TRANSPORTING TOKENS: OPTIMAL-TRANSPORT VIEW OF PARALLEL LLM DECODING

Anonymous authorsPaper under double-blind review

000

001

002

004

006

008 009 010

011 012

013

014

015

016

017

018

019

021

023

025

026

027 028 029

030

032

033

034

037

038

040

041

042

043

044

046

047

048

051

052

ABSTRACT

Autoregressive decoding is a primary bottleneck for large language models (LLMs), as its inherent sequentiality severely limits inference speed. While speculative decoding methods mitigate this via a draft-and-verification pipeline their effectiveness is severely constrained by dependency on draft model quality and availability. We rethink the generation pattern and introduces a novel theoretical perspective by reframing token generation as a predictable state transition process in probability space, formalized through Optimal Transport (OT) theory. We demonstrate that the temporal consistency of hidden states induces a stable transport map, enabling theoretically grounded multi-step prediction. Building on this insight, we develop SHAPE, an OT-based predictor that implements lightweight Sinkhorn iterations. Extensive evaluations across diverse models (e.g., Qwen, Vicuna, LLaMA, DeepSeek) and tasks (text, code, math) show that SHAPE achieves up to $5.23\times$ speedup with minimal quality loss ($\leq 1.2\%$ accuracy drop), empirically validating our distributional transition hypothesis. This work establishes a new theoretical foundation for understanding autoregressive decoding and a practical path toward high-speed generation beyond token-wise limitations.

1 Introduction

Large Language Models (LLMs) have become the cornerstone of modern artificial intelligence, achieving remarkable success across tasks ranging from natural language understanding to text generation (13; 23; 10; 17; 20). LLMs of varying scales have been widely deployed in cloud server clusters (e.g., GPT-4 (14), Llama3 (5), and Grok1 (24)) and edge devices (e.g., the 6B-parameter GPT-3 and 7B-parameter LLaMA-2 variants as lightweight LLMs (12; 19)). With their increasing adoption in search (22) and conversational AI (15), there is a growing demand for low-latency long-sequence generation, making the optimization of effective token generation rate under constrained computational resources a critical research challenge.

Unfortunately, both cloud-based large models and edge-side small models rely on autoregressive token-by-token generation, which requires sequential computation of each token without parallelization. Additionally, the quadratic complexity of attention mechanisms with respect to context length exacerbates the issue. The standard autoregressive decoding used in existing LLMs suffers from inherent inefficiencies (21; 8)—generation time scales linearly with both context length and model size, and its sequential nature leads to cumulative latency. Our benchmarking experiments across diverse models reveal that larger model sizes and longer context lengths lead to significantly higher per-token latency. This cost is compounded by the sequential nature of decoding, highlighting the urgent need for optimization to achieve practical deployment efficiency. Comprehensive results are presented in Appendix B.

Speculative decoding (2) addresses this by introducing a fast draft model to predict multiple tokens in advance, followed by verification from the target model. However this two-stage draft-andverification paradigm still incurs sequential latency and is highly sensitive to the quality of draft models. Lookahead (3) and Medusa (2) reduce decoding time using n-gram heuristics or shallow predictors, but their limited accuracy (e.g., 0.6 for Medusa) results in suboptimal speedup. EA-GLE (11) improves draft accuracy by leveraging hidden-state features, achieving better acceleration,

yet it remains draft-model-dependent, introducing overhead and limiting scalability across diverse model configurations. CLLMs (9) accelerate decoding by directly predicting future token distributions via conditional probabilities, enabling parallel generation. However, they require fine-tuning parts of the original model, increasing training costs, and while particularly effective for mathematical reasoning, they exhibit limited stability in long-form generation.

In contrast, our approach reconceptualizes decoding itself through a distributional lens. In this work, we propose a paradigm shift by reconceptualizing token generation as a *distributional transition process*. Our key insight stems from the empirical observation that hidden states exhibit strong temporal consistency during decoding—consecutive states maintain high semantic similarity with a predictable lower bound. This regularity suggests that token generation follows a structured evolution in probability space, a perspective we formalize through optimal transport (OT) theory.

By modeling the transition between successive token distributions as a mass transport problem, where the semantic similarity between hidden states induces a stable OT map. To empirically validate this theoretical framework, we develop **SHAPE** (Step-ahead Hidden-state Accelerated Prediction Engine) as a concrete instantiation of our OT-based perspective. SHAPE operationalizes the theoretical transport maps by learning lightweight operators between hidden states, enabling parallel token prediction without auxiliary draft models. The empirical success of SHAPE—achieving substantial speedups while maintaining quality—serves as strong evidence for the correctness of our underlying theoretical insight: that token generation can indeed be understood as a predictable transport process in probability space. We evaluated SHAPE on a range of models—including Qwen, Vicuna, LLaMA, and DeepSeek—across general language (WikiText, Alpaca, MT-Bench) and reasoning-heavy tasks (MATH500, AIME24, LiveCodeBench v5). The results show that SHAPE achieves speedups of up to $5.23\times$ while maintaining output quality within a minimal margin of degradation ($\leq 1.2\%$ accuracy drop on reasoning tasks). In comparative experiments, SHAPE consistently outperforms existing acceleration methods: it surpasses EAGLE3 by $1.1\times$, Medusa-1 by $2.1\times$, and Medusa-2 by $1.6\times$ across different models and datasets.

Beyond performance, this work makes the following key contributions:

- A Novel Theoretical Foundation: We introduce a paradigm shift by reconceptualizing
 token generation as a predictable transition of probability distributions. This perspective is
 rigorously formalized through Optimal Transport theory and validated empirically, establishing a new principled understanding of decoding dynamics.
- A Practical, Plug-and-Play Predictor: We develop SHAPE, a lightweight prediction engine that operationalizes this theory. Crucially, SHAPE requires no modifications to the base LLM's parameters, offering a draft-free, plug-and-play solution for immediate deployment that significantly enhances decoding efficiency.
- Scalability to Arbitrary Future Steps. SHAPE generalizes to predict hidden states at arbitrary future time steps (e.g., t+1, t+2, t+3), providing greater flexibility for long-sequence generation tasks. This scalability supports diverse applications with varying sequence lengths and complexity.

By fundamentally rethinking the decoding process rather than optimizing within its constraints, this work opens new directions for efficient LLM inference.

2 From State Similarity to Distributional Transition

2.1 SEMANTIC SIMILARITY OF HIDDEN STATES

Building on recent work that recognizes the regularity of hidden-state sequences (11) and their utility for parallel prediction (2), we systematically analyze the temporal correlations between consecutive hidden states during autoregressive decoding. Let $\mathbf{h}_t \in \mathbb{R}^H$ denote the final-layer hidden state at decoding step t. We quantify the *semantic consistency* between states at steps t and t+n using cosine similarity:

$$\operatorname{sc}(\mathbf{h}_{t}, \mathbf{h}_{t+n}) = \frac{\mathbf{h}_{t} \cdot \mathbf{h}_{t+n}}{\|\mathbf{h}_{t}\|_{2} \cdot \|\mathbf{h}_{t+n}\|_{2}}$$
(1)

As shown in Figure 1, our analysis reveals a consistent pattern across diverse models and datasets: hidden states exhibit strong semantic consistency at offset $n{=}1$, with similarity gradually decaying yet remaining substantial for $n{=}2,3$. Crucially, we observe that there exists a task-agnostic lower bound $\tau \in (0,1)$ such that $SC(\mathbf{h}_t,\mathbf{h}_{t+1}) \geq \tau$ for the vast majority of decoding steps. This empirical finding demonstrates an inherent *temporal smoothness* in the hidden-state trajectory during generation.

