

Token Reduction Should Go Beyond Efficiency in Generative Models – From Vision, Language to Multimodality

Zhenglun Kong^{1,*}, Yize Li^{2,*}, Fanhu Zeng³, Lei Xin^{4,6}, Shvat Messica¹,
Xue Lin², Pu Zhao², Manolis Kellis⁵, Hao Tang⁶, Marinka Zitnik¹

¹Harvard University, ²Northeastern University, ³CAS,

⁴Wuhan University, ⁵MIT, ⁶Peking University,

{zhenglun_kong,marinka}@hms.harvard.edu, li.yize@northeastern.edu,

Abstract

In Transformer architectures, tokens—discrete units derived from raw data—are formed by segmenting inputs into fixed-length chunks. Each token is then mapped to an embedding, enabling parallel attention computations while preserving the input’s essential information. Due to the quadratic computational complexity of transformer self-attention mechanisms, token reduction has primarily been used as an efficiency strategy. This is especially true in single vision and language domains, where it helps balance computational costs, memory usage, and inference latency. Despite these advances, this paper argues that token reduction should transcend its traditional efficiency-oriented role in the era of large generative models. In this paper, we characterize this mechanism as a fundamental principle in generative modeling, critically influencing both model architecture and broader applications. We analyze how token reduction addresses critical challenges in current systems across vision, language, and multimodal, demonstrating its ability to: (i) facilitate deeper multimodal integration and alignment, (ii) mitigate "overthinking" and hallucinations, (iii) maintain coherence over long inputs, and (iv) enhance training stability, etc. We reframe token reduction as more than an efficiency measure. By doing so, we outline promising future directions, including algorithm design, reinforcement learning-guided token reduction, token optimization for in-context learning, agentic framework design, and broader ML and scientific domains.²

1 Introduction

Transformer-based generative models [14, 24, 39, 124] have emerged as dominant deep learning architectures across vision, language, and multimodal tasks, due to their ability to process long sequences of tokens, which are the fundamental representational units derived from raw data such as subwords in language or image patches in vision. As these models are applied to increasingly complex real-world tasks, the input sequence lengths of both the models and their training datasets continue to grow. However, the quadratic computational complexity of the attention mechanism results in high memory usage and slow inference, which hinders the practical deployment of generative models at scale. Token reduction addresses this challenge by reducing the number of tokens processed during inference. By pruning or merging tokens, token reduction [38, 41, 48, 51, 56, 63, 71, 84, 141, 167] reduces computational cost and accelerates runtime, providing a practical solution for enhancing generative efficiency [63, 85, 121, 136, 146, 152, 162, 164, 168, 172].

^{*}Equal contribution.

²We collected a list of token reduction papers at: [Awesome-Collection-Token-Reduction](#).

Token reduction has been widely adopted in computer vision, language processing, and multimodal tasks. In vision transformers, it has primarily been used to reduce computational cost by removing visually redundant tokens [11, 12, 28, 44, 67, 70, 79, 102]. In language models, token reduction has commonly been implemented through early-exit mechanisms and token-skipping strategies [81, 137], which reduce the number of intermediate tokens processed and thus lower computational overhead. Similarly, multimodal large language models (MLLMs) apply visual token pruning primarily during the prefill stage [22], where adaptive attention patterns are learned in the early layers to prune tokens in later stages. Despite progress, token reduction is predominantly viewed as a post-hoc efficiency optimization [86, 106], primarily by reducing the number of tokens to minimize associated computations and accelerate inference. Such an efficiency-only mindset has limitations. Naive pruning methods may discard informative tokens, thereby degrading model understanding and performance [79, 159, 160]. Furthermore, token reduction is commonly treated as a post hoc optimization, rather than being integrated into the core design and training of the model [22].

In this paper, we argue that viewing token reduction purely from an efficiency perspective is fundamentally limited. Instead, we position token reduction as a core design principle in generative modeling, deeply integrated with both training and inference to prioritize tokens that maximize downstream task performance and semantic integrity.

