts. To address this gap, we benchmark existing inference methods with small draft models on 11 datasets across diverse input scenarios and observe scenario-specific performance fluctuations. Motivated by these findings, we propose *Test-time Adaptive Batched Ensemble Drafting (TABED)*, which dynamically ensembles multiple drafts obtained via batch inference by leveraging deviations from past ground truths available in the SD setting. The dynamic ensemble method achieves an average robust walltime speedup of 1.74× over autoregressive decoding and a 5% improvement over single drafting methods, while remaining training-free and keeping ensembling costs negligible through parameter sharing. With its plug-and-play compatibility, we further enhance TABED by integrating advanced verification and alternative drafting methods. Code and custom-trained models are available at <https://github.com/furiosa-ai/TABED>.

## 1 Introduction

Multimodal Large Language Models (MLLMs) (Wu et al., 2023; Zhang et al., 2024a; Jin et al., 2024; Song et al., 2023) are an advanced class of LLMs (Brown et al., 2020; Ouyang et al., 2022; Touvron et al., 2023; Guo et al., 2025) designed to process multiple modalities, such as images, audio, and video, alongside text. In particular, Large Vision Language Models (LVLMs) (Chen et al., 2024c), also known as Large Multimodal Models (Li et al., 2024b; Song et al., 2023), specialize in handling prompts comprised of text and images.

These models have attracted significant attention due to their unique applications, including multi-modal chatbots, visual question answering (VQA), and augmented reality (AR) (OpenAI, 2023; Anthropic, 2024; Gemini Team Google: Anil et al., 2023).

As LVLMs are increasingly deployed, reducing their inference time has become a critical challenge. In addition to the standard LLM autoregressive decoding process, LVLMs must (1) pre-process each image in the input into several hundred image tokens (Radford et al., 2021; Liu et al., 2023, 2024a) and (2) process both text and image tokens, resulting in considerably higher inference time.

Recently, methods like token pruning, layer skipping, and key-value cache compression have been proposed to accelerate LVLM inference (Shang et al., 2024; Chen et al., 2024b; Lin et al., 2024; Liu et al., 2024e; Wan et al., 2024). While effective, these approximation techniques cannot preserve the original LVLM’s output distribution. Moreover, they primarily reduce prompt processing time (i.e., the prefilling stage) and have less impact on response generation time (i.e., the decoding stage).

In contrast, Speculative Decoding (SD) is an acceleration technique for LLM inference that fully preserves the output distribution (Leviathan et al., 2023; Chen et al., 2023). SD speculates a specified number of draft tokens and then uses the original target model to verify these tokens in parallel. For LLM inference, various methods have been proposed to enhance different components of SD (Xia et al., 2024b; Hu et al., 2025b), such as aligning a lightweight draft model more closely with the target model (Zhou et al., 2024; Kang et al., 2025; Hu et al., 2025a), or constructing token trees to verify multiple drafts simultaneously using a target model augmented with additional components (Cai et al., 2024; Li et al., 2024d, 2025). These approaches yield substantial performance gains but often require additional training and reduce compatibility

\*Equal contribution.

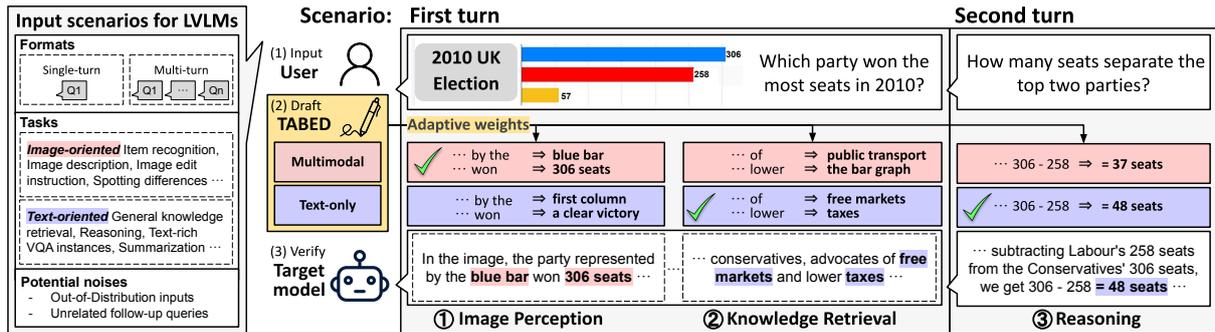


Figure 1: **Overview of TABED.** LVLMs must handle diverse input scenarios involving various combinations of turn-taking, tasks, and potential noise. To effectively accelerate target LVLMs with SD, different drafting methods are required for intra- and inter-response cases. For example, in the first turn, the drafting method needs to identify the visual context referenced by the user based on the input image and then perform text-based knowledge retrieval (① → ②). In the second turn, it performs mathematical reasoning using the accumulated text context (③). Across such varying scenarios, existing single drafting methods with small models—whether multimodal or text-only—exhibit fluctuating performance (Section 4). **TABED** addresses this by dynamically ensembling multiple drafts using past ground truths in SD, achieving robust speedups across diverse scenarios (Section 5). It can be further enhanced by integrating advanced verification and drafting methods (Section 6). See Fig. 2 for details.

with other acceleration methods, thereby limiting deployment accessibility.

Unlike in LLMs, SD for LVLMs has been far less explored. Gagrani et al. (2024) is the first to successfully accelerate LVLM inference via SD using a text-only drafting method that relies solely on the LLM component of an LVLM and text input tokens. While this approach is training-free and its findings are both intriguing and counterintuitive, it is evaluated only on a few single-turn VQA datasets, limiting insight into when the method performs well without image inputs. Another notable work, Kang et al. (2025), makes the draft model vision-aware by compressing image tokens and injecting a global visual context. However, it requires a dedicated data generation using the target model, followed by training the draft model.

This raises the question: “Can we develop a robust, vision-aware, and training-free drafting method that delivers consistent acceleration across diverse real-world input scenarios?” To answer this, we first benchmark existing inference methods for drafting at scale on seven benchmark datasets and two out-of-distribution (OOD) datasets using multi-turn instruction interactions. We find that none of the existing inference methods with a small draft model—whether multimodal or text-only (Gagrani et al., 2024)—effectively handle diverse input scenarios, making the selection of a single drafting method in advance a nontrivial challenge. To address this, we present Test-time Adaptive Batched Ensemble Drafting (TABED),

which (1) obtains multiple drafts simultaneously via batch inference and (2) dynamically ensembles them based on deviations from the ground truth provided by the target model. The method is entirely training-free and fully compatible with enhanced SD approaches in a plug-and-play manner.

As a result, TABED consistently achieves the best or second-best performance across all scenario combinations, with an average robust expected walltime speedup of 1.74× over standard decoding, without any additional tuning. It delivers over 5% improvement compared to existing single drafting methods, which exhibit scenario-specific performance fluctuations. We further enhance TABED by integrating advanced verification and alternative drafting methods, leveraging its plug-and-play compatibility. All components operate in a fully training-free manner with negligible ensembling costs, as demonstrated empirically. To further support LVLM SD research and ensure reproducibility, we open-source our code and draft LVLMs.

In summary, main contributions of our work are:

- We benchmark existing inference methods for drafting and find that their performance fluctuates across diverse LVLM input scenarios (Sec. 4).
- We propose a dynamic ensemble drafting method, TABED, which achieves superior and robust performance across diverse input scenarios (Sec. 5).
- We further enhance TABED by integrating advanced verification and alternative drafting methods via its plug-and-play compatibility (Sec. 6).

## 2 Related Work

### 2.1 Large Vision Language Models

**LVLMS** Frontier proprietary LVLMS (OpenAI, 2023; Anthropic, 2024; Gemini Team Google: Anil et al., 2023) exhibit state-of-the-art performance across multiple modalities beyond just text. Meanwhile, open-source models such as the LLaVA series (Liu et al., 2023, 2024a; Li et al., 2024b,a), Qwen-VL series (Wang et al., 2024b; Bai et al., 2025), and LLaMA 3.2 (Dubey et al., 2024) are also rapidly advancing.

While various methods exist for embedding image inputs (Yin et al., 2024; Jin et al., 2024), one of the most prominent approaches, LLaVA, employs an off-the-shelf vision encoder (Radford et al., 2021; Zhai et al., 2023) and a trainable projector to convert each image into several hundred visual context tokens for an LLM.

**Approximate Inference** To address the inefficiency of handling visual tokens from images, several approaches have been proposed based on a common finding: only a sparse subset of the hundreds of visual tokens is important, enabling reduced computational cost with minimal information loss. Shang et al. (2024); Chen et al. (2024b); Lin et al. (2024) dynamically prune significant visual tokens based on attention sparsity. Further focusing on reducing redundant key-value caches, Liu et al. (2024e); Wan et al. (2024) retain key-value vectors by merging or discarding less critical caches during inference. However, from a latency perspective, these approaches primarily benefit the prefilling stage while providing limited advantages for the decoding stage.

### 2.2 Speculative Decoding

**SD for LLMs** SD accelerates LLM inference using a small draft model while preserving the target model’s output distribution (Leviathan et al., 2023; Chen et al., 2023). To improve the drafting phase, various efforts have been made, including generating multiple draft candidates (Miao et al., 2023; Sun et al., 2024b; Yang et al., 2024; Li et al., 2025), and fine-tuning the draft model with knowledge distillation (Zhou et al., 2024; Hu et al., 2025a). Some studies employ speculative approaches to address exceptionally long sequence lengths that significantly degrade decoding efficiency (Sun et al., 2024a; Chen et al., 2024a; Galim et al., 2025). Separately, other work investigates newer archi-

tectures for SD, including State Space Models (SSMs) (Wang et al., 2024a; Choi et al., 2025).

**SD for LVLMS** Gagrani et al. (2024) introduces text-only drafting, which uses only the language model component of an LVLMS with text input. Without any additional training of the draft model (i.e., training-free), it achieves performance comparable to multimodal drafting. However, their benchmark results and analyses are limited, and they do not explore how to effectively leverage image information for improved drafting. Moreover, whether multiple drafting methods can be effectively combined remains unclear. Kang et al. (2025) present a notable work to make the draft model more vision-aware. Their method involves compressing image tokens and injecting a global visual context into the draft model. However, this approach requires a dedicated data generation phase using the target model, followed by a training process for the draft model. Jang et al. (2024) and Teng et al. (2024) propose SD methods for text-to-image generation, which is distinct from LVLMS approaches.

### 2.3 Test-time Adaptation

To enhance model robustness against distribution shifts during test time, various generalization and adaptation techniques have been proposed (Liang et al., 2025). In situations where data arriving during test time lacks ground truth labels, these techniques update models through explicit gradient steps by formulating loss using pseudo labeling (Hardt and Sun, 2024; Wang et al., 2022), self-supervision (Sun et al., 2020; Krause et al., 2018), or normalization statistics (Schneider et al., 2020; Wang et al., 2021), thereby causing inevitable delays during test time. In contrast, SD allows the use of ground-truth information obtained through verification for test-time adaptation. In this regard, Xia et al. (2024a) enables test-time adaptation for SD but incurs high adaptation cost, limiting its applicability to infrequent domain shifts.

## 3 Preliminaries: Speculative Decoding

**Algorithm** Following (Leviathan et al., 2023; Zhou et al., 2024), let  $M_p$  be the target model whose inference we aim to accelerate, and let  $M_q$  be the draft model for the same task. For a given prefix  $x$ , decoded sequence  $y_{<t}$ , block length  $\gamma$ , and in-block index  $n = 0, \dots, \gamma - 1$ , the following steps are repeated until either an  $\langle \text{EOS} \rangle$  token

is accepted or the maximum sequence length is reached:

1. The *Drafting Phase*, where  $M_q$  sequentially generates  $\gamma$  draft tokens  $y_{t+n} \sim q(\cdot|x, y_{<t+n})$ .
2. The *Verification Phase*, where  $M_p$  reviews these draft tokens in parallel, comparing them to  $p(y_{t+n}|x, y_{<t+n})$ .
3. For sampling, each token  $y_{t+n}$  is sequentially accepted with probability  $\min\left(1, \frac{p(y_{t+n}|x, y_{<t+n})}{q(y_{t+n}|x, y_{<t+n})}\right)$ . If any token is rejected before the end of the block, subsequent tokens are discarded, and the first rejected token is resampled from the adjusted distribution  $\text{norm}(\max(0, p(y) - q(y)))$ .<sup>1</sup>

**Block Efficiency and Walltime Speedup** The *block efficiency*, denoted as  $\tau_{p,q}(\gamma)$ , is defined as the expected number of accepted tokens per block. The expected *walltime speedup* quantifies the ratio between the per-token inference time of autoregressive decoding with  $M_p$  and that of the speculative decoding method. For a given batch size  $B$  and sequence length  $S$ , let  $T_p(B, S)$  and  $T_q(B, S)$  denote the time required for  $M_p$  and  $M_q$  to decode a single token, respectively, and let  $T_p(B, S, \gamma)$  denote the time required for  $M_p$  to verify  $\gamma$  tokens in parallel. For brevity, we omit  $(B, S)$  and use the simplified notations  $T_p$ ,  $T_q$ , and  $T_p(\gamma)$ . When  $B$  and  $S$  are small, the speculative decoding time per block for producing  $\tau_{p,q}(\gamma)$  valid tokens can be approximated as  $T_{\text{SD}} = \gamma \cdot T_q + T_p(\gamma) \approx \gamma \cdot T_q + T_p$  (Chen et al., 2024a). Under this approximation, the expected walltime speedup is given by

$$\text{Speedup} = \frac{T_p}{T_{\text{SD}}/\tau_{p,q}(\gamma)} \approx \frac{\tau_{p,q}(\gamma)}{\gamma \cdot \frac{T_q}{T_p} + 1}. \quad (1)$$

We first empirically measure the draft-to-target latency ratio  $T_q/T_p$ , and then use the computed speedup in Eq. (1) to evaluate the runtime efficiency of speculative decoding. See Section G for empirical measurement of the latency ratio  $T_q/T_p$  and analysis of its effects on latency and throughput.

## 4 Benchmarking SD for LVLMS

In this section, we systematically benchmark speculative decoding for LVLMS by evaluating existing drafting methods, both multimodal and text-only, across a variety of input scenarios.

