

Mitigating Sequential Dependencies: A Survey of Algorithms and Systems for Generation-Refinement Frameworks in Autoregressive Models

Anonymous ACL submission

Abstract

Sequential dependencies present a fundamental bottleneck in deploying large-scale autoregressive models, particularly for real-time applications. While traditional optimization approaches like pruning and quantization often compromise model quality, recent advances in generation-refinement frameworks demonstrate that this trade-off can be significantly mitigated.

This survey presents a comprehensive taxonomy of generation-refinement frameworks, analyzing methods across autoregressive sequence tasks. We categorize methods based on their generation strategies (from simple n-gram prediction to sophisticated draft models) and refinement mechanisms (including single-pass verification and iterative approaches). Through systematic analysis of both algorithmic innovations and system-level implementations, we examine deployment strategies across computing environments and explore applications spanning text, images, and speech generation. This systematic examination of both theoretical frameworks and practical implementations provides a foundation for future research in efficient autoregressive decoding.

1 Introduction

Large Models (LMs) have demonstrated remarkable capabilities across diverse domains, from text generation (Brown et al., 2020; Zhuang et al., 2023; Touvron et al., 2023) and translation (Zhu et al., 2023; Hadi et al., 2023; Huang et al., 2023) to image synthesis (Ho et al., 2020; Yang et al., 2023a; Tian et al., 2024) and video generation (Ding et al., 2023; Wu et al., 2023; ope, 2024). However, these models face a critical challenge: their inherently sequential nature creates significant latency bottlenecks, particularly for real-time applications. While traditional optimization approaches like quantization and pruning often compromise

model quality for speed, recent research has focused on maintaining output quality while breaking sequential dependencies through novel algorithmic and system-level innovations.

Generation-refinement frameworks have emerged as a promising family of solutions that directly address these sequential bottlenecks. These approaches encompass a range of methods, from speculative decoding with draft models to iterative refinement techniques inspired by numerical optimization. The common thread among these approaches is their division of the generation process into two phases: an initial generation step that produces draft tokens in parallel, followed by a refinement step that ensures output quality.

The implementation of these frameworks presents unique system-level challenges across different deployment scenarios. Edge devices require careful optimization of memory usage and computation patterns (Svirschevski et al., 2024; Xu et al., 2024a), while distributed systems must manage complex communication patterns and load balancing. These system-level considerations have driven innovations in areas like kernel design, hardware acceleration, and batch processing optimization, significantly influencing both algorithmic choices and practical performance.

This survey synthesizes research across these approaches, examining both algorithmic innovations and their system implementations. We present a systematic taxonomy of generation-refinement methods, analyze deployment strategies across computing environments, and explore applications spanning text, images (Wang et al., 2024d; Jang et al., 2024), and speech (Li et al., 2024a; Raj et al., 2024). Our primary contributions include comprehensive analysis of system-level implementations and optimizations, detailed examination of applications across modalities, and identification of key research challenges in efficient neural sequence

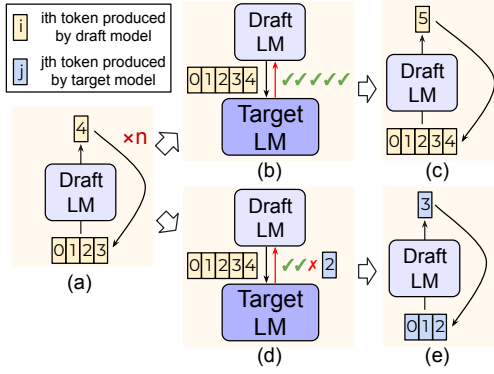


Figure 1: Illustration of speculative decoding workflow.

generation.

2 The Sequential Bottleneck in Large Model Inference

Traditional approaches to accelerating LM inference have focused on reducing computational costs through model compression, knowledge distillation, and architectural optimizations. However, these methods primarily address individual computation costs rather than the fundamental sequential dependency that requires each token to wait for all previous tokens.

Speculative decoding (SD) (Stern et al., 2018) has emerged as a promising solution that directly targets this sequential bottleneck. As illustrated in Figure 1, this approach introduces a two-phase process where a smaller, faster *draft model* first predicts multiple tokens in parallel, followed by verification using the target model. The draft model enables parallel token generation, breaking away from traditional token-by-token generation, while the target model’s verification step maintains output quality through accept/reject decisions.

This strategy has proven particularly valuable for real-time applications like interactive dialogue systems, where response latency directly impacts user experience. The verification mechanism provides a crucial balance between generation speed and output quality, accepting correct predictions to maintain throughput while falling back to sequential generation when necessary to preserve accuracy.

While SD represents one successful approach to breaking sequential dependencies in autoregressive (AR) models, it belongs to a broader family of *generation-refinement* methods. The following sections present a systematic taxonomy of these approaches, examining how different techniques balance the trade-offs between generation paral-

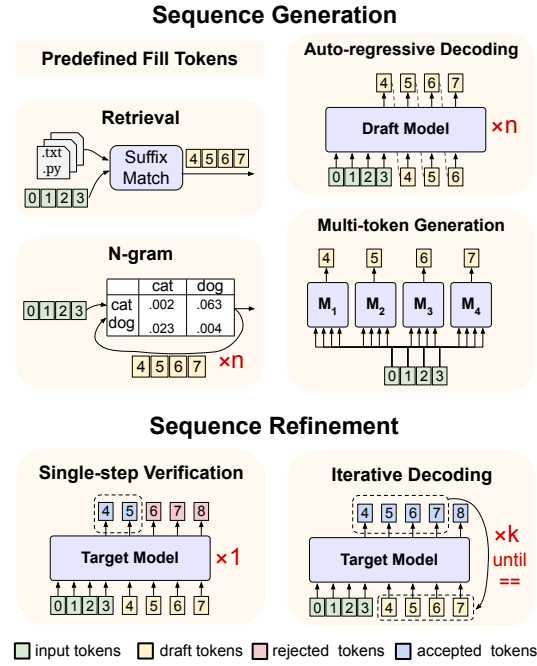


Figure 2: A taxonomy of generation-refinement frameworks, showing two phases: (1) Generation of draft tokens through various methods and (2) Refinement through verification strategies.

lelism and output quality.

3 A Taxonomy for Generation and Refinement Frameworks

To systematically analyze approaches for breaking sequential dependencies in large models, we propose a unified taxonomy that categorizes methods based on their generation and refinement strategies. As shown in Figure 2, our taxonomy decomposes these frameworks into two fundamental phases: *Sequence Generation* and *Sequence Refinement*. This decomposition not only encompasses traditional SD approaches but also captures a broader range of emerging methods that trade off between generation parallelism and output quality.