Modern generative tasks present numerous challenges that highlight the need for thoughtful token selection: (i) Ultra-long contexts in language modeling require selective retention of relevant segments to preserve coherence. (ii) LLMs frequently exhibit overthinking, repeatedly attending to low-value tokens and producing redundant or contradictory outputs. (iii) Multimodal generation tasks often face issues of visual redundancy, where background tokens overshadow salient visual features critical for accurate understanding. (iv) Noisy or irrelevant tokens introduced during training slow down convergence and harm model stability. By learning to intelligently select, merge, or compress tokens based on their contribution to generation objectives, rather than solely on raw redundancy, models can simultaneously reduce computational load, improve robustness, and enhance interpretability and alignment. This paper makes the following three key contributions:

- We categorize existing token reduction methods by their functional objective, identifying a transition from efficiency-centric optimizations to task-aware enhancements in vision, language, and multimodal domains.
- We identify core challenges faced by modern generative models including insufficient visual representation, semantic misalignment, overthinking in reasoning, and training instability. We then demonstrate how principled token reduction strategies can effectively mitigate these issues.
- We outline a roadmap for future research on token reduction, including directions for method design, reinforcement learning-guided token selection, adaptive in-context compression, and hardware-algorithm co-design, etc. These directions aim to support the development of next-generation generative architectures that are both robust and efficient.

This position paper is organized as follows: Sec. 2 reviews prior token reduction methods across various modalities. Sec. 3 introduces the problem formulation, Sec. 4 formalizes the identified challenges and demonstrates how informed token reduction strategies can address them. Sec. 5 proposes promising research directions for advancing token reduction as well as broader implications.

2 Related Work

2.1 Token Reduction in Vision Models

Image Classification. Classification serves as a fundamental task for vision models and token reduction techniques have been widely applied in it due to its simplicity and versatility. It has been widely explored from various aspects [12, 79, 102, 142, 156, 68, 157]. Specifically, DynamicViT [102] devises a lightweight module to predict the importance score of each token, thereby pruning unimportant tokens. SPViT [67] introduces a soft pruning technique, which integrates the less informative tokens generated by the selector module into a package token that will participate in subsequent calculations rather than being completely discarded. EViT [79] identifies attentive tokens from the attention map, enabling token pruning without additional parameters. ToMe [12] merges tokens with similarity based on bipartite matching to maintain information utility. PPT [142] analyzes the statistic

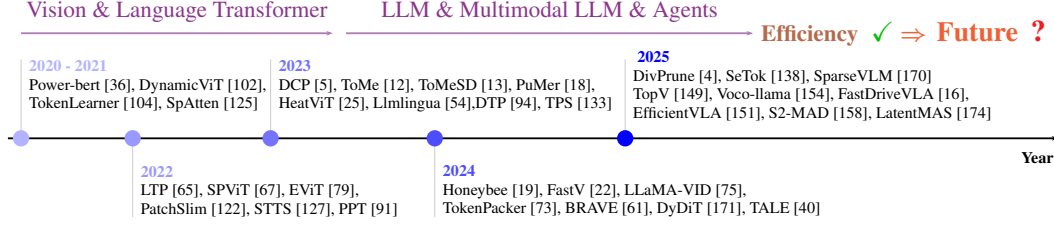


Figure 1: Timeline of notable method developments for token reduction methods with modality shifts (Vision & Language → Multimodal LLMs & Agents). All these strategies aim to speed up inference with negligible performance drops. Conversely, we ask: *What is the next token reduction paradigm in generative model design that goes beyond test-time accelerations?*

data between layers and adaptively employs token pruning and merging within layers to achieve higher acceleration performance.

Video Compression. Unlike token reduction in image classification, video compression focuses more on the temporal redundancy within videos, and algorithms are developed to reduce the number of tokens with less computational overhead. Various token reduction methods have been investigated for different tasks, including video understanding [104, 118, 127], video editing [74], video-text retrieval [87, 108], video action detection [21], and so on. Specifically, STTS [127] introduces a lightweight framework that dynamically selects the most informative spatial-temporal tokens in video transformers. Tokenlearner [104] proposes an adaptive tokenization module that learns a handful of informative spatial-temporal tokens, significantly reducing computational costs. EVAD [21] selectively drops irrelevant spatial-temporal tokens in non-keyframes while preserving keyframe and motion-relevant tokens, and then refines actor features using a context-aware decoder to maintain accuracy with reduced computations.