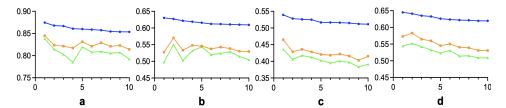


Figure 1: Semantic similarity between tokens at t and t+n across different datasets and models. Color denotes offset: blue for n=1, orange for n=2, green for n=3. Subfigures: (a) Qwen–Chinese, (b) Qwen–English, (c) Vicuna–English, (d) LLaMA–English.

The observed consistency suggests that transitions between consecutive token distributions are both small and structured. Rather than treating token generation as a sequence of discrete sampling operations, we therefore conceptualize it as a *continuous evolution of probability distributions*. This distributional perspective, formalized next using Optimal Transport, provides the theoretical foundation for predicting future decoding steps.

2.2 Modeling Token Generation with Optimal Transport

We introduce a novel perspective that reframes token generation as a *structured probability distribution transition*. Unlike the standard view of sampling discrete tokens, our approach analyzes the continuous evolution of distributions \mathbf{p}_t within the model's geometric landscape. At each step t, the distribution $\mathbf{p}_t = \operatorname{softmax}(W\mathbf{h}_t/\tau_s)$ is mapped to the token embedding space, forming a discrete measure $\mu_t = \sum_{i=1}^V \mathbf{p}_t(i) \, \delta_{E_i}$ to enable geometric comparison. The transition from μ_t to μ_{t+1} is then formulated as an *entropic-regularized optimal transport* (OT) problem, where we seek the coupling Π_t^* that minimizes transport cost—defined by the squared Euclidean distance between token embeddings—regularized by the KL divergence from the product distribution $\mathbf{p}_t\mathbf{p}_{t+1}^{\mathsf{T}}$. The row-normalized optimal coupling K_t yields a principled stochastic transition matrix such that $\mathbf{p}_{t+1} \approx K_t^{\mathsf{T}} \mathbf{p}_t$, casting generation as a path-following process in the probability simplex. Under mild assumptions, the stability of the hidden states ensures the OT distance between consecutive steps is bounded, guaranteeing the existence and uniqueness of this optimal path. This theoretical insight forms the cornerstone of the SHAPE method, offering a powerful framework for analyzing and intervening in the generation process. Complete proofs are provided in Appendix C.

3 SHAPE: AN OT-GUIDED MULTI-TOKEN PREDICTOR

To validate and operationalize the OT-based transition view, we propose **SHAPE** (Step-ahead Hidden-state Accelerated Prediction Engine), a **draft-free**, plug-and-play framework for parallel decoding. As shown in Figure 2. SHAPE consists of two key components: **Step-ahead Hidden State Prediction** and **Tree Rejection Sampling**. The core design of the framework focuses on capturing the semantic correlation of hidden states by capturing temporal features and training a predictor to approximate future hidden states. With tree reject sampling select the longest accepted prefix in parallel dynamically, so we can get α accepted token in one LLM forward to achieve parallel acceleration.

3.1 HIDDEN STATE SEMANTIC CORRELATION MODELING

The main structure of step-ahead hidden state prediction is shown in Figure 3, with three main trainable components. To first extract features in hidden state temporal modeling, the hidden states

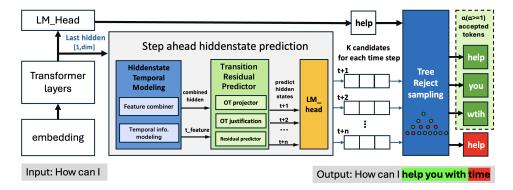


Figure 2: Illustration of the SHAPE (Step-ahead Hidden-state Accelerated Prediction Engine) framework. SHAPE leverages strong temporal correlations in hidden states to predict multiple future tokens by modeling hidden state transitions. It includes temporal modeling and residual predictors for hidden state prediction, followed by edge-to-edge LM head training to generate multiple candidates for each future step. SHAPE uses tree-based rejection sampling to select optimal token candidates at each time step, enabling efficient multi-token generation without a draft model.

from the current and previous three-time steps are concatenated and passed through a series of transformations, including linear projections, layer normalization, activation functions, and dropout. These steps capture temporal dependencies and refine the features, resulting in a final representation that effectively encodes the relationships between the time steps.

3.2 STEP AHEAD HIDDEN STATE RESIDUAL PREDICTOR

3.2.1 Predictor Construction

The high dimension of hidden states makes direct prediction challenging. To address this, we reformulate the task as modeling residual changes between consecutive hidden states. This delta-based approach incorporates an adaptive gating mechanism that dynamically scales the predicted changes based on both the original state and predicted delta. The gate network, implemented as a Linear-Sigmoid structure, controls how much of the predicted change is applied, effectively managing uncertainty and maintaining stability across multiple prediction steps.

We propose modeling the hidden state transition as an optional optimal transport step. When enabled, the transition between hidden states $H_t, H_{t+n} \in \mathbb{R}^H$ is processed through three key components:

Dimensionality Reduction: A learned linear projection $P_1 : \mathbb{R}^H \to \mathbb{R}^d$:

$$h_t^d = P_1(H_t), \quad h_{t+n}^d = P_1(H_{t+n})$$
 (2)

Optimal Transport: Converting low-dimensional representations to probability distributions through softmax:

$$p = \operatorname{softmax}(h_t^d), \quad q = \operatorname{softmax}(h_{t+n}^d) \tag{3}$$

followed by solving the transport problem:

$$\min_{T} \langle T, C \rangle + \varepsilon H(T) \quad \text{s.t.} \quad T\mathbf{1} = p, \quad T^{\top}\mathbf{1} = q \tag{4}$$

The cost matrix C quantifies semantic differences between source and target hidden states, ensuring high-activation neurons align with high-probability regions when minimized with the transport matrix. The entropy term H(T) prevents overly sparse solutions by encouraging uniformity in the transport matrix, enhancing robustness to noise. The marginal constraints (p,q) assume uniform distribution across dimensions, treating each dimension equally. The regularization coefficient ϵ balances alignment accuracy with computational efficiency, where mid-range values achieve low perplexity while maintaining moderate sparsity.