The sequence generation phase focuses on different strategies for producing draft tokens more efficiently than conventional auto-regressive decoding using a single larger model. These strategies range from simple approaches like random token sampling (used in conjunction with iterative decoding) to more sophisticated methods like retrieval-based generation and draft model prediction. Each generation method offers trade-offs in terms of computational cost and prediction quality. The sequence refinement phase then determines how these candidates are processed - either accepting them directly

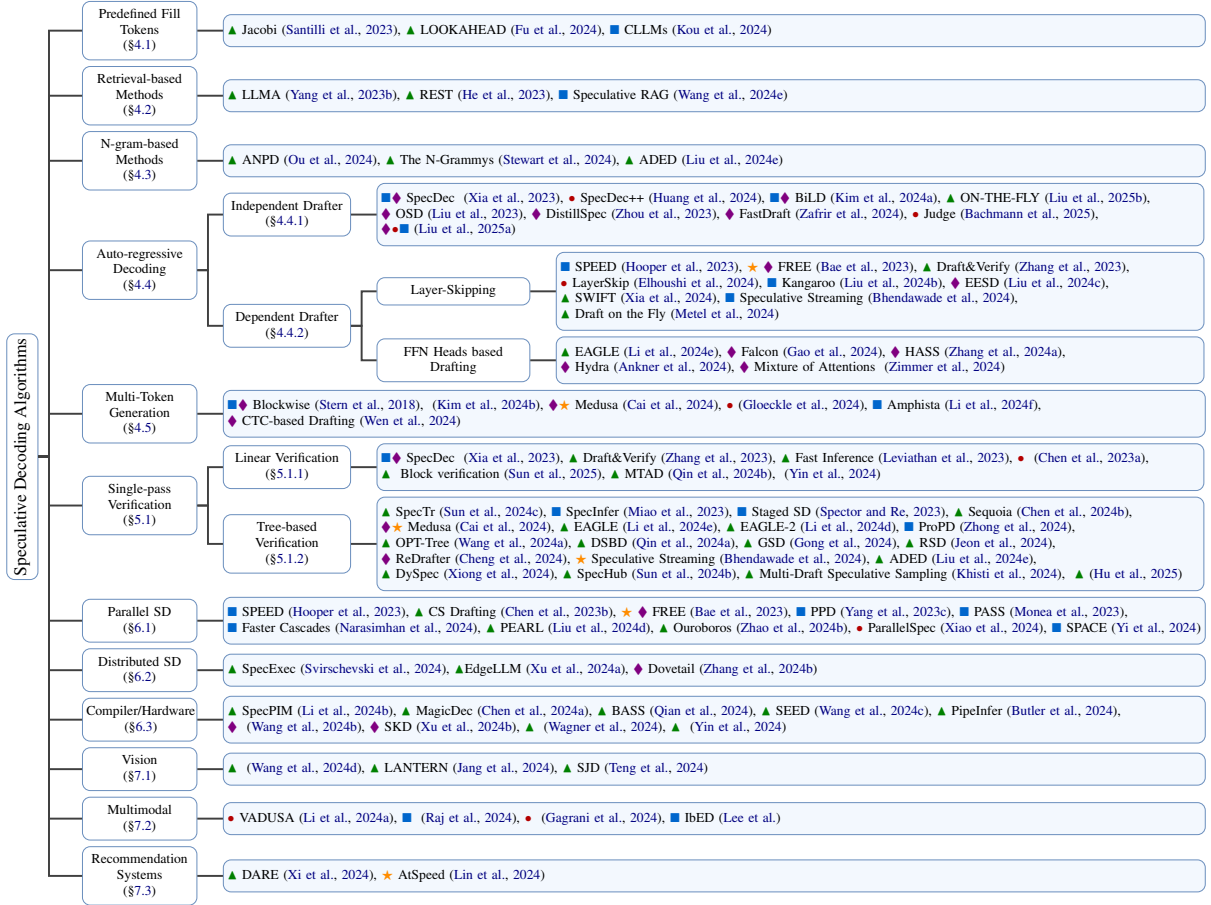


Figure 3: Taxonomy of Speculative Decoding Algorithms. Symbols indicate implementation approach: ▲ Direct application (no training required), ● Full model training from scratch, ■ Model fine-tuning, ★ Parameter-efficient fine-tuning (PEFT), ◆ Knowledge distillation from target model.

(with possible poorer quality), verifying a subset of tokens in a single pass, or refining the draft tokens through multiple iterations until convergence.

4 Sequence Generation Methods

4.1 Predefined Fill Tokens

The simplest approach uses random initialization or predefined tokens (e.g., PAD). While computationally free, these methods provide poor initialization points, requiring multiple refinement iterations as discussed in Section 5.2.

4.2 Retrieval-based Methods

LLMA (Yang et al., 2023b) first proposed exploiting overlaps between LLM outputs and reference documents to accelerate inference through parallel token verification while maintaining identical generation results. In retrieval-based approaches, REST (He et al., 2023) replaces smaller language models with exact suffix matching from a datastore to generate draft tokens. It builds a Trie (pre-

fix tree) from retrieved continuations, where node weights reflect token sequence frequencies. Speculative RAG (Wang et al., 2024e) use a fine-tuned specialist LM to generate complete answer drafts with supporting rationales. It clusters retrieved documents by similarity, generates diverse drafts from different document subsets, and employs self-consistency and self-reflection scores for draft evaluation instead of token-level verification.

4.3 N-gram-based Methods

Several approaches leverage n-gram patterns for efficient token generation. ANPD (Ou et al., 2024) replaces traditional draft models with an adaptive N-gram system that updates predictions based on context. LOOKAHEAD (Fu et al., 2024) uses n-gram verification by collecting and utilizing n-grams from previous iterations as draft tokens. The N-Grammys (Stewart et al., 2024) further develops this idea by creating a dedicated n-gram based prediction system that can operate without requiring a separate draft model.

4.4 Auto-regressive Generation

Most sequence generation methods employ auto-regressive drafting, where a smaller model generates draft tokens that are verified by a larger target model. This drafting paradigm has spawned numerous techniques that vary in how the draft model interacts with the target model.

4.4.1 Independent Drafters

Independent drafters use smaller models to generate tokens sequentially while a larger target model verifies them in parallel. SpecDec (Xia et al., 2023) pioneered this approach with an independent draft model using distinct attention queries for masked positions. SpecDec++ (Huang et al., 2024) improves SpecDec (Xia et al., 2023) by training a prediction head on top of the draft model that estimates the probability of token acceptance by the target model. Based on these predictions, it dynamically determines when to stop generating tokens and trigger verification.

Recent works focus on dynamic adaptation and confidence monitoring. BiLD (Kim et al., 2024a) triggers target model verification when draft confidence falls below a threshold, while ON-THE-FLY (Liu et al., 2025b) dynamically adjusts window sizes based on prediction accuracy. OSD (Liu et al., 2023) enables online adaptation through knowledge distillation during inference, and Distill-Spec (Zhou et al., 2023) extends this by accessing target model logits for improved alignment. (Liu et al., 2025a) introduces special tokens for draft models to autonomously determine target model consultation, eliminating separate verification at some performance cost. For mathematical applications, Judge (Bachmann et al., 2025) adds a learned verification layer atop the target model’s embeddings, using contextual correctness assessment to reduce strict output alignment requirements.

4.4.2 Dependent Drafters

The main drawbacks of independent drafting approaches are that (1) the computation required to generate the draft tokens is fixed per tokens, meaning that computation is over-provisioned for many “easy” tokens and (2) the target model cannot reuse the features of the drafting process, increasing the amount of compute required. Self-speculative decoding approaches generate draft tokens by relying directly on a subset (**Layer Skipping**) or extension (**Dependent Heads**) of the target model.