Generative Tasks. Token reduction in generative tasks [58] aims to accelerate generative models through the efficient utilization of tokens. It can be applied to both diffusion models [13, 89, 126] and diffusion transformers [31, 173]. Specifically, ToMeSD [13] exploits natural redundancy in generated images by merging redundant tokens, successfully extending token merging to stable diffusion with simple unmerging. DyDiT [171] reduces redundancy with a Timestep-wise Dynamic Width approach to adopt model width conditioned on the generation timesteps, and a Spatial-wise Dynamic Token strategy to avoid redundant computations at unnecessary spatial locations.

2.2 Token Reduction in Language Models

Token reduction strategies in language modeling have evolved from early optimizations for BERT [50, 62, 65, 66, 153] to techniques specifically designed for LLMs. PoWER-BERT [36] introduces progressive word-vector elimination by removing redundant token representations based on self-attention dynamics, improving inference efficiency. Learned token pruning [65] extends this approach by learning attention-based thresholds to adaptively prune uninformative tokens, thereby reducing computational costs while preserving model performance. In LLMs, token reduction must account for the constraints of autoregressive decoding across diverse downstream tasks. Dynamic pooling methods [5, 123] adjust token representations on the fly during inference to reduce redundancy. Prompt compression techniques [30, 34, 54] aim to reduce computational overhead by compressing the input prompt before generation. Selective decoding approaches [30, 135] reduce per-step inference costs by computing key-value pairs only for tokens critical to predicting the next token. In multi-agent systems, S²-MAD [158] proposes a sparsification mechanism that limits unnecessary token exchanges between agents, reducing communication costs and improving the efficiency of collaborative reasoning.

2.3 Token Reduction in Multimodal LLMs

Recent work has explored visual token pruning by addressing attention inefficiencies of deep transformer layers in MLLMs [4, 8, 81, 105]. Specifically, FastV [22] shows that deeper vision-language

layers expend significant computations on redundant image tokens. To address this, a lightweight module is adopted to adaptively prune these tokens, reducing inference overheads in subsequent stages. A complementary approach modifies vision feature extractors or projectors to output a smaller set of highly informative image tokens, effectively distilling the input into a compressed representation [15, 19, 61, 73, 76, 154]. However, efficiency gains from these prefill stages often fade during the decoding phase, where per-token computations dominate. To overcome this, recent methods jointly optimize token reduction during both prefill and decoding stages, ensuring sustained speedups throughout inference [47, 111]. Furthermore, the scope of token reduction has recently expanded to Vision-Language-Action (VLA) models, optimizing efficiency for real-time robotic manipulation and autonomous driving [16, 57, 100, 151].

In Fig. 1, we present a timeline of notable developments in token reduction methods, illustrating the shift from early applications in ViT and BERT-based models to more recent advances in LLMs, MLLMs and Agent systems.

3 Problem Formulation

In modern generative models [10, 37, 82, 96, 99], a token denotes one fundamental unit of input or representation, typically encoded as a vector. For example, a token might correspond to a subword in language, a patch in an image, or an embedding of a time step in audio. We denote a sequence of N input tokens as $X = [x_1, \dots, x_N] \in \mathbb{R}^d$. Token reduction refers to any operation that compresses the token sequence to M tokens (with $M < N$) by removing or consolidating tokens while aiming to preserve the original information.

Broadly, token reduction methods fall into four categories: 1) Token pruning methods [12] that remove entire unimportant tokens, simply dropping them from the sequence; 2) Token merging methods [12, 13] which fuse information from multiple tokens into fewer tokens, effectively compressing the sequence by merging similar or related tokens; 3) Hybrid strategies [18, 64, 142] that combine pruning and merging within a unified framework; 4) Token distillation approaches [15, 93] which integrate rich information across longer input sequences or multiple modalities into fewer condensed tokens, enabling efficient cross-modal interactions and long-context reasoning in LLMs and MLLMs.