Dimension Recovery: A learned inverse projection $P_2: \mathbb{R}^d \to \mathbb{R}^H$ to recover aligned states:

$$H_{t+n}^{\mathrm{OT}} = P_2(T^{\top} \mathbf{1}) \tag{5}$$

The final hidden state is obtained through a weighted combination of the raw prediction and the OT-aligned state:

$$H_{t+n} = (1 - \alpha)H_{t+n}^{\text{raw}} + \alpha H_{t+n}^{\text{OT}}$$

$$\tag{6}$$

where H_{t+n}^{raw} is the direct residual prediction output without OT justification, $\alpha \in [0,1]$ controls the influence of the OT alignment.

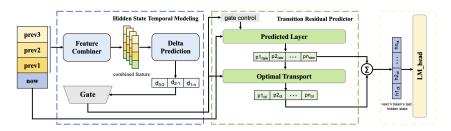


Figure 3: Step-ahead hidden state predictor: a three-part architecture comprising hidden state temporal modeling (blue), transmission residual predictor (green), and edge-to-edge LM head (yellow).

3.2.2 PREDICTOR TRAINING

The hidden state predictor architecture is designed to maintain dimensional consistency with the source large language model, preserving the original hidden state dimensionality. The training procedure utilizes optimal transport learning ($\alpha=0.5,\,\epsilon=0.1$) to enhance multi-step prediction accuracy. The training corpus consists of preprocessed hidden states extracted from both English ShareGPT conversational data and Chinese THUC_News articles, enabling bilingual prediction capabilities. The training was conducted using AdamW optimization with mixed-precision computation, incorporating uniform noise augmentation (std = 0.2) to improve model robustness. Input sequences were truncated at 2048 tokens to maintain computational efficiency with batch size = 16. The training objective combined two loss terms:

Hidden State Loss This loss optimizes the consistency between predicted hidden states $\hat{h}t + n$ and target hidden states ht + n using mean squared error:

$$\mathcal{L}_{\text{hidden}} = \frac{1}{N} \sum_{i=1}^{N} \left\| \hat{h}_{t+n}^{i} - h_{t+n}^{i} \right\|_{2}^{2}$$
 (7)

where N is the sample size.

Token Distribution Loss This cross-entropy loss ensures alignment between predicted and target token distributions:

$$\mathcal{L}_{\text{token}} = -\frac{1}{N} \sum_{i=1}^{N} \sum_{j=1}^{V} p_{\text{target}}^{i}(j) \log p_{\text{output}}^{i}(j)$$
 (8)

where $p_{\text{target}}(j)$ and $p_{\text{output}}(j)$ represent the target and predicted token distributions respectively, and V is the vocabulary size.

3.3 TREE REJECT SAMPLING

Tree Rejection Sampling generates multiple candidate paths for the next N tokens at time step t, forming a tree structure of width k and depth N (thus producing k^N candidate paths). The model then computes the joint probabilities of these paths in parallel. Low-probability paths are rejected based on a predefined acceptance threshold, and the remaining paths are merged by selecting the longest valid prefix. This design balances generation diversity and quality by exploring multiple future branches in a single forward pass. Detailed algorithm implementation is shown in Appendix.

4 EXPERIMENTS

We conduct comprehensive evaluations of SHAPE across multiple mainstream LLMs and diverse benchmarks to assess both efficiency and generation quality. Our evaluation covers a range of model families including traditional generation models: Vicuna (7B/13B), LLaMA2-Chat (7B/13B), Qwen (7B/14B), and more recent long-COT such as Qwen3 and DeepSeek-R1. To thoroughly evaluate SHAPE's effectiveness, we employ multiple benchmark categories. For general text generation, we use Alpaca, WikiLingua, measuring perplexity (PPL) to compare with vanilla decoding. For knowledge and reasoning assessment, we report accuracy on MMLU and directly compute evaluation scores on MT-Bench. Additionally, to specifically test long-context and reasoning-intensive capabilities, we conduct experiments on three challenging benchmarks: MATH500, AIME24, and LiveCodeBench v5. All experiments are conducted on NVIDIA 80GB A100 GPUs, ensuring consistent hardware conditions for fair comparison.

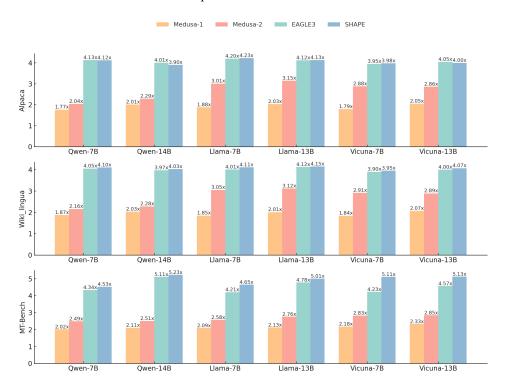


Figure 4: Speedup ratio of Qwen, Vicuna and LLaMA2-Chat inference latency on datasets Alpaca and Wiki_lingua and MT-Bench across different methods (EAGLE3, Medusa-1, Medusa-2, and SHAPE)

4.1 EFFICIENCY

We evaluate inference efficiency across diverse models and tasks, comparing SHAPE against state-of-the-art acceleration methods including EAGLE3 and Medusa 1/2. As shown in Figure 4, SHAPE achieves superior speedups ranging from 1.77× to 5.23× using original configurations from respective papers for fair comparison.

Further validation on contemporary architectures (Qwen3, DeepSeek-R1) in Table 1 shows consistent 4x-5x speedups across MT-Bench, WikiText, and Alpaca, confirming SHAPE's architectural agnosticism and practical utility without quality degradation.

We evaluate SHAPE's inference efficiency under varying batch sizes to assess its practicality in real-world deployment scenarios. Using the MT-Bench dataset on the Qwen2-7B model, we compare SHAPE against EAGLE-3, with vLLM without speculative sampling as the baseline. As shown in Table 2, SHAPE consistently outperforms EAGLE-3 across all batch sizes, demonstrating superior scalability and efficiency in practical batch processing environments. The results confirm that

Table 1: Acceleration ratios of SHAPE on recent LLMs

Model	MT-Bench	WikiText	Alpaca
Qwen3-32B	5.12×	4.35×	4.23×
Qwen3-8B	5.23×	4.67×	4.12×
DeepSeek-R1-Distill-Qwen-32B	5.05×	4.91×	$4.85 \times$
DeepSeek-R1-Distill-Qwen-8B	5.15×	4.72×	4.26×

while both methods exhibit reduced relative gains at larger batch sizes due to increased baseline parallelism, SHAPE maintains a consistent performance advantage.

Table 2: Speedup ratios at different batch sizes

Method	BS = 2	BS = 4	BS = 8	BS = 16	BS = 24
EAGLE-3	1.73×	1.65×	1.52×	1.43×	1.39×
SHAPE	1.92×	1.75×	1.61×	1.52×	1.41×

4.2 QUALITY EVALUATION

We evaluate generation quality from three perspectives: (1) general performance on standard benchmarks, (2) token-level prediction accuracy and semantic consistency, and (3) performance on reasoning-intensive and long-context tasks.