Layer Skipping Draft&Verify (Zhang et al., 2023), SWIFT (Xia et al., 2024), and Draft on the Fly (Metel et al., 2024) achieves fast draft token generation by selectively skipping some intermediate layers in the Draft process, and then verifies these drafts using the full LLM. In order to achieve good draft accuracy, they also designed an intermediate layer selection algorithm based on Bayesian optimization. LayerSkip (Elhoushi et al., 2024) uses an early exiting (Teerapittayanon et al., 2016) approach to dynamically output tokens at different depths of the target model. Kangaroo (Liu et al., 2024b) also applied early exit by adopting a shallow sub-network to generate drafts and using a lightweight adapter module to bridge the performance gap with the full model, achieving efficient and accurate decoding. EESD (Liu et al., 2024c) use Thompson Sampling Control (Slivkins et al., 2019) Mechanism to adaptively determines how many draft token will be generated. SPEED (Hooper et al., 2023) combines speculative execution with parameter sharing, using early predictions to process multiple tokens in parallel through shared decoder layers, rather than waiting for each token to complete sequentially.

Dependent Heads Dependent head-based drafting eliminates the need for a separate draft model by adding lightweight feed-forward prediction heads using the hidden states of the target model. The main idea is that the first token in sequence generation block uses the target model as usual but the features at the end of the model are fed into additional heads to predict subsequent tokens without passing back through the entire target model.

EAGLE (Li et al., 2024e) uses a trained head that takes in hidden states from the target model and generates subsequent draft tokens in an AR manner. Hydra (Ankner et al., 2024) use multiple decoding, one for each draft token position.

EAGLE extensions have focused on improving parallel token generation and attention mechanisms. Falcon (Gao et al., 2024) introduces a semi-autoregressive framework combining LSTM layers and relaxed causal-masked self-attention to generate k tokens per forward pass, while HASS (Zhang et al., 2024a) enhances knowledge distillation by prioritizing high-probability tokens during training. Mixture of Attention (Zimmer et al., 2024) incorporates multiple attention types (LSA, SA, and CA) for improved token prediction, and DeepSeek-V3 (Liu et al., 2024a) adapts (Gloeckle et al.,

2024)’s multi-token approach (discussed in Section 4.5) while maintaining complete causal attention during inference.

4.5 Multi-token Prediction

Stern et al. (2018) proposes adding multiple decoding heads on top of a model to predict k future tokens in parallel, requiring training the entire model from scratch. Medusa (Cai et al., 2024) introduces a parameter-efficient approach, where lightweight decoding heads are fine-tuned on top of pre-trained language models. Each head is trained to predict a specific future position in the sequence without modifying the target model.

Recent approaches improve Medusa’s independent draft heads by modeling inter-token relationships. Amphista (Li et al., 2024f) uses bi-directional self-attention to consider both past and future predictions, while CTC Drafting (Wen et al., 2024) employs Connectionist Temporal Classification (CTC) with blank tokens and repetition, followed by duplicate removal to generate draft sequences.

5 Sequence Refinement Methods

5.1 Single-pass Verification

Single-pass verification represents the most common refinement strategy in draft-and-verify approaches, where drafted tokens are verified exactly once by the target model.

5.1.1 Linear Verification

Linear verification sequentially validates draft tokens against the target model’s logit distributions, with early works like SpecDec (Xia et al., 2023) and Draft&Verify (Zhang et al., 2023) comparing drafted tokens against the target model’s predictions. When a token fails verification (i.e., when the draft output doesn’t match the target model’s distribution), the system falls back to standard AR generation from that point.

Fast Inference (Leviathan et al., 2023) and (Chen et al., 2023a) introduced speculative sampling to improve acceptance rates while approximately maintaining the target distribution. Their method accepts a token if the target model assigns equal or higher probability; otherwise, it accepts with probability $p(x)/q(x)$ or resamples from an adjusted distribution.

Block Verification (Sun et al., 2025) and MTAD (Qin et al., 2024b) improve upon linear ver-

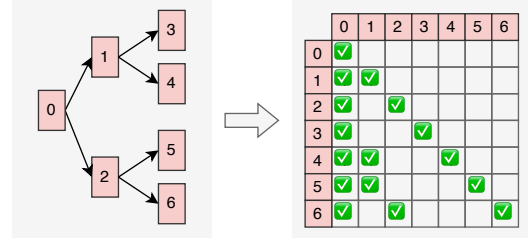


Figure 4: Illustration of tree-based speculative decoding, with token tree construction on the left and tree attention mask on the right.

ification by examining the joint probability distribution of draft tokens as a chain of conditional probabilities. This block-based evaluation approach typically results in higher acceptance rates compared to token-by-token verification for similar quality.

5.1.2 Tree-based Verification

Tree-based verification extends the single-pass paradigm by enabling parallel exploration of multiple completion paths. Unlike linear verification that processes a single sequence, tree-based methods construct and verify a tree of possible completions simultaneously, making more efficient use of parallel compute resources.

SpecInfer (Miao et al., 2023) pioneered this approach by developing an efficient tree-based attention masking scheme that enables parallel verification while maintaining proper token dependencies. This innovation maintains generation quality while significantly increasing the number of tokens that can be verified in parallel.

Recent works have focused on optimizing tree structure and size to maximize computational efficiency. Sequoia (Chen et al., 2024b) introduces a hardware-aware tree optimizer that can maximize inference performance by selecting appropriate tree dimensions based on available computing resources. OPT-Tree (Wang et al., 2024a) searches for optimal tree structures to maximize expected acceptance length per decoding step. DSBD (Qin et al., 2024a) uses a small model to generate multiple candidate sequences via beam search, then the large model verifies these sequences layer by layer while dynamically adjusting the beam width based on acceptance probabilities to balance efficiency and quality. DySpec (Xiong et al., 2024) enables dynamic tree expansion during runtime based on prediction confidence, while EAGLE2 (Li et al., 2024d) incorporates context-aware tree construction to improve acceptance rates. DDD (Brown

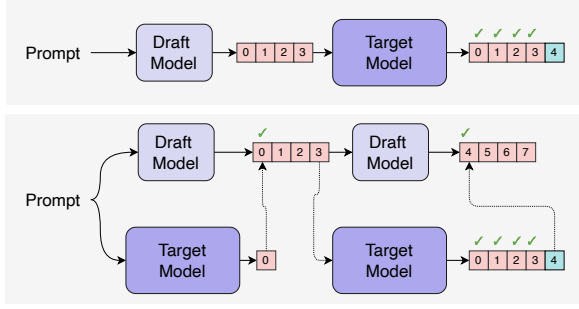


Figure 5: Comparison of speculative decoding approaches: (a) Sequential processing where draft generates tokens (0-3) before target verification. (b) Parallel processing where draft generates new tokens while target simultaneously verifies previous ones.

et al., 2024) optimizes EAGLE2 (Li et al., 2024d)’s tree drafting method by making the depth dynamic based on draft model confidence.