A core challenge in token reduction is the determination of tokens to be pruned or merged. There are various importance criteria and scoring mechanisms to rank token significance, including attention-based heuristics [79], gradient or loss-based criteria [47], clustering [42], and learned predictors [102].

From a purely efficiency-oriented perspective, token reduction delivers substantial computational efficiency gains by reducing the quadratic computation cost from $\mathcal{O}(N^2)$ to $\mathcal{O}(M^2)$ in attention mechanisms. By eliminating redundant tokens and processing fewer computations during inference, it effectively accelerates the inference speed and improves the model throughput, which is crucial for latency-sensitive tasks or real-time applications. Furthermore, it reduces the memory footprint of activations and gradients (e.g., key/value caches), alleviating memory usage for both inference and training, which is particularly beneficial for wide-scale deployments on resource-limited platforms. We present more theoretical detail in the Appendix A

However, as stated in this position paper, token reduction can benefit models in multiple ways beyond efficiency, which will be introduced in detail in the following sections.

4 Core Roles and Challenges

In this section, we discuss token reduction as a foundational mechanism for addressing critical challenges in modern generative systems. We categorize five core challenges across modalities: visual representation sparsity, semantic misalignment, reasoning redundancy, training instability, and long-context overload. We demonstrate how principled token reduction strategies intrinsically address these issues through dynamic token-semantic co-optimization. We position token reduction not only as an efficiency tool, but as an essential paradigm for enhancing semantic coherence and enabling sustainable scaling of generative systems.

4.1 Obtain Informative Visual Representation

MLLMs often suffer from noisy visual inputs that impede fine-grained understanding. We outline key challenges in MLLM visual reasoning: ① *Text-Visual Attention Shift*: Due to the rotary positional embeddings in LLM decoders, later text tokens disproportionately attend to spatially lower image regions [45], shifting attention away from semantically important areas (e.g., objects at the top of an image); ② *Visual Redundancy*: Empirical studies [81, 137] show that beyond the first few layers, many image tokens contribute little new information, ③ *Task-Guided Focus in VQA*: In multimodal question answering, the question itself pinpoints relevant image regions (e.g., "kitten color" directs focus to the kitten patch), implying that many image tokens are unnecessary for correct answers [111].

Therefore, token reduction can serve as a representation-learning optimization: selecting the subset of tokens that preserves informative visual representation. For example, VisPruner [166] identifies high-value tokens using visual-encoder attention and removes duplicates via clustering to ensure diversity. VTW [81] observes that visual information migrates into text tokens within early layers; it therefore withdraws all visual tokens after a chosen layer based on KL-divergence criteria. TRIM [111] leverages the CLIP metric and IQR scoring function to adaptively select image tokens that are crucial for answering questions, while an aggregated token is used to retain additional image information.

4.2 Better Multimodal Token Alignment

Despite their impressive capabilities, MLLMs continue to face challenges in semantic alignment. Standard vision tokenizers typically split images into fixed-size patches, which can fragment coherent visual entities (e.g., objects or regions) across multiple tokens. This fragmentation weakens the alignment between visual and linguistic representations. Token reduction offers a promising solution by selecting visual tokens based on semantic importance, thereby producing a compact set of tokens that better align with language representations. Specifically, SeTok [138] dynamically clusters visual features into semantically meaningful tokens using a density-peak algorithm, which determines both the number and structure of token groupings per image. This approach preserves both high- and low-frequency semantics, substantially improving concept-level alignment and downstream task performance. M3 [15] introduces a hierarchical token structure that captures coarse-to-fine semantic granularity, allowing different levels of abstraction to be selectively retained depending on task needs.