General Performance Evaluation. We assess generation quality on Alpaca and WikiText using perplexity (PPL) to evaluate language modeling capability, along with accuracy on MMLU for knowledge reasoning and MT-Bench scores for conversational ability. As shown in Table 3, SHAPE preserves output quality while delivering substantial speedups. Larger models generally maintain PPL comparable to baseline decoding, while smaller models exhibit minor variations. MMLU accuracy shows fluctuations within 1-2%, and MT-Bench evaluations confirm conversational performance remains stable. These results demonstrate SHAPE's ability to maintain quality across diverse tasks and model architectures.

Table 3: Performance comparison between original and SHAPE-accelerated models

Model	Type	Alpaca(PPL)	WikiText(PPL)	MMLU-5 shot(Acc)	MT-Bench(Score)
Qwen-7B	Vanilla	11.49	11.89	70.5	8.41
	SHAPE	11.9	12.1	68.79	8.56
Qwen-14B	Vanilla	12.30	11.92	66.3	9.08
	SHAPE	11.8	11.7	64.78	8.85
Llama-7B	Vanilla	11.76	12.77	46.2	6.27
	SHAPE	12.2	12.4	44.32	6.43
Llama-13B	Vanilla	12.67	11.94	55.0	7.05
	SHAPE	12.3	12.5	56.38	6.89
Vicuna-7B	Vanilla	11.58	12.83	48.2	6.69
	SHAPE	12.1	12.4	48.55	6.88
Vicuna-13B	Vanilla	12.06	11.63	55.28	6.81
	SHAPE	11.7	12.0	58.42	6.97

Reasoning-Intensive Task Evaluation. To evaluate SHAPE's effectiveness on challenging reasoning tasks, we conducted experiments on Qwen3 and DeepSeek-R1 using MATH500, AIME24, and LiveCodeBench v5. These benchmarks involve multi-step reasoning, long dependency chains, and domain-specific logic, providing rigorous tests for generation fidelity under complex conditions.

As shown in Table 4, SHAPE maintains near-identical accuracy compared to vanilla decoding across all models and tasks. On MATH500, accuracy differences are within 0.3%, while AIME24 and

LiveCodeBench v5 show maximum deviations of 1.3% and 0.8% respectively. These results confirm SHAPE's robustness on reasoning-heavy, long-chain tasks while delivering 4-5× speedups.

Table 4: Accuracy comparison on reasoning-intensive tasks (Vanilla / SHAPE)

Model	MATH500	AIME24	LiveCodeBench v5
Qwen3-32B	97.2 / 97.16	81.4 / 80.8	65.7 / 65.3
Qwen3-8B	97.4 / 97.1	76.0 / 74.7	57.5 / 56.9
DeepSeek-R1-Distill-Qwen-32B	94.3 / 93.2	72.6 / 72.2	54.5 / 54.1
DeepSeek-R1-Distill-Qwen-14B	93.9 / 92.1	69.7 / 68.9	45.5 / 44.7

Token-Level Analysis. Supplementary evaluations (Appendix Tables 12 and 10) show SHAPE achieves token prediction accuracy of 0.85-0.92 for 1-3 token lookahead, with lower perplexity compared to alternative acceleration methods. Semantic similarity metrics (BERTScore and embedding cosine distance) confirm strong alignment with standard decoding outputs.

4.3 ABLATION STUDY

4.3.1 EFFECTIVENESS OF OPTIMAL TRANSPORT

SHAPE employs optimal transport (OT) to model hidden state transitions, motivated by our observation that transformer hidden states maintain a minimum level of similarity between tokens at t and t+n. This "baseline similarity" indicates a theoretically valid pathway for transferring hidden states through optimal transport. Unlike autoregressive models that predict step-by-step, our OT approach captures global transition patterns by finding the optimal mapping to future states (t+n). To validate the effectiveness of OT over simpler alternatives, we conducted comparative experiments replacing the OT mapping with an affine transformation of the same dimensionality (d=128). Table 5 presents the results on Llama-7B using the Alpaca dataset with TRS configuration (N=3,K=3).

Table 5: Comparison of OT with affine transformation and analysis of different dimensionalities on Llama-7B (Alpaca). Baseline AR decoding achieves PPL=11.9.

		~ .	d Value	PPL	Speedup
Method	PPL	Speedup	32	17.3	3.87×
Affine (<i>d</i> =128)	18.4	3.87×	64	16.1	3.66×
OT (<i>d</i> =128)	12.2	3.21×	128	12.2	3.21×
			4096 (full)	12.1	2.67×

The results demonstrate that OT significantly outperforms simple affine transformations, reducing perplexity from 18.4 to 12.2 - approaching the baseline AR performance of 11.9. This validates our hypothesis that OT's ability to find optimal global mappings is crucial for accurate multi-step prediction. Furthermore, we analyzed the impact of dimensionality d on OT performance. As shown in Table 5 (right), increasing d from 32 to 128 consistently improves perplexity, with the most significant gains occurring at d=128. Interestingly, using the full dimensionality (d=4096) provides minimal perplexity improvement (12.1 vs 12.2) while reducing speedup by 17%, confirming that our low-dimensional OT approach effectively captures essential transition patterns.

4.3.2 EFFECTIVENESS OF TREE REJECTION SAMPLING

We evaluate the proposed Tree Rejection Sampling (TRS) mechanism by analyzing the speed-accuracy trade-off across various configurations with $K, N \in [1,5]$. As shown in Table 6, increasing N improves speedup but raises perplexity due to multi-step prediction errors, while increasing K reduces perplexity by providing more candidate choices at the cost of verification overhead. The configuration (K=3, N=3) achieves the optimal balance with $3.21\times$ speedup and 12.2 PPL, representing the best efficiency-quality trade-off for practical deployment. Comparative analysis with conventional decoding methods (Table 7) confirms TRS's superiority. While beam search marginally improves perplexity (11.7 vs. 11.9), it incurs a 10% slowdown. In contrast, TRS delivers substantial acceleration while maintaining competitive generation quality, successfully addressing efficiency challenges in decoder-only LLM inference.

