

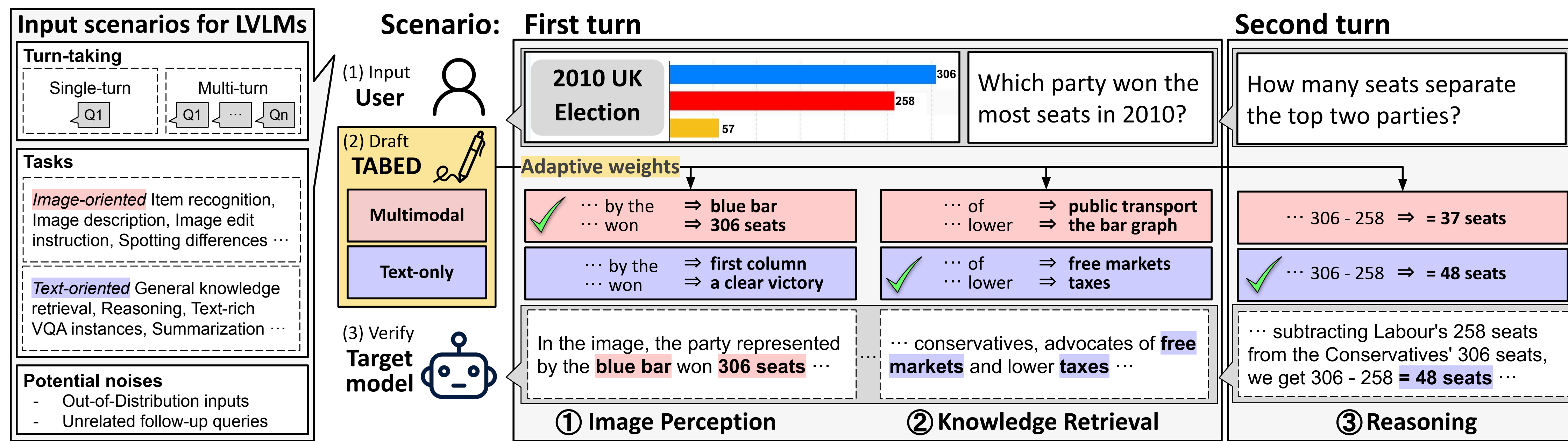
# TABED: Test-Time Adaptive Ensemble Drafting for Robust Speculative Decoding in LVLMs



Minjae Lee\*, Wonjun Kang\*, Byeongkeun Ahn, Christian Classen,  
Kevin Galim, Seunghyuk Oh, Minghao Yan, Hyung Il Koo, Kangwook Lee  
FuriosaAI, University of Wisconsin-Madison, KRAFTON