4.3 Reduce Overthinking in Reasoning

LLM reasoning. In the context of language models, overthinking refers to generating excessively long or convoluted chains of reasoning that go beyond what is necessary to reach a correct answer. An LLM may produce verbose, repetitive, or even self-contradictory explanations when it fails to converge on a solution—often due to uncertainty [117, 129]. Such extended reasoning trajectories are inefficient and recent studies show that state-of-the-art reasoners can consume over 15,000 tokens to solve math problems that could be addressed with a concise chain-of-thought (CoT) of just a few hundred tokens [46]. This issue is particularly acute in LLM agents, where internal reasoning alternates with external tool use [32, 128]; excessive steps can obscure logical clarity and lead to error accumulation. Mitigating overthinking is thus crucial. By trimming unnecessary tokens, LLMs can focus on salient steps, aligning generation with a more concise trajectory.

CoT-Influx [49] introduces a CoT pruning strategy in which concise reasoning examples are included in the prompt. By pruning unimportant tokens from these examples, more reasoning demonstrations can fit into the context window, surprisingly leading to improved math reasoning accuracy. Token-Skip [143] enables LLMs to skip less important tokens within CoT sequences and learn shortcuts between critical reasoning steps. This allows for controllable CoT compression with adjustable compression ratios, enabling models to automatically trim redundant tokens during reasoning.

MLLM reasoning. MLLMs, which reason over text and other modalities, face similar overthinking issues. In vision-language tasks, overthinking often manifests as excessive processing of visual tokens or overly detailed image descriptions, resulting in inefficiency and potential confusion [20]. Token reduction techniques in MLLMs aim to promote more focused and sparse reasoning over multimodal inputs. For example, FAST [144] rewards shorter-than-average token sequences for correct answers, while allowing longer reasoning for more complex tasks. It also adjusts policy optimization constraints to tighten output exploration for simple tasks (thus reducing unnecessary tokens) and loosen it for harder ones to allow deeper reasoning.

Together, these strategies reduce overthinking in straightforward cases, boosting efficiency while preserving effective reasoning depth for complex scenarios.

4.4 Improve Training Stability & Efficiency

While token reduction has traditionally been employed as a post-training optimization to enhance inference efficiency, recent research indicates its potential to significantly improve training stability when integrated into the pre-training phase [31, 77, 80], suggesting that selective token utilization during training can lead to more robust model learning.

One notable approach is Rho-1 [80], which involves scoring tokens based on their alignment with a desired distribution using a reference model and then focusing the training loss on tokens with higher scores. Therefore, it effectively filters out noisy or less informative tokens, leading to faster convergence and improved performance. UPFT [53] emphasizes the importance of initial reasoning steps in training. By reducing the number of training tokens, UPFT encourages the model to focus on the initial prefix substrings of reasoning trajectories, which are often more stable and contain crucial information. This focus helps the model avoid being influenced by subsequent complex or potentially erroneous information, thereby improving training stability.

Additionally, integrating token reduction with training procedures like GRPO [107] is gaining traction; for instance, recent work reveals that optimizing only a subset of high-entropy "forking tokens" matches full-gradient updates [130], suggesting that entropy patterns can effectively guide efficient policy learning. Future research should investigate specialized approaches that incorporate token reduction directly into training objectives, enabling models to learn to prioritize or discard tokens in a task-aware and gradient-aligned manner.

4.5 Enhance Long Context & Video Understanding

Long-context LLMs. Long-context language modeling presents unique challenges: ① Long texts often contain raw tokens that exhibit repetitive descriptions and irrelevant details that strain the attention mechanism; ② LLM-based agent systems use input data as sequential prompts for reasoning or for switching between multiple tasks, which can lead to overload when the prompt grows too large; ③ It is very difficult to scale up to even longer content for learning more information. Token reduction techniques directly address these issues by distilling extensive input sequences into compact summary vectors or representative tokens. By doing so, models preserve core information such as key events, central themes, or task-specific facts, while significantly decreasing cognitive load. For example, AutoCompressors [23] trains pre-trained LLMs to compress long contexts into compact summary tokens, reducing token length by orders of magnitude to extend context windows and speed up inference. TokenSwift [140] reduces the effective number of tokens that the model dynamically processes during generation by using multi-token parallel generation and n-gram retrieval for token reutilization, therefore enabling efficient ultra-long sequence generation (up to 100K tokens).