Table 6: TRS performance under different tree configurations (Speedup/PPL)

K			Depth (N)		
	1	2	3	4	5
1	1.91/17.3	2.71/16.7	3.40/16.1	3.80/16.5	4.10/17.0
2	1.90/16.9	2.65/16.5	3.30/14.0	3.50/14.5	3.70/15.0
3	1.88/16.5	2.59/16.3	3.21/12.2	3.40/13.0	3.55/14.0
4	1.83/16.3	2.46/15.8	3.15/12.15	3.28/12.8	3.40/13.2
5	1.75/16.1	2.33/15.2	3.09/12.1	3.25/12.5	3.35/12.9

Table 7: Comparison of TRS with standard decoding methods

Method	Speedup	PPL
AR (Baseline)	1.0×	11.9
Greedy	1.0×	11.9
Beam Search	0.9×	11.7
TRS (3, 3)	3.21×	12.2

5 RELATED WORK

Recent studies have highlighted the significant inference latency of Large Language Models (LLMs), prompting various acceleration strategies that can be categorized by their underlying methodologies. on-autoregressive approaches represent initial attempts at acceleration. Non-autoregressive translation (NAT) techniques have been investigated in translation tasks (6; 18), it performs suboptimally in general LLM scenarios. To address this, Huang et al. (7) proposed a layer-wise iterative methodology that each layer leverages decoding results from preceding layers. Similarly, Santilli et al. (16) formalized autoregressive decoding through parallel Jacobi and Gauss-Seidel fixed-point iteration. However, such methods often degrade accuracy due to their deviation from standard autoregressive architectures. Accuracy-preserving approaches based on model modifications have since emerged. Block-wise parallel decoding (18) leverages an auxiliary transformer with multi-output capabilities for parallel token prediction but suffers from frequent verification failures. Medusa (2) improves robustness with multiple prediction heads, while FREE (1) uses shallow layers for draft generation. However, these techniques require substantial training of additional components, peculative decoding offers an alternative by using smaller models as draft predictors. For example, Bloom 7.1B has served as a draft model for a 176B model (25). Yet, this method faces challenges: suitable smaller models are not always available across model series, and helper models require parallel tuning, increasing deployment complexity. o address these issues, model-free strategies aim to accelerate decoding without auxiliary models. Ge et al. (4) proposed an input-guided method based on prefix matching, extended by LLMA (26) through content retrieval from inputs and external documents. Recently, LookaheadDecoding (7) fused Jacobi iteration with speculative decoding in a multi-branch

6 Conclusion

In this paper, we introduced a novel perspective that reframes autoregressive decoding as a probability distribution transition process governed by optimal transport principles. We validate this theoretical framework through SHAPE, which demonstrates predictable hidden state evolution via transport maps. Experiments across diverse LLMs show speedups of 1.77×-5.23× with maintained quality, confirming token generation can be understood as structured transport in probability space. This work establishes a new paradigm for efficient LLM inference beyond draft-based approaches.

framework, though its draft generation incurs non-negligible overhead.

REFERENCES

- [1] BAE, S., KO, J., SONG, H., AND YUN, S.-Y. Fast and robust early-exiting framework for autoregressive language models with synchronized parallel decoding, 2023.
- [2] CAI, T., LI, Y., GENG, Z., PENG, H., LEE, J. D., CHEN, D., AND DAO, T. Medusa: Simple llm inference acceleration framework with multiple decoding heads, 2024.
 - [3] FU, E. A. Lookahead: An inference acceleration framework for large language model with lossless generation accuracy. *arXiv preprint arXiv:2312.12728* (2023).
 - [4] GE, Y., ET AL. Input-guided non-autoregressive machine translation. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics* (2022), pp. 7300–7312.
 - [5] GRATTAFIORI, A., DUBEY, A., AND ABHINAV JAUHRI, E. A. The llama 3 herd of models, 2024.
 - [6] GU, J., AND KONG, X. Fully non-autoregressive neural machine translation: Tricks of the trade. *arXiv preprint arXiv:2012.15833* (2020).
 - [7] HUANG, C., ZHOU, H., ZAIANE, O. R., MOU, L., AND LI, L. Lookahead: An inference acceleration framework for large language models. *arXiv preprint arXiv:2312.12728* (2023).
 - [8] JIANG, H., WU, Q., LUO, X., LI, D., LIN, C.-Y., YANG, Y., AND QIU, L. Longlimlingua: Accelerating and enhancing llms in long context scenarios via prompt compression. *arXiv* preprint arXiv:2310.06839 (2023).
 - [9] KOU, S., HU, L., HE, Z., DENG, Z., AND ZHANG, H. Cllms: Consistency large language models, 2024.
 - [10] LI, J., TANG, T., ZHAO, W. X., NIE, J.-Y., AND WEN, J.-R. Pre-trained language models for text generation: A survey. *ACM Computing Surveys* 56, 9 (2024), 1–39.
 - [11] LI, Y., WEI, F., ZHANG, C., AND ZHANG, H. Eagle: Speculative sampling requires rethinking feature uncertainty. *arXiv preprint arXiv:2401.15077* (2024).
- [12] Lu, Z., Li, X., Cai, D., Yi, R., Liu, F., Zhang, X., Lane, N. D., and Xu, M. Small language models: Survey, measurements, and insights. *arXiv preprint arXiv:2409.15790* (2024).
- [13] Mo, Y., QIN, H., DONG, Y., ZHU, Z., AND LI, Z. Large language model (llm) ai text generation detection based on transformer deep learning algorithm, 2024.
- [14] OPENAI, ACHIAM, J., ADLER, S., AGARWAL, S., AHMAD, L., AKKAYA, I., ALEMAN, F. L., ALMEIDA, D., ALTENSCHMIDT, J., ALTMAN, S., AND SHYAMAL ANADKAT, E. A. Gpt-4 technical report, 2024.
- [15] OUYANG, L., WU, J., JIANG, X., ALMEIDA, D., AND AL, C. L. W. Training language models to follow instructions with human feedback, 2022.
- [16] SANTILLI, A., ET AL. Parallel fixed-point methods for autoregressive decoding. *arXiv preprint arXiv:2301.13379* (2023).
- [17] SHU, L., LUO, L., HOSKERE, J., ZHU, Y., LIU, Y., TONG, S., CHEN, J., AND MENG, L. Rewritelm: An instruction-tuned large language model for text rewriting. In *Proceedings of the AAAI Conference on Artificial Intelligence* (2024), vol. 38, pp. 18970–18980.
- 532 [18] STERN, M., CHAN, W., KIROS, J., AND USZKOREIT, J. Blockwise parallel decoding for deep autoregressive models. In *Advances in Neural Information Processing Systems* (2018), pp. 10107–10116.
- [19] Sun, Z., Yu, H., Song, X., Liu, R., Yang, Y., and Zhou, D. Mobilebert: a compact task-agnostic bert for resource-limited devices, 2020.
 - [20] THAKUR, S., AHMAD, B., PEARCE, H., TAN, B., DOLAN-GAVITT, B., KARRI, R., AND GARG, S. Verigen: A large language model for verilog code generation. *ACM Transactions on Design Automation of Electronic Systems (TODAES)* 29, 3 (2024), 1–31.

- [21] TOUVRON, H., LAVRIL, T., IZACARD, G., AND XAVIER MARTINET, E. A. Llama: Open and efficient foundation language models, 2023.
- 543 [22] WANG, L., YANG, N., HUANG, X., YANG, L., MAJUMDER, R., AND WEI, F. Large search model: Redefining search stack in the era of llms, 2024.
 - [23] WU, Y. *Large Language Model and Text Generation*. Springer International Publishing, Cham, 2024, pp. 265–297.
 - [24] XAI. Grok os: The future of ai assistants, 2024. Accessed: 2025-01-31.
 - [25] XIA, Q., ET AL. Speculative decoding for large language models. *arXiv preprint* arXiv:2302.01318 (2023).
 - [26] YANG, F., ET AL. Llma: Language model acceleration via content retrieval. *arXiv preprint* arXiv:2303.16827 (2023).