Several works have explored hybrid approaches that combine tree-based verification with other techniques. ProPD (Zhong et al., 2024) integrates progressive refinement into the tree structure, while RSD (Jeon et al., 2024) employs recursive verification strategies. GSD (Gong et al., 2024) and ADED (Liu et al., 2024e) extend tree-based methods to handle more complex dependency structures through graph-based representations and adaptive depth adjustment.

In terms of verifying multiple candidate draft tokens in parallel (also known as Multi-Draft Speculative Decoding, MDSD), (Hu et al., 2025) propose a hybrid sampling strategy that combines deterministic selection of high-probability tokens with random sampling of the final token, improving acceptance rates in certain scenarios. (Khisti et al., 2024) introduce a two-phase verification method that uses importance sampling to select a candidate token before applying single-draft verification, optimizing the process for parallel draft generation.

5.2 Iterative Decoding

Iterative decoding methods extend the single-pass verification paradigm by allowing multiple refinement iterations on draft tokens until convergence. These approaches draw inspiration from classical numerical methods for solving systems of nonlinear equations, particularly the Jacobi and Gauss-Seidel iteration methods.

In Santilli et al. (2023), the authors reframe AR text generation as an iterative optimization problem. Their approach expresses token generation

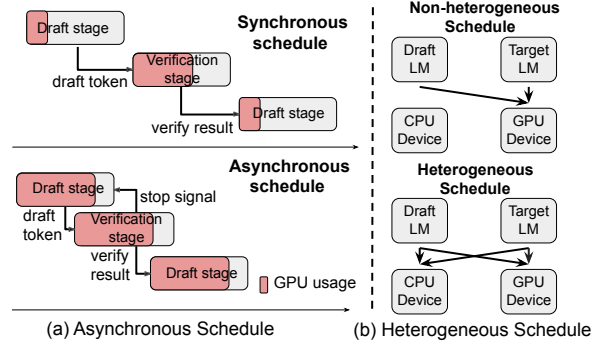


Figure 6: Asynchronous and heterogeneous schedules.

as a system where each position must output the most likely token given the current state of all other positions. Starting with a randomly initialized sequence, they adapt the Jacobi method to update all positions in parallel during each iteration until convergence. The authors prove that this process produces identical output to traditional AR decoding under greedy sampling. Fu et al. (2024) builds upon this framework with LOOKAHEAD decoding, which combines Jacobi iterations with n-gram verification to accelerate convergence by leveraging predictions from earlier steps.

CLLMs (Kou et al., 2024) leverages consistency training to accelerate convergence by enabling better multi-token prediction in early iterations.

6 System-Level Optimizations and Implementation Strategies

6.1 Parallel Speculative Decoding

Traditional SD processes tokens sequentially, with the draft model generating tokens followed by target model verification, creating inherent bottlenecks. As shown in Figure 5, parallel approaches overcome this limitation by enabling simultaneous operation - while the target model verifies earlier tokens, the draft model generates subsequent ones, enabling continuous overlapped execution. Recent methods build upon this paradigm: CS Drafting (Chen et al., 2023b) employs vertical and horizontal cascade structures for 81% speedup, PaSS (Monea et al., 2023) uses look-ahead embeddings for 30% speedup, and Faster Cascades (Narasimhan et al., 2024) incorporates deferral rules for improved cost-quality trade-offs. PEARL (Liu et al., 2024d) further advances this through pre-verify and post-verify strategies with adaptive draft lengths, achieving $4.43\times$ speedup over AR decoding and $1.50\times$ over

standard SD AMUSD (McDaniel, 2024) presents an asynchronous multi-device approach to SD, decoupling the draft and verify phases into continuous, asynchronous operations.

6.2 Distributed Speculative Decoding

Edge computing environments impose stringent constraints on memory, compute power, and latency, necessitating specialized SD approaches to deploy LLMs effectively in resource-constrained settings. SpecExec (Svirschevski et al., 2024) is designed to harness the parallel processing power of consumer GPUs to accelerate LLM inference. By generating multiple tokens per target model iteration and constructing a “cache” tree of probable continuations, SpecExec efficiently validates these continuations with the target model in a single pass. EdgeLLM (Xu et al., 2024a) further optimizes on-device LLM inference through novel techniques for resource allocation and error correction, achieving great token generation speeds and significantly outperforming existing engines. Dovetail (Zhang et al., 2024b) represents a significant advancement in heterogeneous computing for LLM inference. By deploying the draft model on the GPU and the target model on the CPU, Dovetail reduces the granularity of data transfer and enhances the overall inference process. The introduction of Dynamic Gating Fusion (DGF) and optimizations for low-end hardware further improve the balance between latency and performance.

6.3 Compiler and Hardware Optimization for Speculative Decoding

Efficient implementation of SD requires careful optimization of both hardware resources and compiler strategies to maximize throughput and minimize latency. SpecPIM (Li et al., 2024b) presents a novel approach to accelerate speculative inference on a Processing-in-Memory (PIM) system through co-exploration of architecture and dataflow. This method constructs a design space that comprehensively considers algorithmic and architectural heterogeneity, enabling optimal hardware resource allocation for different models and computational patterns. (Wagner et al., 2024) investigates improvements in speculative sampling on GPUs, achieving significant speed gains by parallelizing computations and using sigmoid approximations for softmax, though this comes with a minor reduction in accuracy.

Recent studies have focused on enhancing the

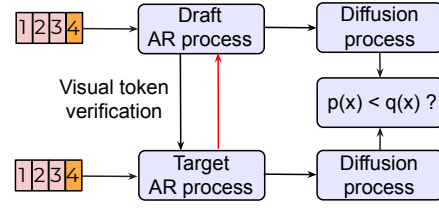


Figure 7: Flow of AR image generation with SD.

throughput of LLMs using SD by optimizing batch processing and scheduling strategies. Figure 6 illustrates two scheduling strategies for SD systems: (a) Asynchronous Schedule: The draft stage is followed by the verify stage, with optional stop signals determining further processing. This non-blocking approach enhances system efficiency. (b) Heterogeneous Schedule: Both CPU and GPU devices are utilized for different stages of the decoding process, enabling parallel processing and optimizing performance through resource allocation. Using Markov chain theory, (Yin et al., 2024) establishes SD’s optimality among unbiased algorithms while highlighting the tradeoff between inference speed and output quality. Their analysis reveals that batch processing benefits are limited by the distribution gap between small and large models. MagicDec (Chen et al., 2024a) identifies the shift from compute-bound to memory-bound bottlenecks as batch size and sequence length increase, using sparse KV caches in draft models to optimize throughput. BASS (Qian et al., 2024) extends SD to a batched setting with customized CUDA kernels for ragged tensors in attention calculations and dynamically adjusts draft lengths for better GPU utilization. SEED (Wang et al., 2024c) accelerates reasoning tree construction through scheduled speculative execution, using a rounds-scheduled strategy for conflict-free parallel processing. PipeInfer (Butler et al., 2024) addresses single-request latency through pipelined speculative acceleration, reducing inter-token latency via asynchronous speculation and early cancellation. TRIFORCE (Sun et al., 2024a) introduces a hierarchical SD mechanism with a dynamic sparse KV cache to achieve lossless acceleration of long sequence generation, significantly improving generation speed and efficiency while maintaining quality. (Zhao et al., 2024a) proposes QSPEC, a novel framework that combines weight-shared quantization schemes with SD, achieving up to 1.55× acceleration without quality loss, paving the way for efficient and high-fidelity quantization deployment in diverse and

memory-constrained settings. (Wang et al., 2024b) introduces a hardware-aware SD algorithm that accelerates the inference speed of Mamba and hybrid models. Inspired by SD, SKD (Xu et al., 2024b) represents a novel, adaptive approach to knowledge distillation. By dynamically generating tokens and using the teacher model to filter or replace low-quality samples, it bridges the gap between supervised KD’s reliance on static data and on-policy KD’s susceptibility to low-quality outputs. This ensures a better alignment between training and inference distributions, and improved performance.