Video-based MLLMs. The necessity of token reduction primarily lies in enhancing the model's effective understanding of video content through: ① *Instruction-guided information filtering*: token reduction prioritizes selecting visual information relevant to user instructions over raw data volume. ② *preserving spatiotemporal structure*: token reduction strategically compresses massive spatiotemporal information to retain spatiotemporal dependencies, ensuring the model can capture dynamic semantics, as well as prevent redundant tokens interfere with long temporal reasoning. ③ *Preserving semantic integrity*: it facilitates feasible processing of extremely long sequences in learning while preserving semantic integrity. ④ *Multi-modal alignment*: token reduction distills visual information into a compact, semantically aligned form, thereby efficiently bridging the gap between language and vision [88]. By doing so, it effectively addresses the challenges posed by the low abstractness and lack of guidance inherent in raw visual inputs, which are the root causes of semantic misalignment and optimization ambiguity in multi-modal models. Recent works illustrate these principles: HICom [88] conducts conditional token compression at local and global levels using user instructions as guidance to retain instruction-relevant visual information while reducing computational burden. Video-XL-Pro [83] employs reconstructive token compression with a dynamic token synthesizer and semantic-guided masking to generate compact yet comprehensive video tokens for improved MLLM performance and efficiency.

5 Future Directions

In this section, we propose eight promising directions for token reduction beyond the efficiency benefits, organized into three categories: (i) Algorithmic Innovations (Sec. 5.1~5.4), (ii) Application Innovations (Sec. 5.5~5.9), and (iii) Hardware-Algorithm Co-Design (Sec. 5.7).

5.1 Design of New Algorithms

Future research on algorithm design should explore holistic and adaptive token reduction strategies. Building on recent advances, we outline six promising directions:

Better Token Importance Metrics. It is critical to re-evaluate how token importance is defined and measured. More robust and unbiased scoring mechanisms can be developed, such as predictors [3] or meta-learning frameworks that go beyond attention-based proxies. These models should capture downstream utility with minimal supervision, enabling adaptive pruning across tasks and domains.

Constructive Token Compression. Token reduction can shift from purely eliminative pruning to strategies that merge spatially or semantically similar tokens into compact summary vectors [73].

Mitigating Position Bias. In MLLMs, attention-based pruning methods (e.g., FastV) often rely on attention scores from a fixed query token, leading to retained tokens concentrating in specific image regions (e.g., lower corner) [134] with potential position bias. Future methods should preserve spatial diversity by enforcing structural uniformity in retained tokens to improve robustness on visual tasks.

Cross-Modal Guided Pruning. Pruning decisions in MLLMs should be guided by inter-modality dependencies, rather than made independently for each modality. For example, text-guided pruning of visual tokens can improve alignment between modalities [17]. The design should account for joint representations and semantic correspondence across all relevant inputs.

End-to-End Sparsification. Token reduction should consider both the prefill stage and decoding phase for LLMs. This includes dynamically managing the sparsity of KV caches and selectively updating generated tokens, sustaining efficiency gains throughout the entire inference process [47].

Hardware-Algorithm Co-Design. Token pruning can explore custom hardware and compiler optimizations that take advantage of dynamic token sparsity patterns (e.g., irregular memory access and conditional computation) to maximize throughput and energy efficiency as detailed in Sec. 5.7.

5.2 From Prompt Tuning to Chain of Thought Reasoning

Current token reduction efforts for prompts have primarily aimed at compressing prompts for efficiency, often with impressive results [54, 93]. Looking forward, token reduction should evolve into enhancing reasoning and maximizing utility per token in context. Instead of focusing solely on making prompts shorter, future research should explore how each remaining token can carry more information or trigger more complex inference during in-context learning and chain of thought reasoning. One direction is to alter the generation paradigm itself, for example, training language models to predict multiple tokens per step [35]. Another idea is to enable deeper internal reasoning without increasing prompt length [1].