<sup>1</sup>When the prefix  $(x, y_{<t})$  is clear from the context, we use  $p(y)$  and  $q(y)$  to denote  $p(y_t|x, y_{<t})$  and  $q(y_t|x, y_{<t})$ , respectively.

## 4.1 Experiment Settings

**Target and Draft Models** We employ 7B and 13B models from the LLaVA-1.5 (Liu et al., 2024a) and LLaVA-NeXT (Liu et al., 2024b) series, both of which are widely adopted public LVLMS, as our target models. We then use three distinct variants of the draft model, considering different model sizes and training strategies to ensure comprehensive benchmarking. Due to the lack of sufficiently small LVLMS that ensure a low *draft-to-target latency ratio* (represented by  $\frac{T_q}{T_p}$  in Eq. (1)), we select the small public LLaMA models of 68M and 160M (Miao et al., 2023) as base models for training. We then train the draft model using different training methodologies, specifically utilizing the procedures for LLaVA-1.5, as well as LLaVA-OV (Li et al., 2024a), which are specialized in multi-image processing. We further verify that the draft model captures multimodality through representative multimodal tasks (see Section H.1 for training details and Section H.2 for evaluation results).

To summarize, we benchmarked draft models by size (68M, 160M) and type (LLaVA-1.5, LLaVA-OV), and target models by size (7B, 13B) and type (LLaVA-1.5, LLaVA-NeXT). For brevity, we primarily report results with the LLaVA-1.5 68M–7B pair, unless otherwise noted.

**Input Scenarios** Systems using SD must deliver consistent performance gains across diverse real-world scenarios, especially due to their application with LVLMS in practical settings. To examine this requirement, we initially curated seven benchmark datasets due to the lack of pre-existing benchmarks specifically for LVLMS SD. These datasets include tasks involving both single-image (Liu et al., 2023; Mathew et al., 2021; Li et al., 2023b; Yu et al., 2023) and multi-image scenarios (Tan et al., 2019; Zhang et al., 2024b; Jhamtani and Berg-Kirkpatrick, 2018). To further challenge the system’s ability to handle unexpected and varied queries while maintaining consistent performance, we added two additional datasets featuring five images per query (Li et al., 2019; Huang et al., 2016), serving as notable Out-of-Distribution (OOD) cases. Moreover, we extend the evaluation from single-turn to multi-turn scenarios using benchmark datasets that include various types of follow-up queries, including those dependent on prior responses (e.g., follow-up requests with images from the same dataset or text-only tasks from multi-turn benchmarks for LVLMS (Liu

Table 1: Block efficiency results for drafting methods are shown, including existing single-drafting approaches—multimodal (M) and text-only (T)—and our proposed method, **TABED<sup>MT</sup>**, which dynamically ensembles M and T through test-time adaptation. Across diverse combinations of turn-taking settings, datasets (both benchmark and OOD), image inclusion, and contextual relatedness, **TABED<sup>MT</sup>** consistently achieves either the **best** or **second-best** performance, while M and T vary depending on the scenario. Even when ranked second, its performance remains close to the best, unlike the larger gaps observed in other methods.

Drafting		Benchmark Datasets (First Turn)								OOD Datasets	
Type	Method	LLaVA-W	DocVQA	POPE	MMVet	IEdit	MB	Spot	Avg.	PSV	VIST
Single	M	<b>2.28</b>	<u>2.15</u>	<b>2.56</b>	<b>2.21</b>	2.19	1.96	<u>2.34</u>	<u>2.24</u>	1.19	1.16
	T [16]	2.19	2.08	2.31	2.16	<u>2.23</u>	<u>2.34</u>	2.27	2.23	<b>2.05</b>	<b>2.05</b>
Ensemble	<b>TABED<sup>MT</sup></b>	<u>2.26</u>	<b>2.16</b>	<u>2.52</u>	<b>2.21</b>	<u>2.29</u>	<b>2.39</b>	<b>2.36</b>	<b>2.31</b>	<u>2.02</u>	<u>2.04</u>
Drafting		Benchmark Datasets (Second Turn)								NLP Datasets	
Type	Method	LLaVA-W	DocVQA	POPE	MMVet	IEdit	MB	Spot	Avg.	NQ	GSM8K
Single	M	2.10	1.96	2.78	2.18	1.61	1.53	1.83	2.00	1.98	2.25
	T [16]	<b>2.32</b>	<b>2.23</b>	<u>2.91</u>	<b>2.56</b>	<b>1.87</b>	<b>2.01</b>	<b>2.08</b>	<b>2.28</b>	<b>2.03</b>	<b>2.30</b>
Ensemble	<b>TABED<sup>MT</sup></b>	<u>2.29</u>	<b>2.23</b>	<b>2.93</b>	<b>2.56</b>	<u>1.85</u>	<u>1.99</u>	<u>2.05</u>	<u>2.27</u>	<b>2.03</b>	<u>2.29</u>

et al., 2024c)) and distinct text-only reasoning tasks (Kwiatkowski et al., 2019; Cobbe et al., 2021).

#### Drafting Methods: Multimodal and Text-only

The multimodal drafting method follows the standard LVLM process, receiving both images and text as input. In contrast, the text-only drafting method, first explored in Gagrani et al. (2024), uses only textual data, adhering to the standard LLM process. We primarily report results with  $\gamma = 5$  and greedy decoding with a maximum of 128 new tokens, unless otherwise noted.

## 4.2 Experimental Results

Table 1 presents the block efficiency results for multimodal (M) and text-only (T) drafting methods. Across all scenarios, neither drafting method consistently outperformed the other. Notably, the multimodal drafting method demonstrated higher block efficiency than the text-only approach in single-turn scenarios across most benchmark datasets. However, the text-only drafting performed comparably or slightly sub-optimally overall in single-turn situations, often surpassing the multimodal drafting in subsequent turns—particularly in tasks dependent on prior responses and independent reasoning tasks, and significantly in OOD cases. This trend continued regardless of the model size or multi-image awareness capabilities. Existing drafting methods, both multimodal and text-only, show scenario-specific performance fluctuations with draft models small enough for SD, with neither method consistently

outperforming the other across various input scenarios for LVLM inference. Before execution, it is challenging to determine in advance which method is superior, and even if known, resolving inconsistencies with a single drafting method is difficult. See Sections A and B for full results for various  $M_q$  and  $M_p$  sizes and types, and different  $\gamma$ , are provided.

*No one-size-fits-all among existing single drafting methods for LVLM input scenarios.*

## 5 Test-time Adaptive Batched Ensemble Drafting

For LVLM SD, selecting a single drafting method in advance for a small draft model is challenging. In this context, ensemble learning offers a promising solution, as it is often used to reduce both bias and variance in predictions (Dietterich, 2000; Ganaie et al., 2022), particularly for models with limited capacity (Zhou, 2012).

In this section, we introduce Test-time Adaptive Batched Ensemble Drafting (TABED), a method that applies ensemble learning with dynamic weight adaptation to each of the drafts obtained via batch inference. It is fully training-free, highly extensible, and delivers robust performance across diverse input scenarios. Fig. 2 provides an overview of TABED.

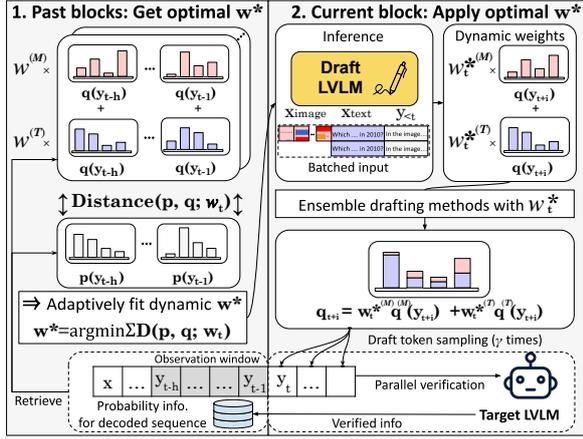


Figure 2: TABED predicts the optimal weight,  $w^*$ , based on the deviation (i.e., distance) of past drafting blocks from the ground truth obtained through the verification in SD. It then dynamically ensembles multiple drafts obtained via batch inference. See Algorithm 1 for further details.

## 5.1 Proposed Method

Unlike typical test-time scenarios where incoming data lacks ground truth labels (Wang et al., 2021, 2022), SD allows access to both hard labels (i.e.,  $y_{<t}$ ) and soft labels (i.e.,  $p(y_{t'})$ ) for  $t' < t$  after verifying all steps prior to  $t$  using the target model  $M_p$ . When drafting is performed in step  $t$ , this information can be leveraged to dynamically adjust ensemble weights  $w_t$ , controlling each drafting method’s influence based on its past performance. This approach can enhance ensemble learning by effectively combining probabilities  $[q_t^{(1)}, \dots, q_t^{(m)}]$  from  $m$  drafting methods, ensuring the resulting ensemble distribution  $q(\cdot|x, y_{<t}; w_t)$  closely aligns with the target distribution  $p(\cdot|x, y_{<t})$ . For instance, the system can adjust the weight assigned to the multimodal drafting method by recognizing varying needs for visual context for a specific sample  $x$  and timestep  $t$ .

**Ensemble learning via Batched Inference** To obtain multiple predictions from draft model  $M_q$ , we utilize batched inference tailored for LVLMLs. At each decoding timestep  $t$  to  $t + \gamma - 1$ , all  $m$  drafting methods share the parameters of  $M_q$ , and their distributions are ensembled to sample the next token (Algorithm 1). We employ a weighted averaging ensemble method, sampling a token from the ensemble distribution to continue drafting. This process incurs no additional costs during training (i.e., training-free) and keeps ensembling costs negligible, unlike typical ensemble learning that re-

quires multiple extra models and computational overhead for deployment.

**Test-time Adaptive Ensemble Weights** To effectively combine drafting methods based on varying needs, we explore test-time adaptation to dynamically weight each method. Specifically, at the beginning of each new drafting block, we sample a list of ensemble weights  $W_t = [w_t^1, \dots, w_t^n]$  using the weight sampling policy  $S_W(t)$  (Algorithm 1). For each  $w_t^j \in \mathbb{R}^m$ , we compute ensemble draft probabilities  $q_{t'}^j = q(\cdot|x, y_{<t'}; w_t^j)$  over the past time step window  $t' \in [t - h, t - 1]$ . We then select the optimal  $w_t^*$  from  $W_t$  to be used at timestep  $t$  by utilizing hard labels  $y_{<t}$  (by maximizing the number of token matches sampled from each  $q_{t'}^j$ ), or soft labels  $p(y_{t'})$  (by minimizing the accumulated error  $e_t^j$  over previous steps  $t'$ ).

$$e_t^j = \sum_{t'} D_{\text{KL}}(p(\cdot | x, y_{<t'}) \parallel q(\cdot | x, y_{<t'}; w_t^j))$$

where  $D_{\text{KL}}$  is the KL divergence between  $p$  and  $q_{t'}^j$  at each of the previous steps  $t'$ . This weight  $w_t^*$  is used throughout the current drafting block of  $\gamma$  tokens (i.e., from timestep  $t$  to  $t + \gamma - 1$ ). To avoid redundant computation during drafting, we cache and reuse the distributions  $q_{t'}^j$  from the  $M_q$  and the labels  $p(y_{t'})$  from the  $M_p$ , computed over the past time-step window  $t' \in [t - h, t - 1]$ .

Drafting methods with higher weights indicate closer alignment to the target model. In our experiments, we also explored the effects of employing various weight sampling policies  $S_W$ , varying the window size  $h$  for past time steps, and using Total Variation Distance (TVD) instead of KL divergence.

## 5.2 Experimental Results

Table 1 presents the block efficiency results for TABED, which is denoted by TABED<sup>MT</sup>. TABED<sup>MT</sup> consistently achieved either the best or second-best performance across all input scenarios. These scenarios include turn-taking, image inclusion, prior context relevance, and both benchmark and OOD datasets. Even when TABED<sup>MT</sup> ranked second, its performance remained close to the best, unlike the larger performance gaps observed for other methods. Without incurring additional costs during training or inference, this resulted in and a 5% improvement compared to existing single drafting methods M and T, when averaged across all

---

**Algorithm 1** TABED

---

**Parameter:**  $M_q$ , Prefix  $X = [x^{(1)}, \dots, x^{(m)}]$  for  $m$  drafting methods, Weight sampling policy  $S_W$ , local window length  $h$

▷ W-AVG stands for Weighted Average  
▷  $D$  stands for Distance

**Input:** Batch  $b_{<t} := [(x^{(1)}, y_{<t}), \dots, (x^{(m)}, y_{<t})]$

**Output:**  $\gamma$  draft tokens  $y_t, \dots, y_{t+\gamma-1}$  and ensemble probabilities  $q_t, \dots, q_{t+\gamma-1}$

```
1: procedure TABED( $b_{<t}; \gamma, S_W$ )
2:    $W_t = [w_t^1, \dots, w_t^n] \sim S_W(t)$ 
3:    $w_t^* \leftarrow \arg \min_j \sum_{t' \in [t-h, t-1]} D(p_{t'}, q_{t'}^j)$ 
4:   for  $i \leftarrow 0$  to  $\gamma - 1$  do
5:      $[q_{t+i}^{(1)}, \dots, q_{t+i}^{(m)}] \leftarrow M_q(b_{<t+i})$ 
6:      $q_{t+i} \leftarrow \text{W-AVG}([q_{t+i}^{(1)}, \dots, q_{t+i}^{(m)}]; w_t^*)$ 
7:      $y_{t+i} \leftarrow \text{SAMPLE}(q_{t+i})$ 
8:   end for
9:   return  $[y_t, \dots, y_{t+\gamma-1}], [q_t, \dots, q_{t+\gamma-1}]$ 
10: end procedure
```

---

benchmark datasets and turns. This leads to an average robust walltime speedup of 1.74 $\times$ , computed using Eq. (1), compared to standard decoding. These consistent gains highlight TABED’s ability to effectively prioritize stronger drafting methods by dynamically assigning optimal weights, even when individual methods exhibit scenario-specific performance fluctuations. For full results on varying the parameters of TABED (sampling policy  $S_W$ , window size  $h$ , and distance  $D$ ), as well as a detailed derivation of the walltime speedup with empirical analysis on latency and throughput, see Sections B and G, respectively.

## 6 Extensions and Analyses of TABED

Since TABED is plug-and-play and training-free, it remains fully compatible with and can benefit from enhanced SD (Xia et al., 2023; Hu et al., 2025b). In this section, we extend TABED by incorporating an enhanced verification method and integrating newly explored drafting candidates. We then provide in-depth analyses and ablations to qualitatively validate the contribution of the proposed components.