7 Multimodal Models and Applications

7.1 Speculative Decoding for Visual Output Generation

Researchers are now using SD to improve the efficiency of AR image generation (Ding et al., 2021; Yu et al., 2022; Li et al., 2024c). As shown in Figure 7, this method greatly speeds up the process by reducing the inference steps needed for generating visual tokens. For instance, (Wang et al., 2024d) proposes a novel continuous SD method that designs a novel acceptance criterion for the diffusion distributions, significantly improving the efficiency of AR image generation. Similarly, LANTERN (Jang et al., 2024) presents a relaxed acceptance condition for the SD strategy to substantially speed up the inference process in visual AR models. Additionally, Speculative Jacobi Decoding (SJD) (Teng et al., 2024) offers a training-free speculative Jacobi decoding technique that effectively accelerates text-to-image generation tasks.

7.2 Speculative Decoding for Multimodal Output Generation

Recent advancements in SD have substantially improve the efficiency and quality of AR generation across various modalities. In the domain of speech synthesis, VADUSA (Li et al., 2024a) leverages SD to accelerate the inference process in AR text-to-speech (TTS) systems, which enhances the quality speech synthesis as well. Inspired by the flavor of SD, (Raj et al., 2024) introduces a multi-token prediction mechanism, offering substantial improvements in inference efficiency for speech generation.

In the context of multimodal large language models, (Gagrani et al., 2024) investigates the integration of SD into the LLaVA 7B model to optimize inference efficiency. Their findings indicate that employing a lightweight, language-only draft model

facilitates a memory-constrained acceleration of up to 2.37 \times . Besides, IbED (Lee et al.) proposes the "In-batch Ensemble Drafting" method to further enhance the robustness and efficiency of SD. It adopts the ensemble techniques during batch-level inference, requires no additional model parameters and significantly increases the validation probability of draft tokens, thereby improving performance and robustness across diverse input scenarios.

7.3 Recommendation Systems

LLM-based recommendation systems have shown great potential in enhancing personalized recommendations, but their high inference latency poses a significant challenge for real-world deployment. To address this, recent research has focused on optimizing decoding efficiency to accelerate recommendation generation. (Xi et al., 2024) propose DARE that integrates retrieval-based SD to accelerate recommendation knowledge generation, thereby improving the deployment efficiency of LLM-based recommender systems in industrial settings. AtSpeed (Lin et al., 2024) combines strict top-K alignment (AtSpeed-S) and relaxed sampling verification (AtSpeed-R), to significantly accelerate LLM-based generative recommendation with speedup from 2 \times to 2.5 \times , addressing inference latency challenges in top-K sequence generation.

8 Conclusion

This survey has presented a comprehensive analysis of generation-refinement frameworks for mitigating sequential dependencies in autoregressive models, highlighting how these approaches are fundamentally changing efficient neural sequence generation across text, speech, and visual domains. Through examining both algorithmic innovations and system-level implementations, we have demonstrated their broad applicability while providing crucial deployment insights for practitioners. Moving forward, significant challenges persist in constructing solid theoretical foundations to grasp the balance between parallelism and quality, as well as in developing comprehensive approaches that span different modalities—efforts that could narrow the divide between the capabilities of large models and their actual implementation. Additionally, it remains crucial to examine the scalability of the speculative decoding system as the quantity of draft and target models increases, particularly within large-scale LLM systems.

Limitations

While this survey provides a comprehensive overview of generation-refinement frameworks, some limitations should be acknowledged. Detailed performance comparisons across different approaches are challenging due to varying experimental settings, model architectures, and hardware configurations used in the original papers. The lack of standardized benchmarks for speculative decoding makes it difficult to make definitive claims about the relative efficiency of different methods. Additionally, while we examine applications across different modalities, our analysis may not fully capture all domain-specific challenges and optimizations, particularly for emerging areas like video generation and multimodal reasoning.

References

2024. Open-sora report v1.1. https://github.com/hpcaitech/Open-Sora/blob/main/docs/report_02.md.
- Zachary Ankner, Rishab Parthasarathy, Aniruddha Nrusimha, Christopher Rinard, Jonathan Ragan-Kelley, and William Brandon. 2024. Hydra: Sequentially-dependent draft heads for medusa decoding. *arXiv preprint arXiv:2402.05109*.
- Gregor Bachmann, Sotiris Anagnostidis, Albert Pumarola, Markos Georgopoulos, Artsiom Sanakoyeu, Yuming Du, Edgar Schönfeld, Ali Thabet, and Jonas K Kohler. 2025. *Judge decoding: Faster speculative sampling requires going beyond model alignment*. In *The Thirteenth International Conference on Learning Representations*.
- Sangmin Bae, Jongwoo Ko, Hwanjun Song, and Se-Young Yun. 2023. Fast and robust early-exiting framework for autoregressive language models with synchronized parallel decoding. *arXiv preprint arXiv:2310.05424*.
- Nikhil Bhendawade, Irina Belousova, Qichen Fu, Henry Mason, Mohammad Rastegari, and Mahyar Najibi. 2024. Speculative streaming: Fast llm inference without auxiliary models. *arXiv preprint arXiv:2402.11131*.
- Oscar Brown, Zhengjie Wang, Andrea Do, Nikhil Mathew, and Cheng Yu. 2024. Dynamic depth decoding: Faster speculative decoding for llms. *arXiv preprint arXiv:2409.00142*.
- 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. *Advances in neural information processing systems*, 33:1877–1901.
- Branden Butler, Sixing Yu, Arya Mazaheri, and Ali Janesari. 2024. Pipeinfer: Accelerating llm inference using asynchronous pipelined speculation. In *SC24: International Conference for High Performance Computing, Networking, Storage and Analysis*, pages 1–19. IEEE.
- 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. 2023a. 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*.
- Zhuoming Chen, Avner May, Ruslan Svirschevski, Yuhsun Huang, Max Ryabinin, Zhihao Jia, and Beidi Chen. 2024b. Sequoia: Scalable, robust, and hardware-aware speculative decoding. *arXiv preprint arXiv:2402.12374*.
- Ziyi Chen, Xiaocong Yang, Jiacheng Lin, Chenkai Sun, Kevin Chen-Chuan Chang, and Jie Huang. 2023b. Cascade speculative drafting for even faster llm inference. *arXiv preprint arXiv:2312.11462*.
- Yunfei Cheng, Aonan Zhang, Xuanyu Zhang, Chong Wang, and Yi Wang. 2024. Recurrent drafter for fast speculative decoding in large language models. *arXiv preprint arXiv:2403.09919*.
- Ming Ding, Zhuoyi Yang, Wenyi Hong, Wendi Zheng, Chang Zhou, Da Yin, Junyang Lin, Xu Zou, Zhou Shao, Hongxia Yang, et al. 2021. Cogview: Mastering text-to-image generation via transformers. *Advances in neural information processing systems*, 34:19822–19835.
- Ning Ding, Xingtai Lv, Qiaosen Wang, Yulin Chen, Bowen Zhou, Zhiyuan Liu, and Maosong Sun. 2023. Sparse low-rank adaptation of pre-trained language models. *arXiv preprint arXiv:2311.11696*.
- Mostafa Elhoushi, Akshat Shrivastava, Diana Liskovich, Basil Hosmer, Bram Wasti, Liangzhen Lai, Anas Mahmoud, Bilge Acun, Saurabh Agarwal, Ahmed Roman, et al. 2024. Layer skip: Enabling early exit inference and self-speculative decoding. *arXiv preprint arXiv:2404.16710*.
- Yichao Fu, Peter Bailis, Ion Stoica, and Hao Zhang. 2024. Break the sequential dependency of llm inference using lookahead decoding. *arXiv preprint arXiv:2402.02057*.