As mentioned in Sec. 4, long CoT chains can become verbose: excessive reasoning steps may introduce errors or obscure logical clarity, particularly in LLM agents where internal reasoning alternates with external tool use [32, 128]. Token reduction may serve a critical role in this context by compressing intermediate reasoning into a compact representation. For example, approaches based on next-token prediction [93, 165] can distill intermediate thinking chains into a set of dense, information-rich tokens. These compressed representations can then replace the full intermediate context and serve as inputs for subsequent reasoning steps.

This compressed-thinking strategy has two main benefits: 1) reducing error accumulation and keeping the logic clear by focusing on key information, and 2) allowing more reasoning rounds to fit within a fixed context window, enabling deeper multi-step inference without exceeding length limits.

In summary, the next phase of token reduction research should shift focus from simple prompt compression to reasoning-centric compression. Rather than just trimming prompts, we should ask: *How can we make each token in the prompt or context do more work for us?* This involves training

models with objectives that reward higher-level inference per token, developing architectures that recycle tokens for multi-step thinking, or dynamically selecting the most salient tokens to keep at each step of reasoning.

5.3 Efficient Reasoning with Reinforcement Learning

Reinforcement-learning (RL)-driven token reduction has shown strong promise for improving reasoning efficiency in both LLMs and MLLMs. The key challenge is to balance compute and reasoning quality via dynamic reward design, sparsity-inducing constraints, and adaptive control of effective token length [2, 90, 139].

In language reasoning, length awareness is incorporated either by diversifying reasoning formats during pre-training [113] or by adding explicit length penalties in the RL stage [9, 46]. While effective, most approaches implicitly assume static task complexity or rely on hand-specified length constraints, which can be suboptimal under heterogeneous workloads. A natural next step is *adaptive reasoning*: using RL to learn per-instance budget allocation from intrinsic task difficulty [95, 101]. In parallel, performing reasoning in a compressed latent space can yield substantial computational savings [110], but current methods often degrade due to poorly structured latent representations. A promising direction is to inject *explicit logical structures* into the latent space to enable more controllable and compositional reasoning [43].

Beyond language modality, under the “Fast-Slow Thinking” framework [144], RL can supervise hierarchical selection of high-value tokens: a fast branch applies sparsity signals (e.g., rule-based rewards or information-density scoring) to prune redundant visual/semantic features, while a slow branch allocates computation to refined reasoning. Additionally, RL enables a Think-with-Image paradigm [116] by letting models adapt visual granularity: VisionThink [150] adopts a progressive resolution strategy, using RL to selectively request high-resolution inputs only for fine-grained tasks like OCR. PixelThink [131] addresses visual overthinking in segmentation by adjusting reasoning length based on task difficulty. These methods demonstrate how RL can dynamically calibrate both input *resolution* and reasoning *depth* across various tasks. Looking ahead, the integration of such approaches could enhance cross-modal alignment and inference efficiency in real-time and resource-constrained scenarios, supporting a new generation of lightweight yet capable multimodal language agents.

5.4 Complementary to Other Methods

Token reduction can complement other efficiency techniques, such as quantization. By selectively reducing the number of tokens processed during inference, models can improve both performance and efficiency, particularly when paired with quantization strategies [72]. Traditional key-value cache quantization methods often suffer from accuracy loss due to their inability to handle outlier tokens that carry distinct or rare features. To mitigate this issue, Outlier Token Tracking [115] identifies outlier tokens during decoding and excludes them from quantization, preserving full-precision necessary representations and improving key-value cache quantization accuracy. Similarly, Agile-Quant [109] incorporates token pruning as a preprocessing step to reduce the impact of activation outliers. It prunes tokens based on their attention to the start-of-sequence token, discarding those with low attentiveness, which often appear in adjacent channels and contribute to quantization noise. This targeted pruning reduces interaction distances between salient tokens and helps maintain model accuracy under low-bit quantization settings.