### 6.1 Further Extensions of TABED

**Verification: Integrating Token Tree** We integrate the well-known token-tree verification method (Miao et al., 2023; Yang et al., 2024; Cai

et al., 2024; Li et al., 2025) into TABED, where all draft sequences from the draft model  $M_q$  are represented as a token tree and verified in parallel by the target model  $M_p$ . Specifically, we employ tree widths  $d = 2$  and 3 and report the results. Table 2 presents block efficiency results averaged across datasets in each category for single drafting methods (M and T) and the dynamic ensemble method (TABED) combined with token-tree verification. TABED<sup>MT</sup> consistently achieves the best or second-best performance, while M and T vary across scenarios. See Section D for full results.

Verification	Drafting		Block Efficiency	
	Tree width	Type Method	Benchmark	OOD
$d = 2$	Sgl.	M T [16]	2.89 2.85	1.30 <b>2.72</b>
		Ens. <b>TABED<sup>MT</sup></b>	<b>2.99</b>	<u>2.64</u>
$d = 3$	Sgl.	M T [16]	3.28 3.24	1.38 <b>3.19</b>
		Ens. <b>TABED<sup>MT</sup></b>	<b>3.39</b>	<u>3.09</u>

Table 3: Averaged block efficiency results for single drafting methods (M, T, caption-based (C), and pooled-multimodal (P)) and the dynamic ensemble method TABED<sup>MTCP</sup>. TABED<sup>MTCP</sup> consistently achieves the best performance in both the first and second turns across all benchmark and OOD datasets, followed by C, which consistently achieves the second-best performance. See Section A for full results.

	Drafting		Benchmark		OOD
	Type	Method	First	Second	
Single		M	2.24	2.00	1.18
		T [16]	2.23	2.28	2.05
		C	<u>2.29</u>	<u>2.30</u>	<u>2.09</u>
		P	2.23	2.25	2.08
Ensemble	<b>TABED<sup>MTCP</sup></b>	<b>2.32</b>	<b>2.32</b>	<b>2.13</b>	

et al., 2024; Li et al., 2025) into TABED, where all draft sequences from the draft model  $M_q$  are represented as a token tree and verified in parallel by the target model  $M_p$ . Specifically, we employ tree widths  $d = 2$  and 3 and report the results. Table 2 presents block efficiency results averaged across datasets in each category for single drafting methods (M and T) and the dynamic ensemble method (TABED) combined with token-tree verification. TABED<sup>MT</sup> consistently achieves the best or second-best performance, while M and T vary across scenarios—reproducing the trend observed in Table 1 with an even larger block-efficiency gap. Full results are provided in Section D.

### Drafting: Exploring Additional Candidates

Since handling sparsely important image tokens is a unique challenge for LVLMs—absent in typical LLMs (Shang et al., 2024; Chen et al., 2024b)—we explore effective strategies to address this issue and extend TABED accordingly.

*Pooled Multimodal Drafting (P)* To condense

Table 4: Averaged block efficiency of ensemble drafting methods with different weighting strategies (Random<sup>MT</sup>, Static<sup>MT</sup>, TABED<sup>MT</sup>) applied to M and T. **TABED<sup>MT</sup>** consistently achieves performance that is equal to or the **best**. The second-best result is underlined.

Drafting		Benchmark		OOD
Type	Weights	First	Second	
Ensemble	Random <sup>MT</sup>	2.29	2.17	1.83
	Static <sup>MT</sup>	<u>2.31</u>	<u>2.25</u>	<u>1.93</u>
	<b>TABED<sup>MT</sup></b>	<b>2.31</b>	<b>2.27</b>	<b>2.03</b>

sparingly important image tokens while preserving their 2D spatial structure, we apply average pooling during inference immediately before the projector maps them into the text embedding space. Specifically, we use a  $2 \times 2$  pooling kernel, reducing the number of visual tokens from 576 to 144 in our default configuration.

*Caption Drafting (C)* To convert sparsely important image tokens into textual descriptions, we use a lightweight captioning model, Florence-2 (Xiao et al., 2024), executed once in parallel with the pre-filling stage of the target model  $M_p$ . This makes its overhead negligible and amortized over the entire decoding process (Fig. 9). For details on replacing image tokens in caption drafting and an analysis of captioning-model latency, see Sections E and I.

Table 3 presents block efficiency results averaged across datasets in each category for the newly explored single drafting methods (P and C) and the dynamic ensemble method TABED<sup>MTCP</sup>. TABED<sup>MTCP</sup> consistently achieves the highest speedup across all benchmark and OOD datasets, outperforming all single drafting methods (M, T, C, and P) in both the first and second turns. This seamless integration improves generalization (Zhou, 2012) and better aligns the ensemble distribution with the target distribution, without additional training cost, while keeping ensembling costs negligible through parameter sharing. While P and C introduce image awareness into the drafting process, C in particular retains T’s robustness and performs second-best on both standard and OOD datasets. See Sections A and G for full results and analysis on latency and throughput of TABED<sup>MTCP</sup>.

## 6.2 Analysis and Ablations of TABED

**Effectiveness of Dynamic Weights** To assess the role of dynamic ensemble weights, we conduct an ablation study in which they are replaced with non-negative, sum-to-one random weights (Random<sup>MT</sup>) or static equal weights (Static<sup>MT</sup>), used to combine

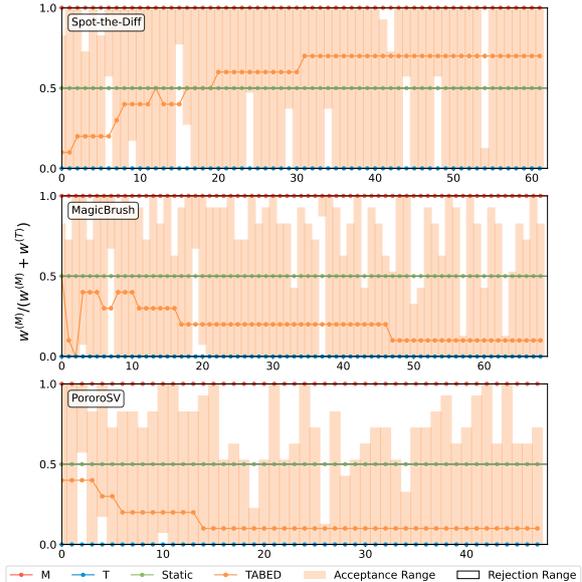


Figure 3: Qualitative samples visualizing dynamic ensemble weights across datasets. The x-axis and y-axis represent decoding steps and the proportion of  $w^{(M)}$  relative to  $w^{(M)} + w^{(T)}$ , respectively. The acceptance range (shaded) comprises weights under which  $M_p$  accepts the drafted tokens, while the rejection range (unshaded) comprises weights under which it rejects them. See Section C for more examples.

multiple drafts generated via batched inference. Table 4 presents block efficiency results averaged across datasets in each category for these ablations, and TABED<sup>MT</sup> achieves performance equal to or better than both Random<sup>MT</sup> and Static<sup>MT</sup>. In particular, TABED<sup>MT</sup> clearly outperforms Static<sup>MT</sup> in OOD scenarios by effectively down-weighting the influence of underperforming drafting methods. See Section A for full results across different  $M_q$  and  $M_p$  sizes.

**Adaptation Behavior of Dynamic Weights** We provide an in-depth view of how the adaptive weights  $w^*$  predicted by TABED evolve during decoding. Fig. 3 visualizes the weights used by single drafting methods (M and T at  $y = 1$  and 0) and ensemble methods (Static<sup>MT</sup> and TABED<sup>MT</sup> at  $y = 0.5$  and the orange line) across decoding steps. The weights  $w^*$  predicted by TABED<sup>MT</sup> effectively stay within the shaded acceptance range while avoiding the unshaded rejection range, unlike single drafting methods (M and T) or the static ensemble Static<sup>MT</sup>. This adaptive behavior enables TABED<sup>MT</sup> to discriminate effectively among multiple drafting methods, explaining the large performance gap observed under OOD settings in Table 4. More examples are provided in Section C.

## 7 Conclusion

We benchmark existing inference methods for LVLM drafting and observe scenario-specific performance fluctuations. To address this, we propose TABED, that dynamically ensembles multiple drafts obtained via batch inference by leveraging deviations of the drafts from past ground truths. TABED consistently achieves the best or second-best performance across all scenarios, with an average robust expected walltime speedup of 1.74× and over 5% improvement over single drafting methods, all without additional tuning. TABED can be further enhanced through advanced verification or alternative drafting methods via its plug-and-play compatibility and remains fully training-free with negligible ensembling costs.

## 8 Limitations and Future Work

While we focus on a plug-and-play, training-free inference method for open-source LVLMs, TABED can also leverage trained draft models (Zhou et al., 2024; Yang et al., 2024; Hu et al., 2025a), as it is fully compatible with such approaches. Moreover, it would be valuable to evaluate a broader range of LVLM families (e.g., Qwen-VL and LLaMA 3.2), beyond LLaVA, as well as to explore scaled-up draft models  $M_q$  with billions of parameters rather than limiting experiments to small-capacity draft LVLMs. A more theoretically grounded approach to modeling the relationship between the past time-step window and the dynamic ensemble weights could further improve the method. We further expect TABED to generalize to other MLLMs (e.g., those handling audio (Fu et al., 2024)).

## Acknowledgement

Kangwook Lee is supported by NSF Award DMS-2023239, NSF CAREER Award CCF-2339978, Amazon Research Award, and a grant from FuriosaAI. We thank all members of the FuriosaAI team for their support, with special thanks to Hanjoon Kim and June Paik for their vision and commitment to research. We also thank the Lee Lab at UW–Madison, led by Kangwook Lee, for their collaboration. Finally, we thank Seongsu Bae, Eunbyeol Cho, Gyubok Lee, and Sungjin Park at KAIST EDLAB, led by Edward Choi, for their valuable comments and inspiration.

## Reproducibility Statement

The code used for all experiments in this paper is publicly available. To support reproducibility, both the code and checkpoints can be accessed at the following GitHub repository: <https://github.com/furiosa-ai/TABED>.

## References

- Anthropic. 2024. *The Claude 3 Model Family: Opus, Sonnet, Haiku*.
- Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Siboz Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, et al. 2025. Qwen2. 5-vl technical report. *arXiv preprint arXiv:2502.13923*.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. In *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901.
- 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*.
- Jian Chen, Vashisth Tiwari, Ranajoy Sadhukhan, Zhuoming Chen, Jinyuan Shi, Ian En-Hsu Yen, and Beidi Chen. 2024a. MagicDec: Breaking the Latency-Throughput Tradeoff for Long Context Generation with Speculative Decoding. *arXiv preprint arXiv:2408.11049*.
- Liang Chen, Haozhe Zhao, Tianyu Liu, Shuai Bai, Junyang Lin, Chang Zhou, and Baobao Chang. 2024b. An image is worth 1/2 tokens after layer 2: Plug-and-play inference acceleration for large vision-language models. *arXiv preprint arXiv:2403.06764*.
- Yangyi Chen, Karan Sikka, Michael Cogswell, Heng Ji, and Ajay Divakaran. 2024c. Dress: Instructing large vision-language models to align and interact with humans via natural language feedback. *Preprint, arXiv:2311.10081*.
- Daewon Choi, Seunghyuk Oh, Saket Dingliwal, Jihoon Tack, Kyuyoung Kim, Woomin Song, Seojin Kim, Insu Han, Jinwoo Shin, Aram Galstyan, et al. 2025. Mamba drafters for speculative decoding. *arXiv preprint arXiv:2506.01206*.

- 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*.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186.
- Thomas G Dietterich. 2000. Ensemble methods in machine learning. In *International workshop on multiple classifier systems*, pages 1–15. Springer.
- 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*.
- Chaoyou Fu, Haojia Lin, Zuwei Long, Yunhang Shen, Meng Zhao, Yifan Zhang, Xiong Wang, Di Yin, Long Ma, Xiawu Zheng, et al. 2024. Vita: Towards open-source interactive omni multimodal llm. *arXiv preprint arXiv:2408.05211*.
- Yao Fu. 2024. Challenges in Deploying Long-Context Transformers: A Theoretical Peak Performance Analysis. *arXiv preprint arXiv:2405.08944*.
- Mukul Gagrani, Raghav Goel, Wonseok Jeon, Junyoung Park, Mingyu Lee, and Christopher Lott. 2024. On speculative decoding for multimodal large language models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8285–8289.
- Kevin Galim, Ethan Ewer, Wonjun Kang, Minjae Lee, Hyung Il Koo, and Kangwook Lee. 2025. Draft-based approximate inference for llms. *arXiv preprint arXiv:2506.08373*.
- Mudasir A Ganaie, Minghui Hu, Ashwani Kumar Malik, Muhammad Tanveer, and Ponnuthurai N Suganthan. 2022. Ensemble deep learning: A review. *Engineering Applications of Artificial Intelligence*, 115:105151.
- Rohan Gemini Team Google: Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. 2023. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*.
- 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*.
- Moritz Hardt and Yu Sun. 2024. [Test-time training on nearest neighbors for large language models](#). *Preprint*, arXiv:2305.18466.
- Yuezhou Hu, Jiaxin Guo, Xinyu Feng, and Tuo Zhao. 2025a. Adaspec: Selective knowledge distillation for efficient speculative decoders. *arXiv preprint arXiv:2510.19779*.
- Yunhai Hu, Zining Liu, Zhenyuan Dong, Tianfan Peng, Bradley McDanel, and Sai Qian Zhang. 2025b. Speculative decoding and beyond: An in-depth survey of techniques. *arXiv preprint arXiv:2502.19732*.
- Ting-Hao Huang, Francis Ferraro, Nasrin Mostafazadeh, Ishan Misra, Aishwarya Agrawal, Jacob Devlin, Ross Girshick, Xiaodong He, Pushmeet Kohli, Dhruv Batra, et al. 2016. Visual storytelling. In *Proceedings of the 2016 conference of the North American chapter of the association for computational linguistics: Human language technologies*, pages 1233–1239.
- Doohyuk Jang, Sihwan Park, June Yong Yang, Yeonsung Jung, Jihun Yun, Souvik Kundu, Sung-Yub Kim, and Eunho Yang. 2024. Lantern: Accelerating visual autoregressive models with relaxed speculative decoding. *arXiv preprint arXiv:2410.03355*.
- Harsh Jhamtani and Taylor Berg-Kirkpatrick. 2018. Learning to Describe Differences Between Pairs of Similar Images. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing (EMNLP)*.
- Yizhang Jin, Jian Li, Yexin Liu, Tianjun Gu, Kai Wu, Zhengkai Jiang, Muyang He, Bo Zhao, Xin Tan, Zhenye Gan, Yabiao Wang, Chengjie Wang, and Lizhuang Ma. 2024. [Efficient Multimodal Large Language Models: A Survey](#). *Preprint*, arXiv:2405.10739.
- Jialiang Kang, Han Shu, Wenshuo Li, Yingjie Zhai, and Xinghao Chen. 2025. Vispec: Accelerating vision-language models with vision-aware speculative decoding. *arXiv preprint arXiv:2509.15235*.
- Ben Krause, Emmanuel Kahembwe, Iain Murray, and Steve Renals. 2018. [Dynamic evaluation of neural sequence models](#). In *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 2766–2775. PMLR.
- Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, et al. 2019. Natural questions: a benchmark for question answering research. *Transactions of the Association for Computational Linguistics*, 7:453–466.
- Yaniv Leviathan, Matan Kalman, and Yossi Matias. 2023. Fast inference from transformers via speculative decoding. In *International Conference on Machine Learning*, pages 19274–19286. PMLR.

- Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Yanwei Li, Ziwei Liu, and Chunyuan Li. 2024a. LLaVA-OneVision: Easy Visual Task Transfer. *arXiv preprint arXiv:2408.03326*.
- Feng Li, Renrui Zhang, Hao Zhang, Yuanhan Zhang, Bo Li, Wei Li, Zejun Ma, and Chunyuan Li. 2024b. LLaVA-NeXT-Interleave: Tackling Multi-image, Video, and 3D in Large Multimodal Models. *arXiv preprint arXiv:2407.07895*.
- Feng Li, Renrui Zhang, Hao Zhang, Yuanhan Zhang, Bo Li, Wei Li, Zejun Ma, and Chunyuan Li. 2024c. Llava-next-interleave: Tackling multi-image, video, and 3d in large multimodal models. *Preprint*, arXiv:2407.07895.
- Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. 2023a. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. In *International conference on machine learning*, pages 19730–19742. PMLR.
- Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. 2022. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In *International conference on machine learning*, pages 12888–12900. PMLR.
- Yifan Li, Yifan Du, Kun Zhou, Jinpeng Wang, Wayne Xin Zhao, and Ji-Rong Wen. 2023b. Evaluating Object Hallucination in Large Vision-Language Models. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 292–305.
- Yitong Li, Zhe Gan, Yelong Shen, Jingjing Liu, Yu Cheng, Yuexin Wu, Lawrence Carin, David Carlson, and Jianfeng Gao. 2019. Storygan: A sequential conditional gan for story visualization. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 6329–6338.
- Yuhui Li, Fangyun Wei, Chao Zhang, and Hongyang Zhang. 2024d. Eagle: Speculative sampling requires rethinking feature uncertainty. *arXiv preprint arXiv:2401.15077*.
- Yuhui Li, Fangyun Wei, Chao Zhang, and Hongyang Zhang. 2025. Eagle-3: Scaling up inference acceleration of large language models via training-time test. *arXiv preprint arXiv:2503.01840*.
- Jian Liang, Ran He, and Tieniu Tan. 2025. A comprehensive survey on test-time adaptation under distribution shifts. *International Journal of Computer Vision*, 133(1):31–64.
- Zhihang Lin, Mingbao Lin, Luxi Lin, and Rongrong Ji. 2024. Boosting Multimodal Large Language Models with Visual Tokens Withdrawal for Rapid Inference. *arXiv preprint arXiv:2405.05803*.
- Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. 2024a. Improved baselines with visual instruction tuning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 26296–26306.
- Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee. 2024b. Llava-next: Improved reasoning, ocr, and world knowledge.
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023. Visual Instruction Tuning. In *Advances in Neural Information Processing Systems*.
- Shuo Liu, Kaining Ying, Hao Zhang, Yue Yang, Yuqi Lin, Tianle Zhang, Chuanhao Li, Yu Qiao, Ping Luo, Wenqi Shao, and Kaipeng Zhang. 2024c. Convbench: A multi-turn conversation evaluation benchmark with hierarchical capability for large vision-language models. *Preprint*, arXiv:2403.20194.
- Yuliang Liu, Zhang Li, Mingxin Huang, Biao Yang, Wenwen Yu, Chunyuan Li, Xucheng Yin, Cheng lin Liu, Lianwen Jin, and Xiang Bai. 2024d. Ocrbench: On the hidden mystery of ocr in large multimodal models. *Preprint*, arXiv:2305.07895.
- Zuyan Liu, Benlin Liu, Jiahui Wang, Yuhao Dong, Guangyi Chen, Yongming Rao, Ranjay Krishna, and Jiwen Lu. 2024e. Efficient Inference of Vision Instruction-Following Models with Elastic Cache. *arXiv preprint arXiv:2407.18121*.
- Ahmed Masry, Xuan Long Do, Jia Qing Tan, Shafiq Joty, and Enamul Hoque. 2022. ChartQA: A Benchmark for Question Answering about Charts with Visual and Logical Reasoning. In *Findings of the Association for Computational Linguistics: ACL 2022*, pages 2263–2279.
- Minesh Mathew, Dimosthenis Karatzas, and CV Jawahar. 2021. Docvqa: A dataset for vqa on document images. In *Proceedings of the IEEE/CVF winter conference on applications of computer vision*, pages 2200–2209.
- Xupeng Miao, Gabriele Oliaro, Zhihao Zhang, Xinhao Cheng, Zeyu Wang, Rae Ying Yee Wong, Zhuoming Chen, Daiyaan Arfeen, Reyna Abhyankar, and Zhihao Jia. 2023. Specinfer: Accelerating generative llm serving with speculative inference and token tree verification. *arXiv preprint arXiv:2305.09781*.
- OpenAI. 2023. GPT-4 technical report. *arXiv preprint arXiv:2303.08774*.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Gray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. 2022. Training language models to follow instructions with human feedback. In *Advances in Neural Information Processing Systems*.

- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. 2021. [Learning Transferable Visual Models From Natural Language Supervision](#). In *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages 8748–8763. PMLR.
- Steffen Schneider, Evgenia Rusak, Luisa Eck, Oliver Bringmann, Wieland Brendel, and Matthias Bethge. 2020. Improving robustness against common corruptions by covariate shift adaptation. *Advances in neural information processing systems*, 33:11539–11551.
- Yuzhang Shang, Mu Cai, Bingxin Xu, Yong Jae Lee, and Yan Yan. 2024. Llava-prumerge: Adaptive token reduction for efficient large multimodal models. *arXiv preprint arXiv:2403.15388*.
- Oleksii Sidorov, Ronghang Hu, Marcus Rohrbach, and Amanpreet Singh. 2020. [Textcaps: a dataset for image captioning with reading comprehension](#). *Preprint*, arXiv:2003.12462.
- Shezheng Song, Xiaopeng Li, Shasha Li, Shan Zhao, Jie Yu, Jun Ma, Xiaoguang Mao, and Weimin Zhang. 2023. [How to bridge the gap between modalities: A comprehensive survey on multimodal large language model](#). *Preprint*, arXiv:2311.07594.
- 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*.
- Yu Sun, Xiaolong Wang, Zhuang Liu, John Miller, Alexei A. Efros, and Moritz Hardt. 2020. [Test-time training with self-supervision for generalization under distribution shifts](#). *Preprint*, arXiv:1909.13231.
- Ziteng Sun, Ananda Theertha Suresh, Jae Hun Ro, Ahmad Beirami, Himanshu Jain, and Felix Yu. 2024b. [SpecTr: Fast Speculative Decoding via Optimal Transport](#). *Preprint*, arXiv:2310.15141.
- Hao Tan, Franck Dernoncourt, Zhe Lin, Trung Bui, and Mohit Bansal. 2019. Expressing Visual Relationships via Language. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 1873–1883.
- Yao Teng, Han Shi, Xian Liu, Xuefei Ning, Guohao Dai, Yu Wang, Zhenguo Li, and Xihui Liu. 2024. Accelerating auto-regressive text-to-image generation with training-free speculative jacobi decoding. *arXiv preprint arXiv:2410.01699*.
- 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*.
- Zhongwei Wan, Ziang Wu, Che Liu, Jinfa Huang, Zhihong Zhu, Peng Jin, Longyue Wang, and Li Yuan. 2024. LOOK-M: Look-Once Optimization in KV Cache for Efficient Multimodal Long-Context Inference. *arXiv preprint arXiv:2406.18139*.
- Dequan Wang, Evan Shelhamer, Shaoteng Liu, Bruno Olshausen, and Trevor Darrell. 2021. Tent: Fully test-time adaptation by entropy minimization. In *International Conference on Learning Representations*.
- Junxiong Wang, Daniele Paliotta, Avner May, Alexander Rush, and Tri Dao. 2024a. The mamba in the llama: Distilling and accelerating hybrid models. *Advances in Neural Information Processing Systems*, 37:62432–62457.
- Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, et al. 2024b. Qwen2-vl: Enhancing vision-language model’s perception of the world at any resolution. *arXiv preprint arXiv:2409.12191*.
- Qin Wang, Olga Fink, Luc Van Gool, and Dengxin Dai. 2022. Continual test-time domain adaptation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7201–7211.
- Jiayang Wu, Wensheng Gan, Zefeng Chen, Shicheng Wan, and Philip S. Yu. 2023. [Multimodal Large Language Models: A Survey](#). *Preprint*, arXiv:2311.13165.
- Heming Xia, Tao Ge, Peiyi Wang, Si-Qing Chen, Furu Wei, and Zhifang Sui. 2023. Speculative decoding: Exploiting speculative execution for accelerating seq2seq generation. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 3909–3925.
- Heming Xia, Yongqi Li, Jun Zhang, Cunxiao Du, and Wenjie Li. 2024a. Swift: On-the-fly self-speculative decoding for llm inference acceleration. *arXiv preprint arXiv:2410.06916*.
- Heming Xia, Zhe Yang, Qingxiu Dong, Peiyi Wang, Yongqi Li, Tao Ge, Tianyu Liu, Wenjie Li, and Zhifang Sui. 2024b. [Unlocking efficiency in large language model inference: A comprehensive survey of speculative decoding](#). *Preprint*, arXiv:2401.07851.
- Bin Xiao, Haiping Wu, Weijian Xu, Xiyang Dai, Houdong Hu, Yumao Lu, Michael Zeng, Ce Liu, and Lu Yuan. 2024. Florence-2: Advancing a unified representation for a variety of vision tasks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4818–4829.
- Sen Yang, Shujian Huang, Xinyu Dai, and Jiajun Chen. 2024. Multi-Candidate Speculative Decoding. *arXiv preprint arXiv:2401.06706*.
- Shukang Yin, Chaoyou Fu, Sirui Zhao, Ke Li, Xing Sun, Tong Xu, and Enhong Chen. 2024. [A Survey on Multimodal Large Language Models](#). *Preprint*, arXiv:2306.13549.

- Weihao Yu, Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin Lin, Zicheng Liu, Xinchao Wang, and Lijuan Wang. 2023. Mm-vet: Evaluating large multimodal models for integrated capabilities. *arXiv preprint arXiv:2308.02490*.
- Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. 2023. [Sigmoid Loss for Language Image Pre-Training](#). *Preprint*, arXiv:2303.15343.
- Duzhen Zhang, Yahan Yu, Jiahua Dong, Chenxing Li, Dan Su, Chenhui Chu, and Dong Yu. 2024a. [MM-LLMs: Recent Advances in MultiModal Large Language Models](#). *Preprint*, arXiv:2401.13601.
- Kai Zhang, Lingbo Mo, Wenhui Chen, Huan Sun, and Yu Su. 2024b. Magicbrush: A manually annotated dataset for instruction-guided image editing. *Advances in Neural Information Processing Systems*, 36.
- 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](#). In *The Twelfth International Conference on Learning Representations*.
- Zhi-Hua Zhou. 2012. *Ensemble methods: foundations and algorithms*. CRC press.

## Appendix

We structure the appendix as follows.

In Section A, we present additional generalization results. In Section B, we provide implementation details. In Section C, we show qualitative examples of dynamic ensemble weights. In Section D, we report additional results with token-tree verification. In Section E, we describe the caption-drafting setup. In Section F, we evaluate the target model and provide sample outputs. In Section G, we provide supporting experiments for remarks in the main text. In Section H, we document draft-model training details. In Section I, we list dataset- and method-specific prompt formats. In Section J, we describe benchmark and OOD datasets. In Section K, we present additional runtime analyses.

### A Comprehensive Benchmarking Across Draft and Target Models

In this section, we present the full results corresponding to the main experiments shown in Tables 1 and 3 for all drafting methods considered in this paper, including both single (multimodal and text-only) and ensemble approaches (MT\* and MTCP\*), across various sizes and types of draft and target models. Specifically, we benchmark draft models (LLaVA-1.5 68M, LLaVA-1.5 160M, LLaVA-OV 68M) and target models (LLaVA-1.5 7B, LLaVA-1.5 13B, LLaVA-NeXT 7B, LLaVA-NeXT 13B). As shown in Tables 18 to 22, the findings from Sections 4, 5 and 6.1 hold across different combinations of draft and target models.

### B Implementation Details of TABED

In this section, we examine two key factors influencing the dynamic ensembling approach: the observation window length  $h$  used for computing adaptive weights (Section B.1), the weight sampling policy  $S_w$  that determines how ensemble weights are selected in our experiments (Section B.2), and the draft block length  $\gamma$  in the sequential draft phase (Section B.3). Finally, we complement our empirical ablations with a theoretically grounded variant, TABED-AdaBoost, a multi-class AdaBoost-inspired weighting scheme (Section B.4).

#### B.1 Observation Window Length $h$

The responsiveness of the dynamic weights predicted by TABED can be controlled by the observation window length  $h$ . In Algorithm 1, we

introduce an observation window of length  $h$  for the drafting method to compute dynamic weights based on information within this window. By varying  $h$ , we further examine how this approach addresses the inherent nature of autoregressive decoding, which relies solely on past context for next-token prediction. For both benchmark and OOD datasets, we explore the effect of varying  $h$  by setting it to 1, 4, 16, and to include all previous tokens. As shown in Table 5, there is no clear winner among the different window lengths  $h$ . Therefore, by default, using all previous information (i.e., ALL) before the current timestep is a reasonable choice.