A TOKEN-LEVEL AUTOREGRESSIVE GENERATION

A.1 SINGLE-STEP GENERATION PROCESS

In autoregressive language models, token generation follows a step-by-step process. At each time step t, given the sequence of previous tokens (x_1, x_2, \ldots, x_t) , the probability of generating the next token x_{t+1} is:

$$P(x_{t+1}|x_1,\ldots,x_t) \tag{1}$$

A.2 OUTPUT HIDDEN STATE BASED GENERATION

The generation process involves the final layer's hidden states:

$$\mathbf{h}_t = \operatorname{Transformer}(x_1, \dots, x_t) \tag{2}$$

$$P(x_{t+1}|x_1,\dots,x_t) = \text{LLM_head}(\mathbf{h}_t)$$
(3)

where $\mathbf{h}_t \in \mathbb{R}^d$ represents the final layer's hidden state at time step t, and LLM_head is a linear transformation that maps the hidden state to token probabilities over the vocabulary.

B MODEL AR DECODING PERFORMANCE METRICS

Table 8 presents the average token generation time across different model sizes and context lengths. The results clearly demonstrate that larger models and longer contexts significantly increase pertoken latency, which accumulates due to the sequential nature of autoregressive decoding. These findings highlight the importance of optimizing the decoding process to ensure practical deployment efficiency.

Table 8: Autoregressive Decoding Latency across Different Input Lengths and Model Scales

Model (B)	Input Length	ITL (ms)	TTFT (ms)	Duration (s/req)
	256	3.83	24.58	0.56
1.5	512	3.85	33.91	0.80
1.5	1024	3.83	55.05	1.52
	2048	3.98	118.47	1.38
	256	7.16	42.93	1.06
7	512	7.12	73.25	1.67
1	1024	7.15	118.62	2.88
	2048	7.17	274.23	2.65
	256	11.90	76.58	1.78
14	512	11.92	134.24	2.57
14	1024	11.98	253.11	4.92
	2048	12.12	603.10	4.64
	256	22.26	116.42	3.31
32	512	22.28	211.08	4.76
32	1024	22.44	392.29	9.13
	2048	22.55	924.34	8.64

C THEORETICAL ANALYSIS OF OPTIMAL TRANSPORT

Lemma A (Lipschitz map from hidden state to distribution). Let $\ell = W\mathbf{h}$ and $\mathbf{p} = \operatorname{softmax}(\ell/\tau_s)$. If $\|W\|_2 \leq L_W$ and \mathbf{h} is confined to a bounded set, then there exists $L_S > 0$ such that

$$\|\mathbf{p}(\mathbf{h}_1) - \mathbf{p}(\mathbf{h}_2)\|_1 \le \frac{L_S L_W}{\tau_s} \|\mathbf{h}_1 - \mathbf{h}_2\|_2.$$

Sketch. Softmax on bounded domains is Lipschitz in ℓ_2 (or ℓ_∞); composing with the linear map W yields the claim.

Lemma B (Similarity \Rightarrow small **OT move).** Let $\mu_t = \sum_i \mathbf{p}_t(i)\delta_{E_i}$ and μ_{t+1} be defined analogously. Under Lemma A, there exists L' > 0 (depending on W, τ_s, E) such that

$$W_c(\mu_t, \mu_{t+1}) \leq L' \|\mathbf{h}_{t+1} - \mathbf{h}_t\|_2.$$

If we normalize $\bar{\mathbf{h}}_t = \mathbf{h}_t/\|\mathbf{h}_t\|_2$, then $\|\bar{\mathbf{h}}_{t+1} - \bar{\mathbf{h}}_t\|_2 \leq \sqrt{2(1-\cos(\mathbf{h}_t,\mathbf{h}_{t+1}))}$, hence

$$W_c(\mu_t, \mu_{t+1}) \leq \tilde{L} \sqrt{1 - \mathrm{SC}(\mathbf{h}_t, \mathbf{h}_{t+1})}.$$

Sketch. Use the Kantorovich–Rubinstein dual bound with ℓ_1 variation of **p** and the diameter of the embedding support, plus the cosine– ℓ_2 relation.

Proposition C (Existence, stability, and uniqueness of Π_t^{\star}).

$$\Pi_{t}^{\star} = \arg \min_{\Pi \mathbf{1} = \mathbf{p}_{t}, \ \Pi^{\top} \mathbf{1} = \mathbf{p}_{t+1}} \left\langle \Pi, C \right\rangle + \varepsilon \operatorname{KL} \left(\Pi \parallel \mathbf{p}_{t} \mathbf{p}_{t+1}^{\top} \right), \quad \varepsilon > 0,$$
(4)

For any $\varepsilon>0$, the entropic OT problem in equation 4 admits a unique solution Π_t^\star ; moreover, when $W_c(\mu_t,\mu_{t+1})$ is small, Π_t^\star depends smoothly on $(\mathbf{p}_t,\mathbf{p}_{t+1})$ and can be well-approximated by a few Sinkhorn iterations. *Sketch*. Entropic regularization makes the objective strictly convex over the transport polytope; standard Sinkhorn–Knopp scaling solves the KKT system, and continuity follows from the implicit function theorem on the strictly convex objective.

Corollary D (OT-optimal path between successive distributions). By Proposition C, the coupling Π_t^* induces a row-stochastic operator K_t such that $\mathbf{p}_{t+1} = K_t^\top \mathbf{p}_t$ holds exactly at optimality and approximately under finite Sinkhorn iterations, thereby defining the OT-optimal path for the one-step distributional transition. *Sketch*. Row normalization rewrites the marginal constraints; the equality follows from $\Pi_t^{*\top} \mathbf{1} = \mathbf{p}_{t+1}$.

D THEORETICAL ANALYSIS OF HIDDEN STATE PREDICTION VIA OPTIMAL TRANSPORT

We establish a theoretical framework for predicting future hidden states in transformer models through optimal transport theory. Let (Ω, \mathcal{F}, P) be a probability space and $\mathcal{H} \subseteq \mathbb{R}^d$ be the hidden state space. For any time step t, we define $H_t: \Omega \to \mathcal{H}$ as the random variable representing the hidden state at time t, with μ_t as its probability measure. Let $\mathcal{P}(\mathcal{H})$ denote the space of probability measures on \mathcal{H} .

Given the temporal nature of hidden states in transformer models, we first establish their similarity properties. The similarity between hidden states is measured by cosine similarity:

$$sim(x,y) = \frac{\langle x,y \rangle}{\|x\| \|y\|} \tag{5}$$

Based on empirical observations in transformer models, as shown in Table 9, we make the following assumption:

For any adjacent time steps t and t+1, the hidden states maintain a minimum similarity threshold:

$$\forall x, y \in \mathcal{H} : \sin(x, y) > T \tag{6}$$

where x and y are hidden states with positive probability under μ_t and μ_{t+1} respectively.