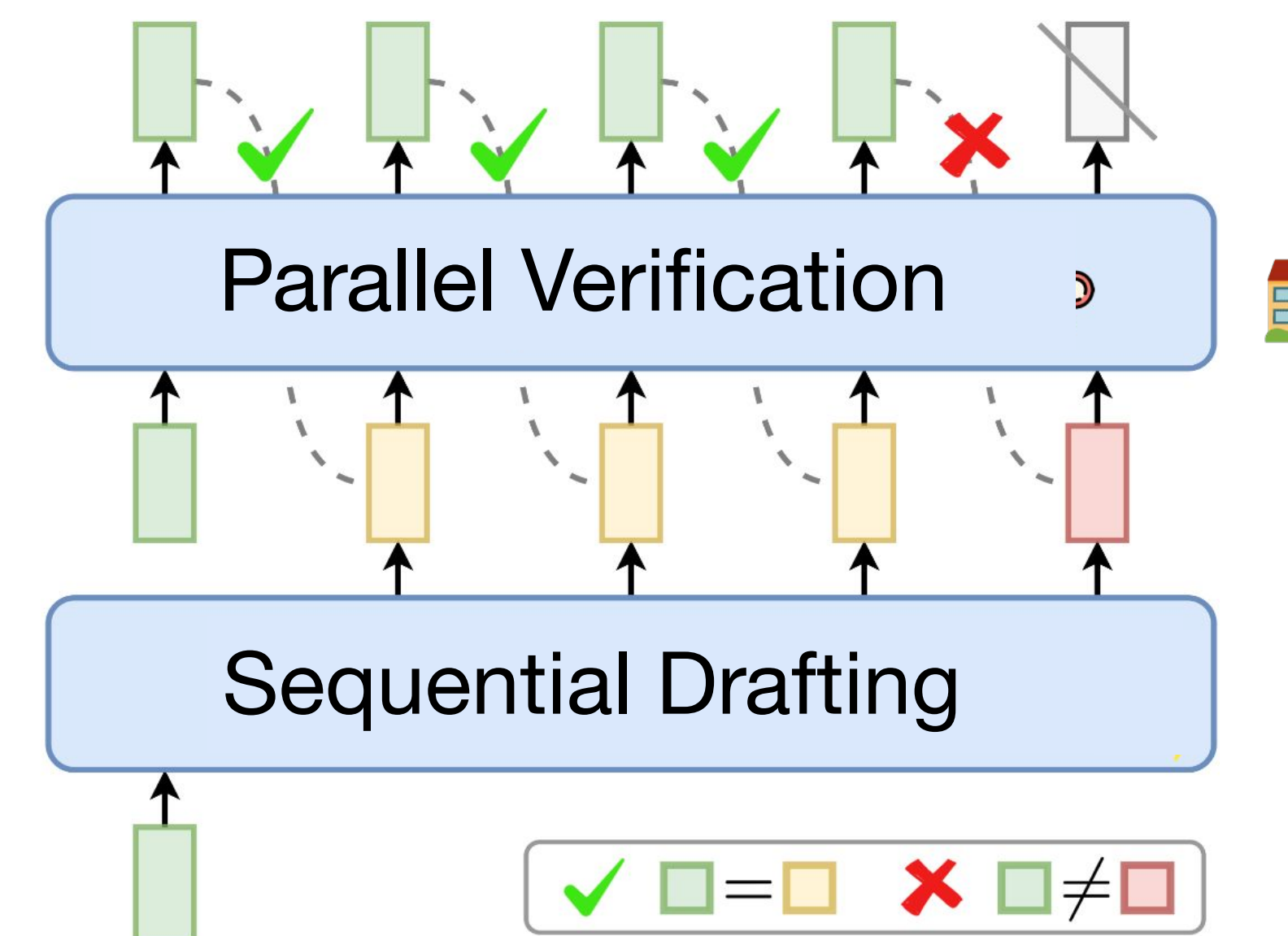
**EACL 2026**  
RABAT • MOROCCO  
Mars • March 24-29, 2026 • مارس