#### B.2 Weight Sampling Policy $S_w$

There are numerous variants of TABED, obtained by varying the weight sampling policy  $S_w$  in Algorithm 1. Here, we describe the representative policy used in our experiments throughout the paper. For  $S_W$  in MT\*, we sample  $W_t = [w_t^j]$ , where  $w_t^j = [1 - \frac{j}{n}, \frac{j}{n}] \in \mathbb{R}^m$  and  $n = 10$ , to succinctly represent linear combinations with non-negative integer coefficients. For MTCP\*, we sample  $W_t = [w_t]$ , where  $w_t = \text{softmax}([1/e_t^{(1)}, \dots, 1/e_t^{(m)}]; \tau) \in \mathbb{R}^m$  using soft labels, temperature  $\tau$ , and  $n = 1$ , offering an efficient alternative as  $m$  increases. We set temperature  $\tau = 1$  by default.

#### B.3 Draft block length $\gamma$

To evaluate the effect of draft block length  $\gamma$ , we conducted additional experiments with draft lengths of 3 and 7, extending our baseline setting of  $\gamma = 5$ . As shown in Tables 6 and 7, the overall trends remain consistent across different values of  $\gamma$ , demonstrating the robustness of our method to the choice of draft block length. Specifically, the dynamic ensemble drafting methods, TABED<sup>MT</sup> and TABED<sup>MTCP</sup>, consistently outperform or closely match the best-performing baselines across the most of tested configurations, confirming that our findings are robust to this hyperparameter.

#### B.4 TABED-AdaBoost: A Theoretically Grounded Variant

To complement our empirical ablations with a more theoretically grounded alternative, we study a multi-class AdaBoost-inspired weighting scheme (TABED-AdaBoost), which assigns a fixed ensemble weight to each drafting method based on its error rate. Concretely, each drafting method  $m$  is

Table 5: Experimental results with different observation window length  $h$ , which constrain the number of previous timesteps used to compute adaptive weights.  $h = 1$  indicates that the weights are based only on the ground truth from the previous decoding step.

Draft Model			Benchmark Datasets							OOD Datasets	
Type	Size	$h$	LLaVA-W	DocVQA	POPE	MMVet	IEEdit	MB	Spot	PSV	VIST
LLaVA 1.5	68M	1	2.26	2.15	2.54	2.21	2.30	2.36	2.35	2.00	2.02
		4	2.25	2.16	2.54	2.20	2.29	2.40	2.36	2.05	2.05
		16	2.26	2.17	2.53	2.21	2.31	2.42	2.36	2.04	2.05
		ALL	2.26	2.16	2.52	2.21	2.29	2.39	2.36	2.02	2.04

assigned

$$w_m = \ln\left(\frac{1}{\epsilon_m} - 1\right) + C \cdot \ln(K - 1), \quad (2)$$

where  $\epsilon_m$  denotes the method’s error rate,  $K$  is the vocabulary size, and  $C$  is a constant.

Table 8 compares TABED<sup>MT</sup>-AdaBoost with the original TABED<sup>MT</sup>. Across five of the seven in-domain benchmarks, TABED<sup>MT</sup> performs comparably to TABED<sup>MT</sup>-AdaBoost, while achieving better robustness on OOD datasets (PororoSV and VIST), suggesting that TABED<sup>MT</sup>’s adaptive weighting quickly upweights the stronger drafting method.

### C More Qualitative Samples of Dynamic Ensemble Weights

In this section, we provide further qualitative examples of dynamic ensemble weights. In Fig. 4, the y-axis denotes the proportion of  $w^{(M)}$  relative to  $w^{(M)} + w^{(T)}$ , and the x-axis indicates decoding steps. In comparison to other single drafting methods, M and T ( $y = 1$  and 0), or the static ensembling method MT ( $y = 0.5$ ), adaptive ensemble weights  $w^*$  predicted by TABED (the blue graph) effectively explore the region shaded, representing the weights that lead to token acceptance by the target model.

### D Full Results for TABED with Token Tree Verification

In this section, we present the full results corresponding to the main experiments in Table 2 for all tree widths  $d$  evaluated in this paper, including both single (multimodal (M) and text-only (T)) and dynamic ensemble approaches (TABED<sup>MT</sup>). As shown in Tables 9 and 10, the findings from Sections 4 and 5 remain consistent when combined with orthogonal enhanced verification methods across different tree widths  $d$ .

## E Details for Caption Drafting

In this section, we describe various types of lightweight image captioning models that can be used for caption drafting (Section E.1). We then demonstrate that captioning model inference completes earlier than the target model’s prefilling by analyzing the captioning model’s latency (Section E.2).

### E.1 Captioning Models

**BLIP (Li et al., 2022)** A vision-language model trained on bootstrapped synthetic captions. It uses a visual transformer and the text encoder of BERT (Devlin et al., 2019) to separately encode image and text.

<https://huggingface.co/Salesforce/blip-image-captioning-base>

**BLIP-2 (Li et al., 2023a)** A vision-language model using a frozen off-the-shelf image encoder and LLM. A querying transformer trained using bootstrapped data is included for cross-modal alignment.

<https://huggingface.co/Salesforce/blip2-opt-2.7b>

**Florence-2 (Xiao et al., 2024)** A vision-language model that is instruction-trained for a variety of tasks. Its architecture consists of a single sequence-to-sequence transformer and a vision encoder.

<https://huggingface.co/microsoft/Florence-2-large-ft>

### E.2 Latency Analysis

It is important to ensure that the captioning model runs fast enough so that it does not delay drafting. In this line, we measure in Table 11 the time taken by the two captioning models, BLIP and Florence-2, to generate captions. The results demonstrate captioning completes earlier than target model’s prefilling.

## F Evaluation of Target Model

In this section, we present both quantitative results (Section F.1) and qualitative samples (Section F.2) produced by LLaVA-1.5 7B, which serves as the target model in our experimental setting.

### F.1 Quantitative Evaluation

We assess the output quality of LLaVA-1.5 7B using benchmark datasets, as shown in Table 12. The results demonstrate that the quality of LLaVA-1.5 7B is comparable to GPT-4o-mini. The injection of multimodal information benefits the target model, LLaVA-1.5 7B, from the perspective of output quality.

### F.2 Generated Samples

#### Spot The Difference [Fig. 8e]

*Q:* What has changed from the first image to the second?

*A:* The image shows a city street with a grassy area in the middle of the road. In the first image, there are several cars parked along the street, and a few people are walking around. The street appears to be relatively empty, with no traffic or significant activity. In the second image, the street is bustling with activity. There are numerous cars parked along the street, and several people are walking around. The scene is more lively and dynamic, with a higher level of activity compared to the first image.

#### MagicBrush [Fig. 8f]

*Q:* Please provide the image edit instruction that can transform the source image to the target image.

*A:* To transform the source image into the target image, follow these steps: 1. Remove the red-haired girl from the scene, as she is not part of the target image. 2. Replace the girl with a backpack, which should be placed on the stairs. 3. Add a bus in the background, making sure it is yellow and occupies the entire width of the image. 4. Place a person in the background, standing near the bus, to give a sense of scale and context. 5. Adjust the lighting and shadows to create a more realistic scene.

#### PororoSV [Fig. 8g]

*Q:* Taking the stories from the initial images into account, can you conclude the story using the elements of the final image?

*A:* The image features a group of cartoon characters, including a penguin, a fox, and a pink character, standing in a snowy environment. The charac-

ters are positioned in various locations, with some standing closer to the foreground and others further back. The scene appears to be a part of a video game or animated series, as the characters are interacting with each other and their surroundings.

#### VIST [Fig. 8h]

*Q:* With the narratives paired with the initial images, how would you conclude the story using the last picture?

*A:* The image captures a thrilling moment during a volley ball game, with a woman leaping into the air to hit the ball. She is wearing a purple shirt and is in the process of jumping to make contact with the sports ball. The scene is filled with excitement and anticipation as the players on the court watch her performance. The volley ball game is taking place on a court, and the players are actively engaged in the match.

## G Empirical Validation of Latency and Throughput Remarks

In this section, we present detailed empirical analyses that support the theoretical assumptions introduced in the main text. We first report measurements of the draft-to-target latency ratio  $T_q/T_p$  across representative LVLMM benchmarks as well as a detailed derivation of the walltime speedup (Section G.1).

We then validate the three latency invariance remarks formulated in Remarks 1 to 3 and discuss their implications for speculative decoding efficiency (Section G.2). Finally, we analyze throughput and computational cost to confirm that the additional compute introduced by TABED is negligible in practice (Section G.3).

### G.1 Empirical Measurement of Draft-to-Target Latency Ratio $T_q/T_p$

We empirically measure the draft-to-target latency ratio ( $T_q/T_p$ ) across well-known LVLMM benchmarks (Li et al., 2024c). This ratio allows us to estimate the expected walltime speedup by substituting it into Eq. (1). All measurements reported in Table 13 are averaged across datasets.

This allows us to calculate the expected speedup by inputting the draft-to-target latency ratio into Eq. (1). Specifically, the draft-to-target latency ratio between the 7B target model and the 68M draft model is measured as 0.063 Table 13. The average block efficiency of TABED is 2.29, computed as the mean of 2.31 and 2.27, and the block length  $\gamma$

used in Table 1 is set to 5. Substituting these values into Eq. (1) yields an estimated wall-time speedup of  $\text{speedup} = 2.29 / (5 \times 0.063 + 1) = 1.74$ . We report speedup in this factorized form (block efficiency, draft-to-target latency ratio, and block length) because it is more robust to measurement noise.

## G.2 Analysis of Latency in TABED

To further analyze and discuss the effects of the TABED when latency ratio  $T_q/T_p$  is small like shown in Table 13 to the latency and throughput of speculative decoding system, we formulate the remarks Remarks 1 to 3 and empirically showed that the effect is negligible.

Focusing on the case where the remarks apply to speculative decoding settings, we use LLaVA-1.5 7B and LLaVA-1.5 68M to measure  $T_p(B, S, \gamma)$  for Remark 1 and  $T_q(B, S, 1)$  for Remarks 2 and 3, respectively. Also, we set  $S$  (e.g.,  $S \leq 3k$ ) and sufficiently small  $B$  (e.g.,  $B \leq 4$ ) and  $\gamma$  (e.g.,  $\gamma \leq 10$ ) to cover all configurations used throughout the paper. All experiments are conducted on an A100 80GB GPU using the fp16 data type for the models. As a result, the decoding phase displays the following observations:

**Remark 1.** For given  $B$  and  $S$ , regardless of  $\gamma$ ,  $T(B, S, \gamma)$  remains approximately constant (e.g.,  $|\Delta T/T| < 0.05$ ).

**Remark 2.** For a given  $B$ , regardless of  $S$ ,  $T(B, S, 1)$  remains approximately constant (e.g.,  $|\Delta T/T| < 0.05$ ).

**Remark 3.** For a given  $S$ , regardless of  $B$ ,  $T(B, S, 1)$  remains approximately constant (e.g.,  $|\Delta T/T| < 0.05$ ).

Fig. 5 shows  $T_q(B, S, 1)$  in milliseconds for sequence lengths up to 3k for each batch size  $B \in \{1, 2, 3, 4\}$ . For moderate sequence lengths  $S \leq 3k$ ,  $T_q$  varies by no more than 5% for each  $B$ , which supports Remark 2. Similarly, when comparing different  $B$ s with a fixed  $S$ ,  $T_q$  varies by no more than 5%, which supports Remark 3.

Fig. 6 shows  $T_p(B, S, \gamma)$  in milliseconds for each  $\gamma \in \{1, 3, 5, 7\}$ . We test the case of  $B=1$ , which aligns with our experimental settings where the target model always performs inference on a single batch. Over the values of  $\gamma$  considered,  $T_p$  varies by no more than 5%.

## G.3 Analysis of Throughput in TABED

We further evaluate the effect of TABED from a throughput perspective rather than latency. Because  $M_q$  is substantially smaller than  $M_p$ , batch inference using  $M_q$  (68M) introduces negligible computational overhead relative to the total workload dominated by  $M_p$  (7B). Prior studies have shown that Transformer decoding is memory-bound rather than compute-bound (Chen et al., 2024a; Fu, 2024), implying that such additional compute does not meaningfully affect walltime.

Table 14 presents the FLOPs analysis supporting this claim: the additional 0.08 GFLOPs per prediction (0.17–0.09) from  $M_q$  is minimal compared to the 66.08 GFLOPs required by  $M_p$ . All measurements are obtained using a single image (576 projected image tokens) and 24 text tokens. This confirms that TABED introduces no meaningful throughput degradation in practice.

## H Training and Evaluation of Draft Models

In this section, we present a more detailed overview of our custom training procedure for the draft models (Section H.1). We then evaluate our primary draft model, LLaVA-1.5 68M, on multimodal tasks to ensure it has the capability to properly perceive multimodality, and we provide some qualitative samples from the draft model (Section H.2).

### H.1 Details of Training

**LLaVA-1.5 (Liu et al., 2024a)** The process for developing draft models with the LLaVA-1.5 (68M, 160M) training recipe was divided into two stages: pre-training and instruction fine-tuning (IFT). Pre-training focused on training the projector while the parameters of the LLM and vision encoder were frozen. During the IFT stage, visual instruction tuning was used to teach the LLM to follow multimodal instructions. The vision encoder remained frozen throughout both stages. The hyperparameters used for each stage are described in Table 15. We trained the draft model using datasets curated by the original author of LLaVA-1.5. For more training details, see <https://github.com/haotian-liu/LLaVA/tree/main>.

**LLaVA-OneVision (Li et al., 2024a)** The development of draft models using the LLaVA-OneVision (LLaVA-OV) training recipe was divided into three stages: language-image alignment,

high-quality knowledge learning, and visual instruction tuning. In the language-image alignment stage, visual features were aligned with the word embedding space of LLMs. High-quality knowledge learning balanced computational efficiency with the integration of new knowledge into LVLMs. Visual instruction tuning consisted of two phases: (i) Single-Image Training, where the model learned to perform visual tasks using instructions from single images, and (ii) OneVision Training, where the model learned to execute multi-image visual tasks using a blend of video, single-image, and multi-image data. During the language-image alignment stage, only the projector for aligning visual features was updated, whereas all components including LLM were updated in the following three stages. We trained the draft model using datasets curated by the original author of LLaVA-OV (Li et al., 2024a). The hyperparameters used for each stage are described in Table 16, and the learning rate for the vision encoder is one-fifth of that for the LLM across all stages. For more details, visit <https://github.com/LLaVA-VL/LLaVA-NeXT>.