This assumption leads to our first key result regarding the bounded evolution of hidden states:

Under Assumption 1, there exists a constant M>0 such that the 2-Wasserstein distance between consecutive hidden state distributions is bounded:

$$W_2(\mu_t, \mu_{t+1}) \le M \tag{7}$$

Consider any hidden states $x, y \in \mathcal{H}$ with positive probability under μ_t and μ_{t+1} respectively. From Assumption 1:

$$1 - \frac{\langle x, y \rangle}{\|x\| \|y\|} \le T \tag{8}$$

This implies:

$$\langle x, y \rangle \ge T \|x\| \|y\| \tag{9}$$

Define the Euclidean metric $d(x,y) = ||x-y||_2$. We can expand:

$$d^{2}(x,y) = ||x||^{2} + ||y||^{2} - 2\langle x, y \rangle$$
(10)

$$\leq ||x||^2 + ||y||^2 - ||x|| ||y|| \tag{11}$$

$$= (\|x\| - \|y\|)^2 \tag{12}$$

Since \mathcal{H} is bounded in \mathbb{R}^d , there exists R > 0 such that $||x|| \leq R$ for all $x \in \mathcal{H}$. Therefore:

$$d^2(x,y) \le 4R^2 \tag{13}$$

Taking M = 2R completes the proof.

This lemma establishes that the evolution of hidden states is well-behaved, allowing us to formulate our main theorem:

There exists a cost function $c: \mathcal{H} \times \mathcal{H} \to \mathbb{R}_+$ such that the hidden state evolution from time t to t+k can be represented as an optimal transport problem:

$$\min_{\pi \in \Pi(\mu_t, \mu_{t+k})} \int_{\mathcal{H} \times \mathcal{H}} c(x, y) d\pi(x, y)$$
 (14)

where $\Pi(\mu_t, \mu_{t+k})$ denotes the set of joint distributions with marginals μ_t and μ_{t+k} . Moreover, this problem admits an optimal solution π^* .

We construct the proof in three steps. First, we define the cost function $c(x,y) = d^2(x,y)$, where d is the Euclidean metric. This choice is natural as it preserves the geometric structure of the hidden state space.

Second, from Lemma 1, we know that for adjacent time steps:

$$W_2^2(\mu_t, \mu_{t+1}) = \inf_{\pi \in \Pi(\mu_t, \mu_{t+1})} \int_{\mathcal{H} \times \mathcal{H}} d^2(x, y) d\pi(x, y) \le M^2$$
(15)

For multi-step evolution (k > 1), we can apply the Chapman-Kolmogorov equation. There exist intermediate measures $\pi_1, ..., \pi_{k-1}$ such that:

$$W_2^2(\mu_t, \mu_{t+k}) \le \left(\sum_{i=0}^{k-1} W_2(\mu_{t+i}, \mu_{t+i+1})\right)^2 \le k^2 M^2 \tag{16}$$

Finally, the existence of an optimal solution follows from three key properties: 1) $\mathcal{P}(\mathcal{H})$ is compact in the weak topology 2) The cost function c(x,y) is lower semi-continuous 3) The objective functional is bounded below

By the Kantorovich duality theorem, there exists an optimal solution $\pi^* \in \Pi(\mu_t, \mu_{t+k})$ achieving:

$$\int_{\mathcal{H}\times\mathcal{H}} c(x,y)d\pi^*(x,y) = \inf_{\pi\in\Pi(\mu_t,\mu_{t+k})} \int_{\mathcal{H}\times\mathcal{H}} c(x,y)d\pi(x,y)$$
(17)

This theoretical framework provides a rigorous foundation for predicting hidden states through optimal transport. Given a hidden state h_t at time t, we can predict h_{t+k} by:

$$h_{t+k} = \int_{\mathcal{H}} y d\pi^*(y|h_t) \tag{18}$$

Moreover, we can establish an error bound for this prediction:

$$||h_{t+k} - h_{t+k}^*||_2 \le kM \tag{19}$$

where h_{t+k}^* denotes the true hidden state at time t+k.

Table 9: Token similarity between the token at time t and t+n across various contexts.

	qwen-zh qwen-en				vicuna			llama				
t+1	t+2	t+3	t+1	t+2	t+3	t+1	t+2	t+3	t+1	t+2	t+3	
0.8744	0.8447	0.8383	0.6304	0.5273	0.4976	0.5392	0.4647	0.4358	0.645	0.5729	0.5443	
0.8677	0.8235	0.814	0.6272	0.5703	0.5503	0.5282	0.4287	0.4062	0.6411	0.5829	0.552	
0.8666	0.8213	0.8027	0.6217	0.5334	0.5027	0.526	0.436	0.4174	0.6351	0.5647	0.5422	
0.8607	0.8171	0.7849	0.6189	0.5478	0.5316	0.5248	0.4285	0.4121	0.6329	0.5598	0.5322	
0.8599	0.8314	0.8195	0.6157	0.5459	0.5435	0.5165	0.4208	0.4023	0.6261	0.5456	0.523	
0.8587	0.821	0.8079	0.6122	0.5373	0.5199	0.5165	0.4185	0.3951	0.6236	0.5503	0.5309	
0.8574	0.829	0.809	0.6118	0.5426	0.5239	0.516	0.422	0.4008	0.622	0.5407	0.5141	
0.8544	0.8209	0.8052	0.6107	0.5384	0.5291	0.5149	0.416	0.3971	0.6214	0.5391	0.5151	
0.8536	0.823	0.8071	0.6098	0.5303	0.5144	0.5121	0.4031	0.3832	0.6197	0.531	0.5091	
0.8535	0.8141	0.7921	0.6093	0.5295	0.5046	0.5113	0.4153	0.3913	0.6197	0.531	0.5091	

E QUALITY EVALUATION

Table 10: Comparison of perplexity (ppl) across different decoding methods (EAGLE, M-1: Medusa-1, M-2: Medusa-2, and SHAPE) on various DS: datasets (A: Alpaca, T: THUC news, W: wiki lingua).