Mukul Gagrani, Raghavv Goel, Wonseok Jeon, Junyoung Park, Mingu Lee, and Christopher Lott. 2024. On speculative decoding for multimodal large language models. <i>arXiv preprint arXiv:2404.08856</i> .	797
Wonseok Jeon, Mukul Gagrani, Raghavv Goel, Junyoung Park, Mingu Lee, and Christopher Lott. 2024. Recursive speculative decoding: Accelerating llm inference via sampling without replacement. <i>arXiv preprint arXiv:2402.14160</i> .	798
Xiangxiang Gao, Weisheng Xie, Yiwei Xiang, and Feng Ji. 2024. Falcon: Faster and parallel inference of large language models through enhanced semi-autoregressive drafting and custom-designed decoding tree. <i>arXiv preprint arXiv:2412.12639</i> .	799
Ashish Khisti, M Reza Ebrahimi, Hassan Dbouk, Arash Behboodi, Roland Memisevic, and Christos Louizos. 2024. Multi-draft speculative sampling: Canonical architectures and theoretical limits. <i>arXiv preprint arXiv:2410.18234</i> .	800
Fabian Gloeckle, Badr Youbi Idrissi, Baptiste Rozière, David Lopez-Paz, and Gabriel Synnaeve. 2024. Better & faster large language models via multi-token prediction. <i>arXiv preprint arXiv:2404.19737</i> .	801
Sehoon Kim, Karttikeya Mangalam, Suhong Moon, Jitendra Malik, Michael W Mahoney, Amir Gholami, and Kurt Keutzer. 2024a. Speculative decoding with big little decoder. <i>Advances in Neural Information Processing Systems</i> , 36.	802
Zhuocheng Gong, Jiahao Liu, Ziyue Wang, Pengfei Wu, Jingang Wang, Xunliang Cai, Dongyan Zhao, and Rui Yan. 2024. Graph-structured speculative decoding. <i>arXiv preprint arXiv:2407.16207</i> .	803
Taehyeon Kim, Ananda Theertha Suresh, Kishore A Papineni, Michael Riley, Sanjiv Kumar, and Adrian Benton. 2024b. Accelerating blockwise parallel language models with draft refinement. In <i>The Thirtieth Annual Conference on Neural Information Processing Systems</i> .	804
Muhammad Usman Hadi, R Qureshi, A Shah, M Irfan, A Zafar, MB Shaikh, N Akhtar, J Wu, and S Mirjalili. 2023. A survey on large language models: Applications, challenges, limitations, and practical usage. <i>TechRxiv</i> .	805
Siqi Kou, Lanxiang Hu, Zhezhi He, Zhijie Deng, and Hao Zhang. 2024. Cllms: Consistency large language models. <i>arXiv preprint arXiv:2403.00835</i> .	806
Zhenyu He, Zexuan Zhong, Tianle Cai, Jason D Lee, and Di He. 2023. Rest: Retrieval-based speculative decoding. <i>arXiv preprint arXiv:2311.08252</i> .	807
Minjae Lee, Wonjun Kang, Minghao Yan, Christian Classen, Hyung Il Koo, and Kangwook Lee. In-batch ensemble drafting: Toward fast and robust speculative decoding for multimodal language models.	808
Jonathan Ho, Ajay Jain, and Pieter Abbeel. 2020. Denoising diffusion probabilistic models. <i>Advances in neural information processing systems</i> , 33:6840–6851.	809
Yaniv Leviathan, Matan Kalman, and Yossi Matias. 2023. Fast inference from transformers via speculative decoding. In <i>International Conference on Machine Learning</i> , pages 19274–19286. PMLR.	810
Coleman Hooper, Sehoon Kim, Hiva Mohammadzadeh, Hasan Genc, Kurt Keutzer, Amir Gholami, and Sophia Shao. 2023. Speed: Speculative pipelined execution for efficient decoding. <i>arXiv preprint arXiv:2310.12072</i> .	811
Bohan Li, Hankun Wang, Situo Zhang, Yiwei Guo, and Kai Yu. 2024a. Fast and high-quality auto-regressive speech synthesis via speculative decoding. <i>arXiv preprint arXiv:2410.21951</i> .	812
Zhengmian Hu, Tong Zheng, Vignesh Viswanathan, Ziyi Chen, Ryan A. Rossi, Yihan Wu, Dinesh Manocha, and Heng Huang. 2025. Towards optimal multi-draft speculative decoding. In <i>The Thirteenth International Conference on Learning Representations</i> .	813
Cong Li, Zhe Zhou, Size Zheng, Jiaxi Zhang, Yun Liang, and Guangyu Sun. 2024b. Specpin: Accelerating speculative inference on pim-enabled system via architecture-dataflow co-exploration. In <i>Proceedings of the 29th ACM International Conference on Architectural Support for Programming Languages and Operating Systems, Volume 3</i> , pages 950–965.	814
Hui Huang, Shuangzhi Wu, Xinnian Liang, Bing Wang, Yanrui Shi, Peihao Wu, Muyun Yang, and Tiejun Zhao. 2023. Towards making the most of llm for translation quality estimation. In <i>CCF International Conference on Natural Language Processing and Chinese Computing</i> , pages 375–386. Springer.	815
Tianhong Li, Yonglong Tian, He Li, Mingyang Deng, and Kaiming He. 2024c. Autoregressive image generation without vector quantization. <i>arXiv preprint arXiv:2406.11838</i> .	816
Kaixuan Huang, Xudong Guo, and Mengdi Wang. 2024. Specdec++: Boosting speculative decoding via adaptive candidate lengths. <i>arXiv preprint arXiv:2405.19715</i> .	817
Yuhui Li, Fangyun Wei, Chao Zhang, and Hongyang Zhang. 2024d. Eagle-2: Faster inference of language models with dynamic draft trees. <i>Preprint</i> , arXiv:2406.16858.	818
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. <i>arXiv preprint arXiv:2410.03355</i> .	819
Yuhui Li, Fangyun Wei, Chao Zhang, and Hongyang Zhang. 2024e. Eagle: Speculative sampling requires rethinking feature uncertainty. <i>arXiv preprint arXiv:2401.15077</i> .	820