## H.2 Evaluation Results

Table 17 presents the evaluation results of our primary draft model, LLaVA-1.5 68M, on the OCRBench (Liu et al., 2024d) and TextCaps (Sidorov et al., 2020) datasets. We assess the output quality of the draft model with and without image inputs and compare the results with those of the target model, LLaVA-1.5 7B. In terms of output quality, the draft model with image inputs consistently outperforms the one without, illustrating that the injection of multimodal information benefits the custom-trained draft model.

Fig. 7 presents qualitative samples from the OCRBench dataset. Both LLaVA-1.5 7B and 68M models provided accurate responses, whereas the text-only LLaVA-1.5 68M model failed to answer correctly due to its lack of image-processing capabilities.

## I Prompts for Each Dataset and Drafting

In this section, we describe the formats of prompts used for inference on each dataset, including system prompts and how to organize prompts with text and image inputs (Section I.1). We then provide details on replacing image tokens in text-only and caption drafting (Section I.2).

### I.1 Prompt Formats for Each Dataset

We use the following prompt formats for their respective tasks. Based on the template for chat (*USER:* and *ASSISTANT:*), each system prompt is prepended with the start token  $\langle s \rangle$ . The  $\langle image \rangle$  token is used to represent image data within a prompt.  $[QUESTION]$  and  $[CAPTION]$  are placeholders denoting information unique to each sample of a dataset.

**LLaVA-Bench (In-the-Wild)**  $\langle s \rangle$  *USER:*  $\langle image \rangle$  *For the following question, provide a detailed explanation of your reasoning leading to the answer. [QUESTION]* *ASSISTANT:*

**DocVQA**  $\langle s \rangle$  *USER:*  $\langle image \rangle$  *For the following question, provide a detailed explanation of your reasoning leading to the answer. [QUESTION]* *ASSISTANT:*

**POPE**  $\langle s \rangle$  *USER:*  $\langle image \rangle$  *For the following question, provide a detailed explanation of your reasoning leading to the answer. [QUESTION]* *ASSISTANT:*

**MMVet**  $\langle s \rangle$  *USER:*  $\langle image \rangle$  *For the following question, provide a detailed explanation of your reasoning leading to the answer. [QUESTION]* *ASSISTANT:*

**IEdit**  $\langle s \rangle$  *USER:* *Please provide instructions for editing the source image to match the target image. Source Image:  $\langle image \rangle$  Target Image:  $\langle image \rangle$  Instruction:* *ASSISTANT:*

**MagicBrush**  $\langle s \rangle$  *USER:* *Please provide instructions for editing the source image to match the target image. Source Image:  $\langle image \rangle$  Target Image:  $\langle image \rangle$  Instruction:* *ASSISTANT:*

**Spot The Difference**  $\langle s \rangle$  *USER:* *Explain the disparities between the first and second image.  $\langle image \rangle$   $\langle image \rangle$  Difference:* *ASSISTANT:*

**PororoSV**  $\langle s \rangle$  *USER:* *Given the progression of the story with the first few images, can you write a fitting end considering the last image?  $\langle image \rangle$  Caption #1:  $[CAPTION]$   $\langle image \rangle$  Caption #2:  $[CAPTION]$ .  $\langle image \rangle$  Caption #3:  $[CAPTION]$   $\langle image \rangle$  Caption #4:  $[CAPTION]$   $\langle image \rangle$  Caption #5:  $[CAPTION]$*  *ASSISTANT:*

**VIST**  $\langle s \rangle$  *USER:* *With the narratives paired with the initial images, how would you conclude the story using the last picture?  $\langle image \rangle$  Caption #1:  $[CAPTION]$   $\langle image \rangle$  Caption #2:  $[CAPTION]$ .*

<image> Caption #3: [CAPTION] <image> Caption #4: [CAPTION] <image> Caption #5: ASSISTANT:

## I.2 Replacing Image tokens in Draftings

For text-only drafting, the <image> token is replaced by the escape character “\n”. We experimented with several replacement methods: (1) tokenizing the string “<image>” into three tokens, and (2) retaining the special token <image> without replacing it with an image embedding. Method (2) resulted in very poor block efficiency, but method (1) showed comparable block efficiency. Our replacement approach is simple because it ensures that the prompt length remains consistent before and after replacement.

For caption drafting, the <image> token is replaced by a generated caption with a prefix. Specifically, after the lightweight captioning model generates a caption based on the image inputs in the sample, we prepend the string “image: ” to the caption and replace the <image> token.

## J Details of Each Dataset

In this section, we describe each of the curated datasets in benchmark (Section J.1) and OOD (Section J.2) datasets and provide links to them for convenience and reproducibility.

### J.1 Benchmark Datasets

#### LLaVA-Bench (In-the-Wild) (Liu et al., 2023)

A dataset for comparing the performance of vision-language models against state-of-the-art proprietary models. Each prompt is provided with an image, a caption, and a reference answer from text-only GPT-4. Prompt styles include question answering, image description, and complex reasoning. In total, the dataset contains 60 unique prompts and 24 unique images.

<https://huggingface.co/datasets/lmms-lab/llava-bench-in-the-wild>

**DocVQA (Mathew et al., 2021)** A dataset designed for visual question answering on document images, comprising 50,000 questions over 12,000+ diverse document images.

<https://huggingface.co/datasets/lmms-lab/LMMS-Eval-Lite>

**POPE (Li et al., 2023b)** A multimodal question answering dataset that asks binary (yes or no) questions about whether certain objects are present in

an image. The subset used for evaluation in our work contains 100 pairs of images and questions.

<https://huggingface.co/datasets/lmms-lab/LMMS-Eval-Lite>

**MMVet (Yu et al., 2023)** A benchmark designed to evaluate large multimodal models on complex tasks, focusing on the integration of six core vision-language capabilities: recognition, OCR, knowledge, language generation, spatial awareness, and math.

<https://huggingface.co/datasets/lmms-lab/MMVet>

#### Spot the Difference (Jhamtani and Berg-Kirkpatrick, 2018)

A dataset of crowd-sourced descriptions of differences between a pair of images. The subset used for evaluation in our work contains 100 annotated image pairs collected using individual frames of security-footage data. This dataset is available under the Creative Commons Attribution 4.0 International license.

<https://huggingface.co/datasets/lmms-lab/LLaVA-NeXT-Interleave-Bench>

**IEEdit (Tan et al., 2019)** A dataset to train models to describe the relationship between images via editing instructions. The subset used for evaluation in our work contains 100 image pairs of a source image and a target image, accompanied by instructions on how to transform the source image into the target. This dataset is available under the Creative Commons Attribution 4.0 International license.

<https://huggingface.co/datasets/lmms-lab/LLaVA-NeXT-Interleave-Bench>

**MagicBrush (Zhang et al., 2024b)** A dataset for text-guided image editing containing manually annotated editing instructions to transform one real image into another. The subset used for evaluation in our work contains 100 triplets of a source image, a target image, and editing instructions. This dataset is available under the Creative Commons Attribution 4.0 International license.

<https://huggingface.co/datasets/lmms-lab/LLaVA-NeXT-Interleave-Bench>

### J.2 OOD Datasets

**Pororo-SV (Li et al., 2019)** A dataset of stories each created by pairing 5 consecutive frames from the animated series *Pororo* with a text description. The subset used for evaluation in our work contains 100 stories. This dataset is available under the

Creative Commons Attribution 4.0 International license.

<https://huggingface.co/datasets/Imms-lab/LLaVA-NeXT-Interleave-Bench>

**VIST (Huang et al., 2016)** A dataset of sequential images paired with three types of descriptions ranging from isolated factual descriptions to causal, narrative interpretations. The subset used for evaluation in our work contains 100 sequences of 3 images. This dataset is available under the Creative Commons Attribution 4.0 International license.

<https://huggingface.co/datasets/Imms-lab/LLaVA-NeXT-Interleave-Bench>

### K Time Analysis of LVLM Inference Stages

In this section, we present the time analysis of the inference stages of LVLMs. Specifically, Fig. 9 analyzes how the number of input images affects the LVLM inference time, we select ChartQA (Masry et al., 2022), Spot the Difference (Jhamtani and Berg-Kirkpatrick, 2018), and PororoSV (Li et al., 2019) datasets representing 1, 2, and 5 images with corresponding visual context lengths of 0.6k, 1.2k, and 3k, respectively. We visualize the generation time by component in Fig. 9 with 100 generated tokens for analysis, with actual average decoding lengths of 92, 117, and 88, respectively. The execution time of the *vision encoder* and *prefilling* stages increases in proportion with the number of input images, as each image is converted into several hundred context tokens. In contrast, the *decoding* stage shows little difference in execution time across varying numbers of input images, while dominating the total generation time. Hence, although reducing the number of visual tokens (Shang et al., 2024; Chen et al., 2024b; Lin et al., 2024) would significantly improve the efficiency of *vision encoder* and *prefilling* stages, it would have only marginal impact on the dominant *decoding* stage.

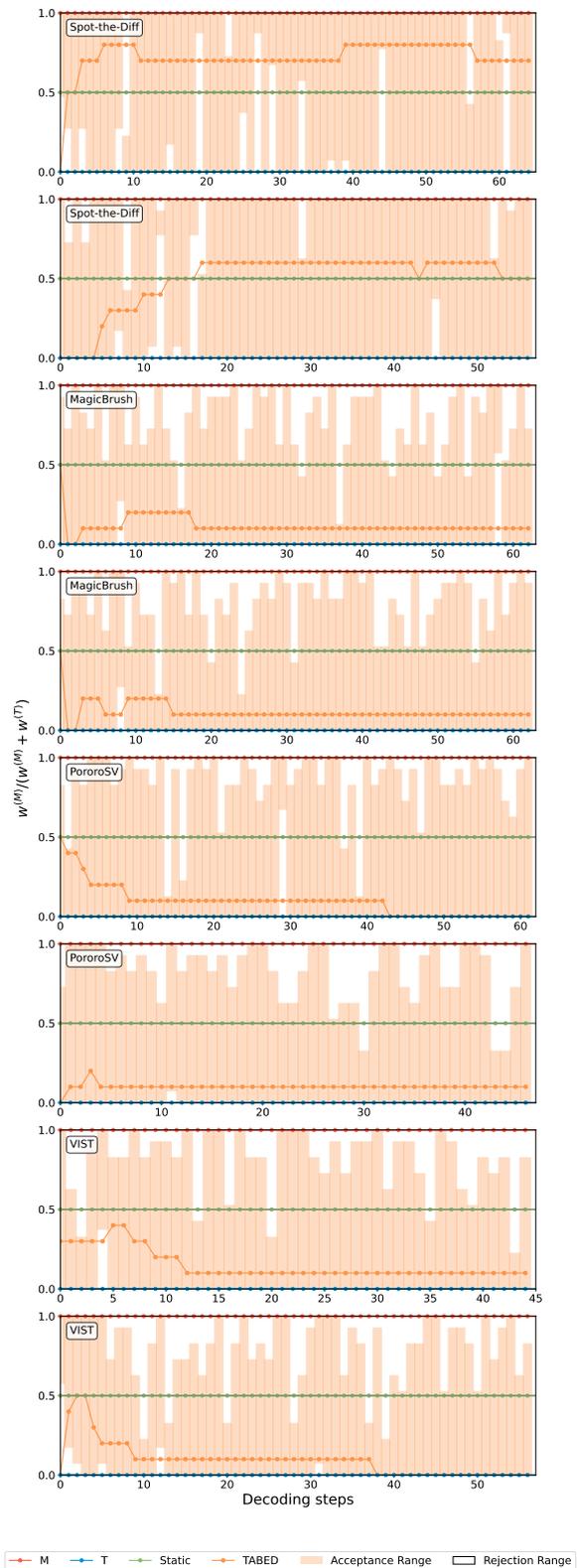


Figure 4: Additional qualitative samples of dynamic ensemble weights

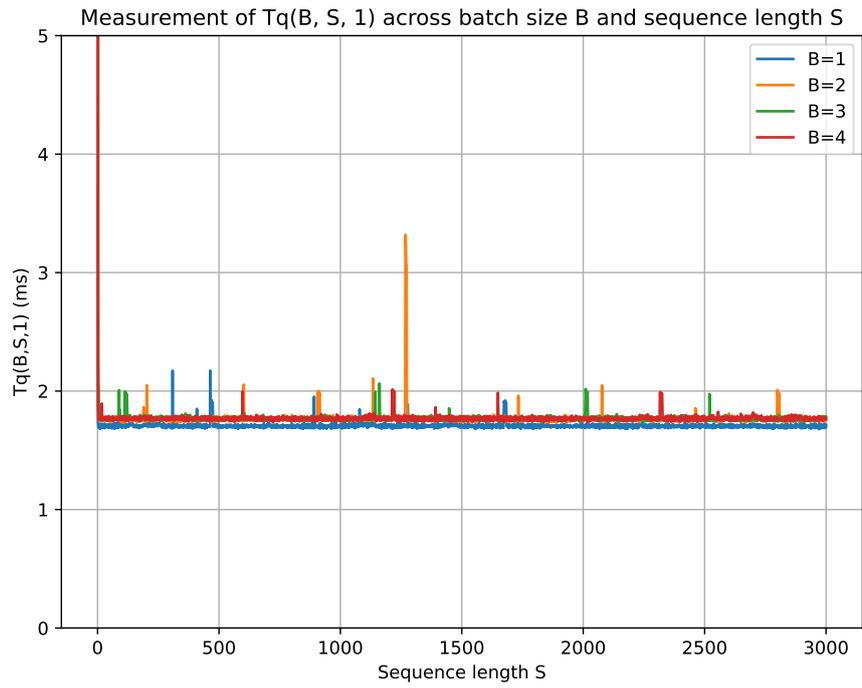


Figure 5: Empirical demonstration of Remarks 2 and 3.