## Motivation



## Preliminary

### Speculative Decoding<sup>[1]</sup>



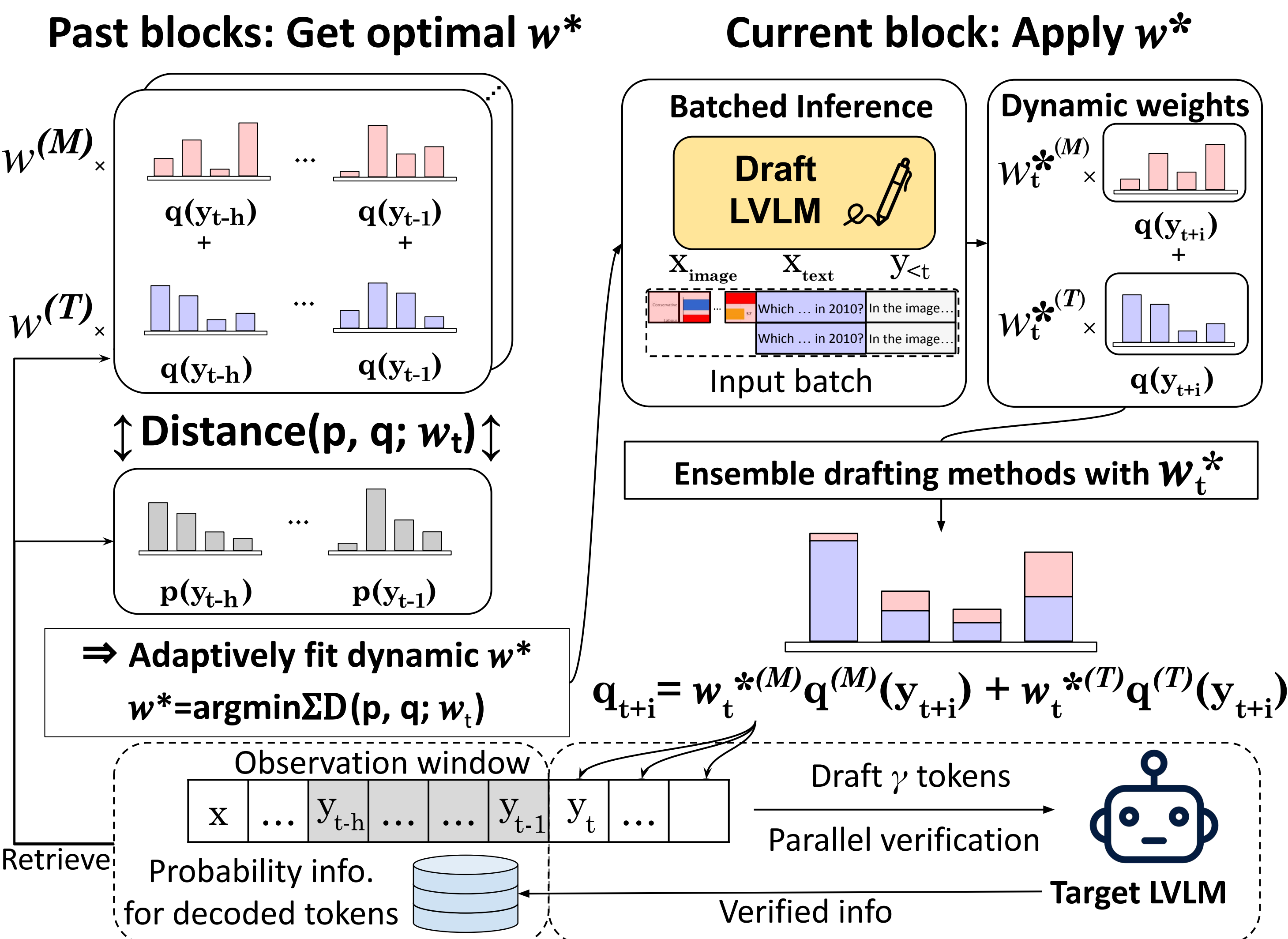
- **Large Vision-Language Models (LVLMs)** must handle **diverse input scenarios**.
- To effectively accelerate LVLMs with **Speculative Decoding (SD)**, different drafting strategies are required for intra-response and inter-response settings.
- While SD has proven effective for LLMs, it remains **underexplored for LVLMs**.

1. A small **draft model** speculates a specified number of draft tokens.
2. The larger **target model** verifies these proposed tokens in parallel.