852	Zeping Li, Xinlong Yang, Ziheng Gao, Ji Liu, Zhuang	speculative inference and verification. <i>arXiv preprint</i>	906
853	Liu, Dong Li, Jinzhang Peng, Lu Tian, and Emad	<i>arXiv:2305.09781</i> .	907
854	Barsoum. 2024f. Amphista: Accelerate llm in-		
855	ference with bi-directional multiple drafting heads	Giovanni Monea, Armand Joulin, and Edouard Grave.	908
856	in a non-autoregressive style. <i>arXiv preprint</i>	2023. Pass: Parallel speculative sampling. <i>arXiv</i>	909
857	<i>arXiv:2406.13170</i> .	<i>preprint arXiv:2311.13581</i> .	910
858	Xinyu Lin, Chaoqun Yang, Wenjie Wang, Yongqi Li,		
859	Cunxiao Du, Fuli Feng, See-Kiong Ng, and Tat-Seng	Harikrishna Narasimhan, Wittawat Jitkrittum,	911
860	Chua. 2024. Efficient inference for large language	Ankit Singh Rawat, Seungyeon Kim, Neha Gupta,	912
861	model-based generative recommendation. <i>arXiv</i>	Aditya Krishna Menon, and Sanjiv Kumar. 2024.	913
862	<i>preprint arXiv:2410.05165</i> .	Faster cascades via speculative decoding. <i>arXiv</i>	914
863	Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang,	<i>preprint arXiv:2405.19261</i> .	915
864	Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi		
865	Deng, Chenyu Zhang, Chong Ruan, et al. 2024a.	Jie Ou, Yueming Chen, and Wenhong Tian. 2024.	916
866	Deepseek-v3 technical report. <i>arXiv preprint</i>	Lossless acceleration of large language model via	917
867	<i>arXiv:2412.19437</i> .	adaptive n-gram parallel decoding. <i>arXiv preprint</i>	918
868	Fangcheng Liu, Yehui Tang, Zhenhua Liu, Yunsheng	<i>arXiv:2404.08698</i> .	919
869	Ni, Kai Han, and Yunhe Wang. 2024b. Kangaroo:		
870	Lossless self-speculative decoding via double early	Haifeng Qian, Sujun Kumar Gonugondla, Sungsoo Ha,	920
871	exiting. <i>arXiv preprint arXiv:2404.18911</i> .	Mingyue Shang, Sanjay Krishna Gouda, Ramesh Nal-	921
872	Guanlin Liu, Anand Ramachandran, Tanmay Gangwani,	lapati, Sudipta Sengupta, Xiaofei Ma, and Anoop De-	922
873	Yan Fu, and Abhinav Sethy. 2025a. Knowledge dis-	oras. 2024. Bass: Batched attention-optimized spec-	923
874	tillation with training wheels.	ulative sampling. <i>arXiv preprint arXiv:2404.15778</i> .	924
875	Jiahao Liu, Qifan Wang, Jingang Wang, and Xunliang		
876	Cai. 2024c. Speculative decoding via early-exiting	Zongyue Qin, Zifan He, Neha Prakriya, Jason Cong,	925
877	for faster llm inference with thompson sampling con-	and Yizhou Sun. 2024a. Dynamic-width speculative	926
878	trol mechanism. <i>arXiv preprint arXiv:2406.03853</i> .	beam decoding for efficient llm inference. <i>arXiv</i>	927
879	Jiesong Liu, Brian Park, and Xipeng Shen. 2025b. A	<i>preprint arXiv:2409.16560</i> .	928
880	drop-in solution for on-the-fly adaptation of specula-	Zongyue Qin, Ziniu Hu, Zifan He, Neha Prakriya, Jason	929
881	tive decoding in large language models.	Cong, and Yizhou Sun. 2024b. Optimized multi-	930
882	Tianyu Liu, Yun Li, Qitan Lv, Kai Liu, Jianchen Zhu,	token joint decoding with auxiliary model for llm	931
883	and Winston Hu. 2024d. Parallel speculative de-	inference. <i>arXiv preprint arXiv:2407.09722</i> .	932
884	coding with adaptive draft length. <i>arXiv preprint</i>		
885	<i>arXiv:2408.11850</i> .	Desh Raj, Gil Keren, Junteng Jia, Jay Mahadeokar,	933
886	Xiaoxuan Liu, Lanxiang Hu, Peter Bailis, Alvin Che-	and Ozlem Kalinli. 2024. Faster speech-llama in-	934
887	ung, Zhijie Deng, Ion Stoica, and Hao Zhang.	ference with multi-token prediction. <i>arXiv preprint</i>	935
888	2023. Online speculative decoding. <i>arXiv preprint</i>	<i>arXiv:2409.08148</i> .	936
889	<i>arXiv:2310.07177</i> .	Andrea Santilli, Silvio Severino, Emiliano Postolache,	937
890	Xukun Liu, Bowen Lei, Ruqi Zhang, and Dongkuan	Valentino Maiorca, Michele Mancusi, Riccardo	938
891	Xu. 2024e. Adaptive draft-verification for efficient	Marin, and Emanuele Rodolà. 2023. Accelerating	939
892	large language model decoding. <i>arXiv preprint</i>	transformer inference for translation via parallel de-	940
893	<i>arXiv:2407.12021</i> .	coding. <i>arXiv preprint arXiv:2305.10427</i> .	941
894	Bradley McDanel. 2024. Amusd: Asynchronous multi-	Aleksandrs Slivkins et al. 2019. Introduction to multi-	942
895	device speculative decoding for llm acceleration.	armed bandits. <i>Foundations and Trends® in Machine</i>	943
896	<i>arXiv preprint arXiv:2410.17375</i> .	<i>Learning</i> , 12(1-2):1–286.	944
897	Michael R Metel, Peng Lu, Boxing Chen, Mehdi Reza-	Benjamin Spector and Chris Re. 2023. Accelerating llm	945
898	gholizadeh, and Ivan Kobayev. 2024. Draft on the	inference with staged speculative decoding. <i>arXiv</i>	946
899	fly: Adaptive self-speculative decoding using cosine	<i>preprint arXiv:2308.04623</i> .	947
900	similarity. <i>arXiv preprint arXiv:2410.01028</i> .	Mitchell Stern, Noam Shazeer, and Jakob Uszkoreit.	948
901	Xupeng Miao, Gabriele Oliaro, Zhihao Zhang, Xin-	2018. Blockwise parallel decoding for deep autore-	949
902	hao Cheng, Zeyu Wang, Zhengxin Zhang, Rae	gressive models. <i>Advances in Neural Information</i>	950
903	Ying Yee Wong, Alan Zhu, Lijie Yang, Xiaoxiang	<i>Processing Systems</i> , 31.	951
904	Shi, et al. 2023. Specinfer: Accelerating genera-	Lawrence Stewart, Matthew Trager, Sujun Kumar	952
905	tive large language model serving with tree-based	Gonugondla, and Stefano Soatto. 2024. The n-	953
		grammys: Accelerating autoregressive inference with	954
		learning-free batched speculation. <i>arXiv preprint</i>	955
		<i>arXiv:2411.03786</i> .	956