Model	Dataset	EAGLE	M-1	M-2	SHAPE
	A	13.2	15.1	14.6	11.9
Qwen-7B	T	_	_	_	12.3
	W	13.5	15.3	14.8	12.1
	Α	12.9	14.1	13.7	11.8
Qwen-14B	T	_	_	_	11.5
	W	13.2	14.2	13.6	11.7
Llama-7B	A	13.1	15.0	14.4	12.2
Liailia-/D	W	13.4	15.2	14.6	12.4
Llama-13B	Α	12.8	14.0	13.5	12.3
Liailia-13B	W	13.1	14.1	13.7	12.5
Vicuna-7B	A	13.0	15.2	14.3	12.1
vicuna-/B	W	13.3	15.3	14.4	12.4
Viouno 12D	Α	12.9	14.2	14.0	11.7
Vicuna-13B	W	13.0	14.3	14.1	12.0

Table 11: SHAPE's Average speed-up ratio compared with vanilla generation on different datasets

Model			S	Speed Up				
(Size)	Datasets	Task	N ste	N step forward				
			1	2	3			
	Alpaca	Intruction Following	1.90	2.47	4.12			
Qwen-7B	THUC_News wiki_lingua	Text Continuation Text Generation	1.96 1.93	2.51 2.55	4.09 4.07			
Alpaca Intruction Following Qwen-14B THUC_News Text Continuation wiki_lingua Text Generation			1.87 1.99 1.89	2.55 2.65 2.77	3.89 4.02 4.05			
Llama-7B	Alpaca wiki_lingua	Intruction Following Text Generation	1.88 1.91	2.59 2.69	3.21 3.34			
Llama-13B	Alpaca wiki_lingua	Intruction Following Text Generation	1.89 1.90	2.51 2.53	3.22 3.20			
Vicina-/B		Intruction Following Text Generation	1.83 1.87	2.57 2.56	3.24 3.29			
Vicuna-13B	Alpaca wiki_lingua	Intruction Following Text Generation	1.89 1.91	2.68 2.77	3.38 3.21			

Table 12: Comparison of performance metrics between vanilla decoding (baseline) and SHAPE across different N-step forward prediction configurations.

					Average en Accu				rage Score		Sei	Aver	age Similari	ty
Model	Size	Datasets	Task					N st	ep forw	ard				
				1	2	3	V	1 S	2 S	3 S	V	1 S	2 S	3 S
Qwen	7B	Alpaca THUC_News wiki_lingua	IF TC TG	0.91 0.92 0.91	0.89 0.87 0.88	0.86 0.85 0.86	0.74 0.75 0.72	0.71 0.70 0.69	0.68 0.67 0.65	0.65 0.66 0.63	0.94 0.93 0.92	0.92 0.90 0.87	0.88 0.89 0.85	0.84 0.82 0.81
	14B	Alpaca THUC_News wiki_lingua	IF TC TG	0.91 0.90 0.91	0.89 0.88 0.87	0.86 0.85 0.85	0.76 0.74 0.75	0.73 0.71 0.70	0.70 0.68 0.66	0.68 0.67 0.64	0.95 0.94 0.93	0.92 0.89 0.88	0.91 0.87 0.85	0.85 0.84 0.80
Llama	7B	Alpaca wiki_lingua	IF TG	0.90 0.89	0.89 0.87	0.86 0.85	0.60 0.63	0.59 0.60	0.55 0.58	0.52 0.54	0.82 0.80	0.78 0.76	0.75 0.73	0.70 0.69
2141114	13B	Alpaca wiki_lingua	IF TG	0.90 0.89	0.88 0.87	0.87 0.86	0.54 0.56	0.53 0.54	0.51 0.50	0.49 0.47	0.76 0.78	0.73 0.74	0.70 0.71	0.67 0.68
Vicuna	7B	Alpaca wiki_lingua	IF TG	0.91 0.89	0.88 0.87	0.86 0.85	0.56 0.58	0.55 0.56	0.52 0.54	0.50 0.52	0.76 0.78	0.74 0.75	0.71 0.72	0.68 0.69
	13B	Alpaca wiki_lingua	IF TG	0.90 0.90	0.88 0.87	0.85 0.87	0.58 0.59	0.57 0.57	0.54 0.53	0.52 0.51	0.79 0.80	0.76 0.77	0.73 0.74	0.70 0.71

F TRAIN DETAILS AND TRAINING COST

F.1 Data Acquisition and Generation

For the English dataset used to train models such as Qwen, Vicuna, and Llama on 7B, 13B, and 14B, we use ShareGPT as the dataset, which contains 96,000 dialogue data samples. For the Chinese dataset, we use the THUCNews training set, consisting of 50,000 news samples, split into two parts: "prompt" and "completion". The target model processes pre-processed data to generate outputs from the transformer layers for each token. The training data includes fields such as input token IDs, hidden states, hidden states from the previous three tokens, target values, attention masks, and loss masks. These components are combined to construct the final training dataset.

F.2 TRAIN CONFIGURATION

The training configuration includes the following settings: Learning rate (1r) and batch size (bs) are dynamically adjusted, with gradient accumulation steps set to ensure stable training. The number of epochs is set to 40, and a warm-up phase of 2,000 steps is applied, targeting a total of 800,000 steps. The configuration employs a maximum sequence length of 2,048 tokens, balancing performance and memory efficiency. To improve robustness, data noise is introduced using a uniform distribution with a mean of 0 and a standard deviation of 0.2. Additional settings include weight decay, gradient clipping, and periodic model saving every epoch. Finally, the optimizer uses momentum parameters (b1 = 0.9, b2 = 0.95) to facilitate effective training convergence.

F.3 TRAINING COST

We compare SHAPE with Medusa (+2 heads) in terms of parameter count and computational efficiency. SHAPE contains approximately 450.98M parameters (6.4% of LLaMA-7B), achieved through parameter sharing and lightweight modules such as OT projection and a gating network. In contrast, Medusa adds 300M parameters per head, reaching 610M (9% of LLaMA-7B) with two heads. Thus, SHAPE uses only 71.5% of the parameters of Medusa+2 heads and requires no per-head adaptation. In training, SHAPE is 3–5× faster than full model fine-tuning and requires only 5 hours for 40 epochs on a single A100 GPU, compared to 8 hours for Medusa. Peak training memory usage is also lower: 41.89GB for SHAPE versus 51.47GB for Medusa. During inference, SHAPE achieves maxium 5.23x speedup through one hidden state correction step (OT module) and tree rejection sampling. In contrast, Medusa incurs higher overhead due to additional multi-head attention and memory usage, limiting its speedup to 1.8–3.1×.

G TREE-BASED REJECT SAMPLING ALGORITHM IMPLEMENTATION

Algorithm 1 Tree Rejection Sampling

```
Require: model: target language model
```

- 1: context: current context or hidden state at time step t
- 2: N: number of future steps (depth)
- 3: k: number of candidates per step (branch factor)

Ensure: selected_prefix: the longest valid prefix among accepted paths

- 4: $candidate_paths \leftarrow generate_candidates(model, context, N, k)$ // Generate k^N candidate sequences
- 5: $path_probs \leftarrow model.get_path_probabilities(candidate_paths)$ // Compute joint probabilities in parallel
- 6: $max_prob \leftarrow max(path_probs)$
- 7: $threshold \leftarrow 0.8 \times max_prob$ // Define acceptance threshold (e.g. 80% of the maximum probability)
- 8: $accepted_paths \leftarrow []$
- 9: for each (path, prob) in (candidate_paths, path_probs) do
- 10: **if** $prob \ge threshold$ **then**
 - 11: accepted_paths.append(path) // Retain paths with sufficiently high probability
 - 12: **end if**
 - 13: **end for**
 - 14: $selected_prefix \leftarrow select_longest_valid_prefix(accepted_paths)$ // Extract the longest prefix common to accepted paths
 - return $selected_prefix$