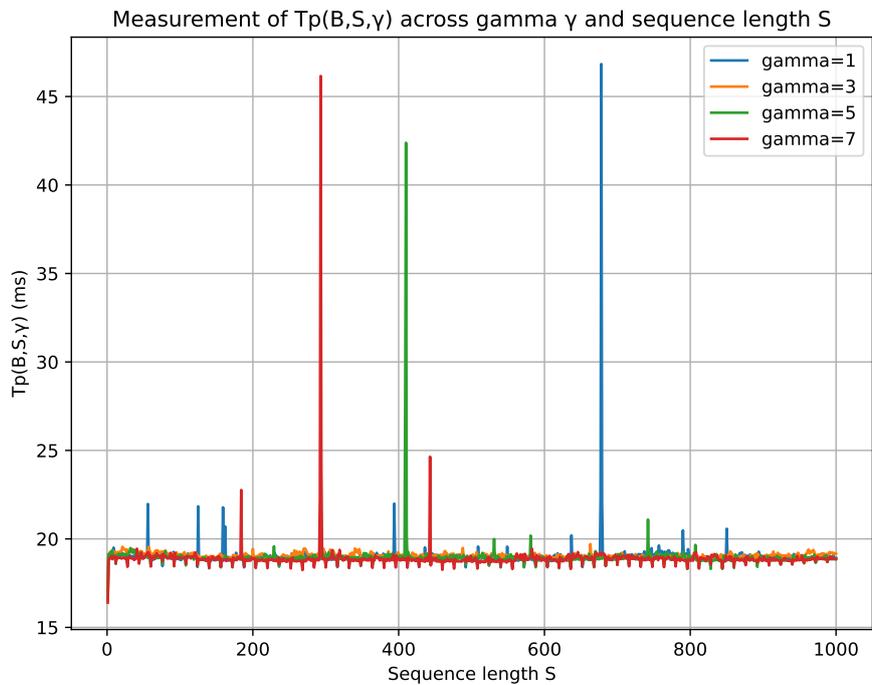


Figure 6: Empirical demonstration of Remark 1.

Table 6: Block efficiency across datasets for  $\gamma = 3$ . **Bold** and underlined denote the best and second-best results.

$\gamma = 3$ Drafting		Benchmark Datasets							OOD Datasets	
Type	Method	LLaVA-W	DocVQA	POPE	MMVet	IEdit	MB	Spot	PSV	VIST
Sgl.	M	2.09	<u>2.03</u>	<b>2.24</b>	<b>2.06</b>	2.06	1.85	2.15	1.19	1.16
	T	2.04	1.97	2.11	2.02	2.06	2.14	2.11	1.95	1.94
	C	2.08	2.01	<u>2.23</u>	2.02	<u>2.11</u>	2.17	2.12	<u>1.98</u>	<u>1.99</u>
	P	2.07	1.97	2.16	2.03	2.07	2.10	2.07	1.97	<u>1.99</u>
Ens.	<b>TABED<sup>MT</sup></b>	<b>2.10</b>	<b>2.04</b>	2.22	2.05	<u>2.11</u>	<u>2.19</u>	<b>2.17</b>	1.93	1.93
	<b>TABED<sup>MTCP</sup></b>	<b>2.10</b>	2.02	2.22	<b>2.06</b>	<b>2.14</b>	<b>2.20</b>	<b>2.17</b>	<b>2.01</b>	<b>2.03</b>

Table 7: Block efficiency across datasets for  $\gamma = 7$ . **Bold** and underlined denote the best and second-best results.

$\gamma = 7$ Drafting		Benchmark Datasets							OOD Datasets	
Type	Method	LLaVA-W	DocVQA	POPE	MMVet	IEdit	MB	Spot	PSV	VIST
Sgl.	M	2.31	<u>2.22</u>	<b>2.59</b>	<u>2.24</u>	2.26	2.01	2.41	1.19	1.16
	T	2.21	2.14	2.34	2.19	2.28	2.42	2.34	2.05	2.05
	C	2.26	2.19	<u>2.56</u>	2.21	<u>2.35</u>	<u>2.48</u>	2.37	<u>2.11</u>	<u>2.12</u>
	P	2.26	2.12	2.45	2.21	2.27	2.35	2.28	2.08	<u>2.12</u>
Ens.	<b>TABED<sup>MT</sup></b>	<b>2.32</b>	<b>2.23</b>	2.54	<u>2.24</u>	2.34	<u>2.48</u>	<u>2.44</u>	2.03	2.04
	<b>TABED<sup>MTCP</sup></b>	<b>2.32</b>	2.21	2.55	<b>2.25</b>	<b>2.39</b>	<b>2.51</b>	<b>2.47</b>	<b>2.15</b>	<b>2.19</b>

Table 8: Comparison between TABED-AdaBoost and TABED.

	LLaVA-W	DocVQA	POPE	MMVet	IEdit	MB	Spot	PSV	VIST
<b>TABED<sup>MT</sup>-AdaBoost</b>	2.28	2.12	2.48	2.23	2.26	2.28	2.33	1.96	1.87
<b>TABED<sup>MT</sup></b>	2.26	2.16	2.52	2.21	2.29	2.39	2.36	2.02	2.04

Table 9: Block efficiency across datasets for  $d = 2$ . **Bold** and underlined denote the best and second-best results.

$d = 2$ Drafting		Benchmark Datasets							OOD Datasets		
Type	Method	LLaVA-W	DocVQA	POPE	MMVet	IEdit	MB	Spot	Avg.	PSV	VIST
Sgl.	M	<b>2.93</b>	<u>2.72</u>	<b>3.48</b>	<b>2.82</b>	2.77	2.37	<u>3.13</u>	<u>2.89</u>	1.32	1.27
	T	2.82	2.64	2.92	2.76	<u>2.84</u>	<u>2.90</u>	3.04	2.85	<b>2.70</b>	<b>2.73</b>
Ens.	<b>TABED<sup>MT</sup></b>	<u>2.92</u>	<b>2.75</b>	<u>3.39</u>	<b>2.82</b>	<b>2.91</b>	<b>2.97</b>	<b>3.17</b>	<b>2.99</b>	<u>2.62</u>	<u>2.66</u>

Table 10: Block efficiency across datasets for  $d = 3$ . **Bold** and underlined denote the best and second-best results.

$d = 3$ Drafting		Benchmark Datasets							OOD Datasets		
Type	Method	LLaVA-W	DocVQA	POPE	MMVet	IEdit	MB	Spot	Avg.	PSV	VIST
Sgl.	M	<u>3.31</u>	<u>3.10</u>	<b>3.98</b>	<b>3.22</b>	3.08	2.66	<u>3.63</u>	<u>3.28</u>	1.41	1.35
	T	3.22	2.98	3.34	3.10	<u>3.20</u>	<u>3.27</u>	3.54	3.24	<b>3.16</b>	<b>3.21</b>
Ens.	<b>TABED<sup>MT</sup></b>	<b>3.32</b>	<b>3.11</b>	<u>3.87</u>	<b>3.22</b>	<b>3.25</b>	<b>3.33</b>	<b>3.64</b>	<b>3.39</b>	<u>3.03</u>	<u>3.14</u>

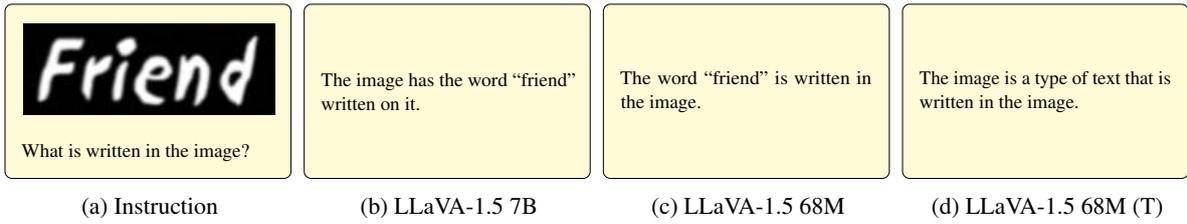


Figure 7: Qualitative evaluation samples from the OCRBench dataset by LLaVA-1.5 7B and 68M. Both the target (b) and the draft (c) models recognize the text “friend” written on the image by multimodal reasoning whereas the text-only model (d) fails, as expected.

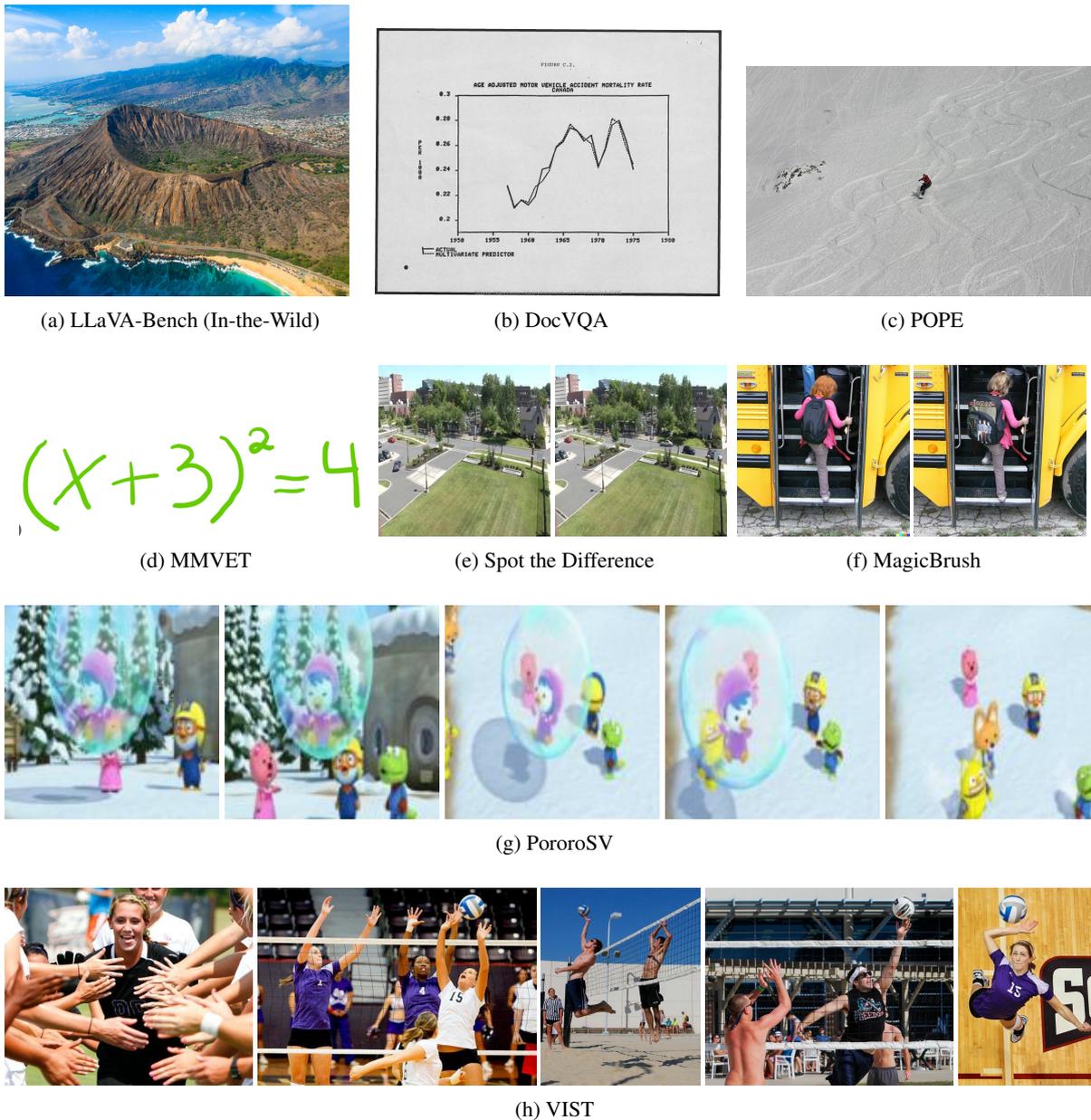


Figure 8: Qualitative image samples of benchmark and OOD datasets. The corresponding questions and answers are presented in Section F.

Model	Type	Latency (s)		
		$n = 1$	$n = 2$	$n = 5$
Target LVLM (prefilling)	LLaVA-1.5 7B	0.112	0.207	0.540
Image Captioning	BLIP	<b>0.054</b>	<b>0.055</b>	<b>0.074</b>
	Florence-2	<u>0.105</u>	<u>0.149</u>	<u>0.292</u>

Table 11: Latency analysis of image captioning models. BLIP and Florence-2 captioning latencies are lower than the target LVLM’s prefilling latency. Parallel processing can therefore hide captioning latency without affecting time to first token.

Model	Datasets	
	Spot-the-Diff	MagicBrush
GPT-4o	19.1	10.5
GPT-4o-mini	12.0	9.1
LLaVA-1.5 7B	13.9	9.6
LLaVA-1.5 7B (text-only)	11.1	8.4

Table 12: Evaluation results for the target model on benchmark datasets. Our target model, LLaVA-1.5 7B, demonstrates performance comparable to GPT-4o-mini. The injection of multimodal information in the input scenarios of the benchmark datasets benefits the target model, LLaVA-1.5 7B.

$M_p$ Size	$M_q$ Size	draft-to-target latency ratio $T_q/T_p$
7B	68M	0.063
	160M	0.206
13B	68M	<b>0.042</b>
	160M	0.137

Table 13: Empirical measurements of the draft-to-target latency ratio ( $T_q/T_p$ ), covering all model sizes. All measurements are averaged across datasets.

Table 14: Computational cost (GFLOPs) of the target model ( $M_p$ ) and the draft model ( $M_q$ ) under different drafting methods (Multimodal and TABED<sup>MT</sup>).