1064	Daliang Xu, Wangsong Yin, Hao Zhang, Xin Jin, Ying	Libo Zhang, Zhaoning Zhang, Baizhou Xu, Songzhu	1119
1065	Zhang, Shiyun Wei, Mengwei Xu, and Xuanzhe Liu.	Mei, and Dongsheng Li. 2024b. Dovetail: A cpu/gpu	1120
1066	2024a. Edgellm: Fast on-device llm inference with	heterogeneous speculative decoding for llm inference.	1121
1067	speculative decoding. <i>IEEE Transactions on Mobile</i>	<i>arXiv preprint arXiv:2412.18934</i> .	1122
1068	<i>Computing</i> .		
1069	Wenda Xu, Rujun Han, Zifeng Wang, Long T Le, Dhruv	Juntao Zhao, Wenhao Lu, Sheng Wang, Lingpeng Kong,	1123
1070	Madeka, Lei Li, William Yang Wang, Rishabh Agar-	and Chuan Wu. 2024a. Qspec: Speculative decoding	1124
1071	wal, Chen-Yu Lee, and Tomas Pfister. 2024b. Spec-	with complementary quantization schemes. <i>arXiv</i>	1125
1072	ulative knowledge distillation: Bridging the teacher-	<i>preprint arXiv:2410.11305</i> .	1126
1073	student gap through interleaved sampling. <i>arXiv</i>	Weilin Zhao, Yuxiang Huang, Xu Han, Wang Xu,	1127
1074	<i>preprint arXiv:2410.11325</i> .	Chaojun Xiao, Xinrong Zhang, Yewei Fang, Kai-	1128
1075	Ling Yang, Zhilong Zhang, Yang Song, Shenda Hong,	huo Zhang, Zhiyuan Liu, and Maosong Sun. 2024b.	1129
1076	Runsheng Xu, Yue Zhao, Wentao Zhang, Bin Cui,	Ouroboros: Generating longer drafts phrase by	1130
1077	and Ming-Hsuan Yang. 2023a. Diffusion models: A	phrase for faster speculative decoding . In <i>Proceed-</i>	1131
1078	comprehensive survey of methods and applications.	<i>ings of the 2024 Conference on Empirical Methods in</i>	1132
1079	<i>ACM Computing Surveys</i> , 56(4):1–39.	<i>Natural Language Processing</i> , pages 13378–13393,	1133
1080	Nan Yang, Tao Ge, Liang Wang, Binxing Jiao, Daxin	Miami, Florida, USA. Association for Computational	1134
1081	Jiang, Linjun Yang, Rangan Majumder, and Furu	Linguistics.	1135
1082	Wei. 2023b. Inference with reference: Lossless ac-	Shuzhang Zhong, Zebin Yang, Meng Li, Ruihao Gong,	1136
1083	celeration of large language models. <i>arXiv preprint</i>	Runsheng Wang, and Ru Huang. 2024. Propd: Dy-	1137
1084	<i>arXiv:2304.04487</i> .	namic token tree pruning and generation for llm par-	1138
1085	Seongjun Yang, Gibbeum Lee, Jaewoong Cho, Dim-	allel decoding. <i>arXiv preprint arXiv:2402.13485</i> .	1139
1086	itris Papailiopoulous, and Kangwook Lee. 2023c.	Yongchao Zhou, Kaifeng Lyu, Ankit Singh Rawat,	1140
1087	Predictive pipelined decoding: A compute-latency	Aditya Krishna Menon, Afshin Rostamizadeh, San-	1141
1088	trade-off for exact llm decoding. <i>arXiv preprint</i>	jiv Kumar, Jean-François Kagy, and Rishabh Agar-	1142
1089	<i>arXiv:2307.05908</i> .	wal. 2023. Distillspec: Improving speculative de-	1143
1090	Hanling Yi, Feng Lin, Hongbin Li, Ning Peiyang, Xi-	coding via knowledge distillation. <i>arXiv preprint</i>	1144
1091	aotian Yu, and Rong Xiao. 2024. Generation meets	<i>arXiv:2310.08461</i> .	1145
1092	verification: Accelerating large language model in-	Wenhao Zhu, Hongyi Liu, Qingxiu Dong, Jingjing Xu,	1146
1093	ference with smart parallel auto-correct decoding .	Lingpeng Kong, Jiajun Chen, Lei Li, and Shujian	1147
1094	In <i>Findings of the Association for Computational</i>	Huang. 2023. Multilingual machine translation with	1148
1095	<i>Linguistics: ACL 2024</i> , pages 5285–5299, Bangkok,	large language models: Empirical results and analy-	1149
1096	Thailand. Association for Computational Linguistics.	sis. <i>arXiv preprint arXiv:2304.04675</i> .	1150
1097	Ming Yin, Minshuo Chen, Kaixuan Huang, and	Yuchen Zhuang, Yue Yu, Kuan Wang, Haotian Sun, and	1151
1098	Mengdi Wang. 2024. A theoretical perspective	Chao Zhang. 2023. Toolqa: A dataset for llm ques-	1152
1099	for speculative decoding algorithm. <i>arXiv preprint</i>	tion answering with external tools. <i>arXiv preprint</i>	1153
1100	<i>arXiv:2411.00841</i> .	<i>arXiv:2306.13304</i> .	1154
1101	Jiahui Yu, Yuanzhong Xu, Jing Yu Koh, Thang Lu-	Matthieu Zimmer, Milan Gritta, Gerasimos Lampouras,	1155
1102	ong, Gunjan Baid, Zirui Wang, Vijay Vasudevan,	Haitham Bou Ammar, and Jun Wang. 2024. Mixture	1156
1103	Alexander Ku, Yinfei Yang, Burcu Karagol Ayan,	of attentions for speculative decoding. <i>arXiv preprint</i>	1157
1104	et al. 2022. Scaling autoregressive models for	<i>arXiv:2410.03804</i> .	1158
1105	content-rich text-to-image generation. <i>arXiv preprint</i>		
1106	<i>arXiv:2206.10789</i> , 2(3):5.		
1107	Ofir Zafrir, Igor Margulis, Dorin Shteyman, and Guy		
1108	Boudoukh. 2024. Fastdraft: How to train your draft.		
1109	<i>arXiv preprint arXiv:2411.11055</i> .		
1110	Jun Zhang, Jue Wang, Huan Li, Lidan Shou, Ke Chen,		
1111	Gang Chen, and Sharad Mehrotra. 2023. Draft		
1112	& verify: Lossless large language model accelera-		
1113	tion via self-speculative decoding. <i>arXiv preprint</i>		
1114	<i>arXiv:2309.08168</i> .		
1115	Lefan Zhang, Xiaodan Wang, Yanhua Huang, and Rui-		
1116	wen Xu. 2024a. Learning harmonized represen-		
1117	tations for speculative sampling. <i>arXiv preprint</i>		
1118	<i>arXiv:2408.15766</i> .		