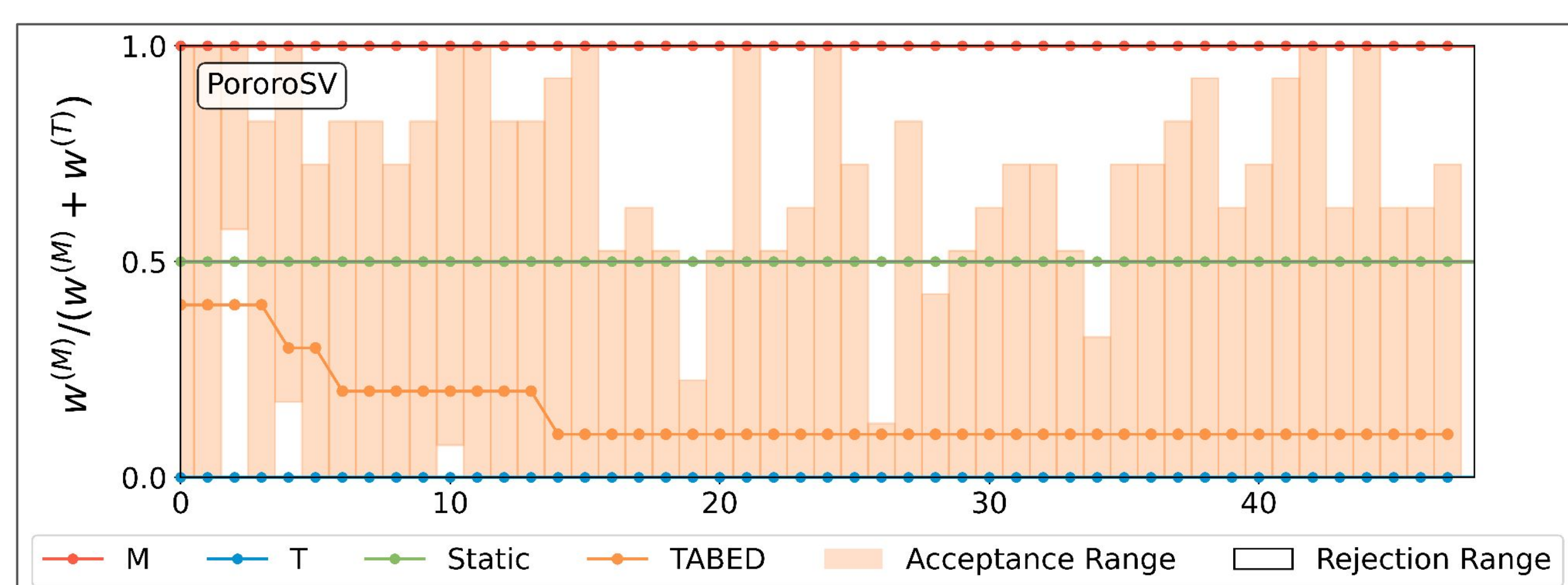
## Contribution

- We **benchmark** existing drafting methods and find that their performance fluctuates across diverse LVLM input scenarios.
- We propose **Test-time Adaptive Batched Ensemble Drafting (TABED)**, which achieves **superior and robust** performance.
- TABED is further enhanced via its **plug-and-play integration** with advanced verification and alternative drafting methods.

## Method



1. **Batched Inference:** TABED generates multiple drafts via batched inference with shared model parameters.
2. **Test-time Adaptive Ensemble Weights:** TABED dynamically ensembles drafts by leveraging deviations from past ground truths available in the SD setting.



### Adaptation Behavior of Dynamic Weights:

- TABED effectively stay within the shaded acceptance range while avoiding the unshaded rejection range

## Experimental Results

Drafting		Benchmark Datasets (First Turn)								OOD Datasets	
Type	Method	LLaVA-W	DocVQA	POPE	MMVet	IEEdit	MB	Spot	Avg.	PSV	VIST
Single	M	<b>2.28</b>	<b>2.15</b>	<b>2.56</b>	<b>2.21</b>	2.19	1.96	2.34	2.24	1.19	1.16
	T [16]	2.19	2.08	2.31	2.16	2.23	2.34	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	IEEdit	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>

### Benchmarking Results

- TABED consistently achieves either the **best** or **second-best** performance across diverse input scenarios.
- TABED achieves an average **robust walltime speedup** of 1.74× over autoregressive decoding and a 5% improvement over single drafting methods (Multimodal and Text-only).
- TABED supports **plug-and-play** compatibility with **no further training**, and maintains **negligible ensembling costs**.

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

Drafting		Benchmark		OOD
Type	Method	First	Second	
Single	M	2.24	2.00	1.18
	T [16]	2.23	2.28	2.05
	C	2.29	2.30	2.09
	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>

### Plug-and-Play Extensions

- (Top) TABED integrated with **token-tree verification** using tree width  $d$ .
- (Bottom) TABED integrated with **alternative drafting methods** (caption-based and pooled-multimodal).