Model	$M_p$	$M_q$ (Multimodal)	$M_q$ (TABED <sup>MT</sup> )
<b>GFLOPs</b>	66.08	0.09	0.17

Hyperparameter	Value	Hyperparameter	Value
Training Epochs	1	Training Epochs	1
Batch Size	256	Batch Size	128
Learning Rate (LR)	1e-3	Learning Rate (LR)	2e-5
LR Schedule Type	Cosine	LR Schedule Type	Cosine
Warm-up Ratio	0.03	Warm-up Ratio	0.03
Weight Decay	0.0	Weight Decay	0.0

(a) Pretraining stage

(b) Instruction fine-tuning stage

Table 15: Details of hyperparameters used in LLaVA-1.5 training

Hyperparameter	Value	Hyperparameter	Value
Training Epochs	1	Training Epochs	1
Batch Size	512	Batch Size	512
Learning Rate (LR)	1e-3	Learning Rate (LR)	1e-5
LR Schedule Type	Cosine	LR Schedule Type	Cosine
Warm-up Ratio	0.03	Warm-up Ratio	0.03
Weight Decay	0.0	Weight Decay	0.0

(a) Image alignment stage

(b) High-quality knowledge learning stage

Hyperparameter	Value	Hyperparameter	Value
Training Epochs	1	Training Epochs	1
Batch Size	512	Batch Size	512
Learning Rate (LR)	1e-5	Learning Rate (LR)	1e-5
LR Schedule Type	Cosine	LR Schedule Type	Cosine
Warm-up Ratio	0.03	Warm-up Ratio	0.03
Weight Decay	0.0	Weight Decay	0.0

(c) Visual instruction tuning stage: Single-image training

(d) Visual instruction tuning stage: OneVision training

Table 16: Details of hyperparameters used in LLaVA-OV training

Model	OCRBench		TextCaps	
	Accuracy	METEOR	ROUGE	
LLaVA-1.5 7B	0.207	0.249	0.480	
LLaVA-1.5 68M	0.048	0.133	0.254	
LLaVA-1.5 68M (text-only)	0.014	0.064	0.132	

Table 17: Evaluation results for the off-the-shelf target model and the custom-trained draft model on MLLM tasks. Since the draft model is trained to perceive multimodality, the injection of multimodal information benefits the custom-trained draft model.

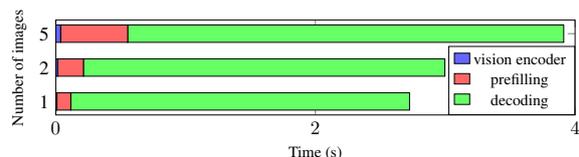


Figure 9: Inference time analysis for the LLaVA-1.5 7B model. Although the time for vision encoder and prefilling increases with the number of images, the decoding stage still dominates.

Draft Model			Benchmark Datasets (First Turn)								OOD Datasets	
Type	Size	Method	LLaVA-W	DocVQA	POPE	MMVet	IEdit	MB	Spot	Avg.	PSV	VIST
LLaVA-1.5	68M	M	2.28	2.15	2.56	2.21	2.19	1.96	2.34	2.24	1.19	1.16
		T	2.19	2.08	2.31	2.16	2.23	2.34	2.27	2.23	2.05	2.05
		C	2.22	2.15	2.50	2.17	2.29	2.36	2.31	2.29	2.08	2.10
		P	2.22	2.08	2.42	2.17	2.22	2.27	2.23	2.23	2.07	2.09
		<b>MT</b>	2.25	2.15	2.47	2.21	2.31	2.37	2.40	2.31	1.94	1.91
		<b>MT*</b>	2.26	2.16	2.52	2.21	2.29	2.39	2.36	2.31	2.02	2.04
		<b>MTCP</b>	2.27	2.15	2.49	2.22	2.32	2.39	2.40	2.32	2.10	2.13
		<b>MTCP*</b>	2.29	2.17	2.51	2.22	2.30	2.40	2.35	2.32	2.11	2.14
Draft Model			Benchmark Datasets (Second Turn)								NLP Datasets	
Type	Size	Method	LLaVA-W	DocVQA	POPE	MMVet	IEdit	MB	Spot	Avg.	NQ	GSM8K
LLaVA-1.5	68M	M	2.1	1.96	2.78	2.18	1.61	1.53	1.83	2.00	1.98	2.25
		T	2.32	2.23	2.91	2.56	1.87	2.01	2.08	2.28	2.03	2.30
		C	2.28	2.29	2.94	2.58	1.87	2.01	2.10	2.30	2.03	2.29
		P	2.28	2.09	2.95	2.49	1.86	2.00	2.08	2.25	2.01	2.26
		<b>MT</b>	2.29	2.24	2.93	2.54	1.81	1.91	2.02	2.25	2.02	2.28
		<b>MT*</b>	2.29	2.23	2.93	2.56	1.85	1.99	2.05	2.27	2.03	2.29
		<b>MTCP</b>	2.32	2.29	3.03	2.59	1.86	2.00	2.08	2.31	2.02	2.27
		<b>MTCP*</b>	2.34	2.31	3.04	2.59	1.86	2.01	2.09	2.32	2.02	2.28

Table 18: Full experimental results with LLaVA-1.5 68M (draft model) and LLaVA-1.5 7B (target model)

Draft Model			Benchmark Datasets (First Turn)								OOD Datasets	
Type	Size	Method	LLaVA-W	DocVQA	POPE	MMVet	IEdit	MB	Spot	Avg.	PSV	VIST
LLaVA-1.5	68M	M	2.01	1.96	2.58	2.14	1.89	1.77	2.05	2.06	1.19	1.15
		T	1.97	1.92	2.27	2.11	1.85	2.12	2.03	2.04	1.79	1.90
		C	2.00	1.97	2.48	2.14	1.92	2.16	2.04	2.10	1.78	1.93
		P	2.00	1.90	2.41	2.11	1.87	2.05	1.97	2.04	1.78	1.92
		<b>MT</b>	2.02	1.97	2.47	2.14	1.94	2.14	2.08	2.11	1.75	1.80
		<b>MT*</b>	2.02	1.97	2.50	2.14	1.94	2.15	2.07	2.11	1.79	1.89
		<b>MTCP</b>	2.02	1.99	2.50	2.15	1.94	2.18	2.09	2.12	1.81	1.93
		<b>MTCP*</b>	2.03	1.99	2.52	2.15	1.94	2.19	2.08	2.13	1.82	1.95
Draft Model			Benchmark Datasets (Second Turn)								NLP Datasets	
Type	Size	Method	LLaVA-W	DocVQA	POPE	MMVet	IEdit	MB	Spot	Avg.	NQ	GSM8K
LLaVA-1.5	68M	M	1.88	1.88	2.78	2.08	1.52	1.43	1.67	1.89	1.94	2.12
		T	2.05	2.09	2.84	2.45	1.76	1.83	1.93	2.14	1.99	2.17
		C	2.07	2.12	2.92	2.48	1.75	1.84	1.94	2.16	1.99	2.17
		P	2.04	1.97	2.89	2.39	1.75	1.82	1.93	2.11	1.96	2.13
		<b>MT</b>	2.03	2.08	2.88	2.41	1.69	1.75	1.86	2.10	1.97	2.15
		<b>MT*</b>	2.04	2.10	2.88	2.42	1.71	1.78	1.87	2.11	1.97	2.15
		<b>MTCP</b>	2.06	2.13	2.95	2.46	1.74	1.82	1.94	2.16	1.98	2.16
		<b>MTCP*</b>	2.07	2.14	2.97	2.47	1.75	1.83	1.95	2.17	1.98	2.16

Table 19: Full experimental results with LLaVA-1.5 68M (draft model) and LLaVA-NeXT 7B (target model)

Draft Model			Benchmark Datasets (First Turn)								OOD Datasets	
Type	Size	Method	LLaVA-W	DocVQA	POPE	MMVet	IEdit	MB	Spot	Avg.	PSV	VIST
LLaVA-1.5	160M	M	2.29	2.29	3.06	2.44	2.17	2.04	2.26	2.36	1.23	1.24
		T	2.25	2.19	2.56	2.36	2.20	2.41	2.24	2.32	1.96	2.08
		C	2.26	2.26	2.87	2.42	2.26	2.46	2.29	2.40	2.05	2.19
		P	2.29	2.23	2.94	2.42	2.27	2.49	2.32	2.42	2.08	2.20
		<b>MT</b>	2.31	2.30	2.95	2.44	2.24	2.46	2.32	2.43	1.89	1.97
		<b>MT*</b>	2.26	2.31	3.01	2.44	2.19	2.49	2.31	2.43	1.95	2.07
		<b>MTCP</b>	2.31	2.30	3.02	2.44	2.30	2.54	2.34	2.46	2.06	2.20
		<b>MTCP*</b>	2.31	2.31	3.06	2.45	2.31	2.54	2.33	2.47	2.07	2.22
Draft Model			Benchmark Datasets (Second Turn)								NLP Datasets	
Type	Size	Method	LLaVA-W	DocVQA	POPE	MMVet	IEdit	MB	Spot	Avg.	NQ	GSM8K
LLaVA-1.5	160M	M	2.15	2.14	3.17	2.45	1.74	1.6	1.94	2.17	2.28	2.48
		T	2.37	2.31	3.04	2.72	2.11	2.14	2.33	2.43	2.3	2.55
		C	2.40	2.38	3.11	2.75	2.12	2.16	2.33	2.46	2.29	2.56
		P	2.39	2.35	3.25	2.76	2.08	2.13	2.33	2.47	2.29	2.51
		<b>MT</b>	2.37	2.37	3.21	2.77	2.04	2.06	2.23	2.44	2.29	2.54
		<b>MT*</b>	2.34	2.36	3.21	2.75	2.11	2.18	2.27	2.46	2.30	2.53
		<b>MTCP</b>	2.42	2.38	3.22	2.78	2.08	2.13	2.30	2.47	2.30	2.54
		<b>MTCP*</b>	2.42	2.39	3.22	2.79	2.09	2.14	2.31	2.48	2.30	2.54

Table 20: Full experimental results with LLaVA-1.5 160M (draft model) and LLaVA-NeXT 7B (target model)

Draft Model			Benchmark Datasets (First Turn)								OOD Datasets	
Type	Size	Method	LLaVA-W	DocVQA	POPE	MMVet	IEdit	MB	Spot	Avg.	PSV	VIST
LLaVA-1.5	160M	M	1.96	1.93	2.22	2.10	1.91	1.98	1.98	2.01	1.65	1.73
		T	1.99	1.91	2.02	2.10	1.86	2.24	1.99	2.02	1.81	1.86
		C	2.00	1.96	2.17	2.13	1.92	2.20	1.98	2.05	1.79	1.84
		P	1.99	1.91	2.13	2.10	1.96	2.19	2.02	2.04	1.87	1.91
		<b>MT</b>	2.01	1.95	2.21	2.14	1.93	2.22	2.02	2.07	1.84	1.89
		<b>MT*</b>	2.01	1.96	2.23	2.13	1.96	2.27	1.99	2.08	1.83	1.88
		<b>MTCP</b>	2.02	1.97	2.28	2.14	1.97	2.26	2.08	2.10	1.88	1.93
		<b>MTCP*</b>	2.02	1.97	2.29	2.14	1.97	2.27	2.07	2.10	1.88	1.93
Draft Model			Benchmark Datasets (Second Turn)								NLP Datasets	
Type	Size	Method	LLaVA-W	DocVQA	POPE	MMVet	IEdit	MB	Spot	Avg.	NQ	GSM8K
LLaVA-1.5	160M	M	1.97	1.96	2.37	2.11	1.63	1.64	1.77	1.92	1.87	2.16
		T	2.06	2.08	2.57	2.39	1.79	1.86	1.96	2.10	1.97	2.23
		C	2.07	2.12	2.67	2.43	1.79	1.86	1.94	2.13	1.96	2.23
		P	2.04	2.09	2.62	2.44	1.78	1.85	1.93	2.11	1.94	2.23
		<b>MT</b>	2.07	2.10	2.60	2.41	1.75	1.81	1.90	2.09	1.95	2.23
		<b>MT*</b>	2.07	2.10	2.60	2.42	1.76	1.83	1.91	2.10	1.95	2.23
		<b>MTCP</b>	2.09	2.14	2.71	2.45	1.78	1.86	1.93	2.14	1.96	2.24
		<b>MTCP*</b>	2.09	2.14	2.69	2.46	1.79	1.86	1.93	2.14	1.96	2.24

Table 21: Full experimental results with LLaVA-OV 68M (draft model) and LLaVA-NeXT 7B (target model)

Draft Model			Benchmark Datasets (First Turn)								OOD Datasets	
Type	Size	Method	LLaVA-W	DocVQA	POPE	MMVet	IEdit	MB	Spot	Avg.	PSV	VIST
LLaVA-1.5	68M	M	2.01	1.97	2.44	2.04	1.87	1.75	1.95	2.00	1.18	1.15
		T	1.98	1.93	2.23	2.02	1.86	2.07	1.94	2.00	1.73	1.89
		C	1.98	1.96	2.40	2.05	1.89	2.09	1.92	2.04	1.74	1.90
		P	1.97	1.91	2.35	2.02	1.87	1.98	1.87	2.00	1.74	1.90
		<b>MT</b>	2.00	1.97	2.36	2.06	1.91	2.07	1.99	2.05	1.69	1.79
		<b>MT*</b>	2.02	1.94	2.39	2.06	1.89	2.07	1.97	2.05	1.73	1.86
		<b>MTCP</b>	2.00	1.99	2.39	2.06	1.93	2.11	2.00	2.07	1.75	1.91
		<b>MTCP*</b>	2.01	1.99	2.40	2.06	1.93	2.11	1.99	2.07	1.75	1.93
Draft Model			Benchmark Datasets (Second Turn)								NLP Datasets	
Type	Size	Method	LLaVA-W	DocVQA	POPE	MMVet	IEdit	MB	Spot	Avg.	NQ	GSM8K
LLaVA-1.5	68M	M	1.85	1.82	2.73	2.02	1.52	1.43	1.72	1.87	1.97	2.18
		T	2.00	2.03	2.70	2.31	1.77	1.81	1.98	2.09	2.01	2.22
		C	2.03	2.08	2.77	2.35	1.77	1.81	2.00	2.12	2.01	2.23
		P	2.01	1.91	2.75	2.27	1.76	1.79	1.98	2.07	1.99	2.19
		<b>MT</b>	2.02	2.02	2.74	2.29	1.69	1.73	1.91	2.06	1.99	2.20
		<b>MT*</b>	2.01	2.01	2.79	2.30	1.74	1.78	1.95	2.08	2.00	2.22
		<b>MTCP</b>	2.04	2.07	2.82	2.32	1.74	1.79	1.97	2.11	2.00	2.21
		<b>MTCP*</b>	2.04	2.07	2.82	2.32	1.75	1.79	1.98	2.11	2.00	2.22

Table 22: Full experimental results with LLaVA-1.5 68M (draft model) and LLaVA-NeXT 13B (target model)