5.5 Towards Dense Prediction Tasks for Vision

Existing works primarily concentrate on compressing the backbone of models to ensure their generalization ability, and few works explore recovering all tokens for dense prediction tasks [13, 78]. It is necessary to develop custom token reduction methods for various downstream dense prediction applications like autonomous driving and robotic control with specific settings and requirements [16, 100]. Lacking these specialized designs would lead to a mismatch and performance drop when deployed in real-world settings. For example, autonomous driving [161] would require displacement and velocity based on occupancy prediction, and robotic control [132] would demand rotation angle according to the grid map. Therefore, how to develop fast and specialized token reduction strategies tailored for downstream dense prediction tasks is crucial for deployment in practical scenarios.

5.6 Towards Long Video Applications

Exploiting long videos holds great potential, as processing hours of footage is significantly more labor-intensive and time-consuming than working with short clips. Due to the inherent complexity and resource demands, most current research on long video learning focuses on discriminative tasks such as video understanding [69, 103]. In contrast, broader applications including long video editing [169], long video-text retrieval, and narrative-level generation remain largely underexplored. Progress in these areas could have a significant impact on scene editing in video clips, character rendering in movies, and retrieving useful information from numerous videos.

Moreover, token reduction offers a path toward interpretability and efficiency in long video processing [55]. This mimics the human visual system, which does not attend to every frame in detail but instead focuses on salient spatiotemporal changes, such as actions or object movement, while filtering out static, redundant content like backgrounds and stationary objects. Future models should similarly prioritize informative frames and temporal segments, allowing them to reason over extended video sequences with greater efficiency and interpretability.

5.7 Algorithm-Hardware Co-Design

While algorithmic advancements in token reduction have achieved impressive computational savings, the next crucial step is to integrate these techniques with hardware-aware design principles. We posit that algorithm-hardware co-design is essential for holistic optimization across the compute stack, considering the interplay between algorithmic choices, hardware architectures (specialized data paths, memory hierarchies, communication fabrics, control logic, etc.), and compiler/runtime support (efficient sparse mapping, dynamic scheduling, irregular-data management, etc.) [25, 98].

Currently, co-design efforts targeted at token reduction lag significantly behind pure algorithmic research. This gap is problematic because hardware design needs to balance PPA (power, performance, and area), platform specifics, data movement costs, control overhead, and scalability/reusability [97]. Algorithms developed in isolation often generate sparse or irregular compute patterns that general-purpose hardware cannot exploit effectively. Therefore, future research should aim to: 1) Design parameterizable, reconfigurable accelerator modules-such as on-the-fly importance-scoring units and sparse-data pipelines-that natively support token-reduced Transformers. 2) Explore Processing-in-Memory (PIM) architectures to alleviate severe memory bottlenecks caused by dynamic token pruning. By executing scoring operations or partial attention mechanisms within or near memory arrays, PIM can drastically reduce data movement costs and improve end-to-end efficiency.

5.8 Towards Efficient Agentic Systems

Dynamic Memory & Context Engineering. AI agents face a severe context bottleneck where accumulating interaction history not only incurs quadratic computational costs but also degrades reasoning through "lost-in-the-middle" phenomena. Transitioning to active context management is essential [7]; systems must distinguish between immutable instructions and transient episodic data. ACON [59] demonstrate that semantic memory compression can reduce peak token usage while maintaining long-horizon performance. Complementary strategies include hierarchical memory abstraction, which offloads older history to retrieval-based storage, and hybrid masking techniques [52]. These approaches treat tokens as a finite resource to be optimized on-the-fly, dynamically adjusting the sparsity of the context window based on the complexity of the current reasoning step.

Observation & Tool Pruning. Agents must often process verbose tool outputs (e.g., massive JSON/HTML). Feeding raw data is inefficient; future research should focus on token-aware interaction, where lightweight scorers filter observations to retain only task-relevant features. This preserves critical context bandwidth without sacrificing trajectory accuracy, as evidenced in recent programming agents [145]. Furthermore, adaptive truncation strategies that prioritize error traces over standard success logs can significantly improve debugging capabilities while minimizing token consumption.

Communication-Efficient Multi-Agent Systems. At the system level, the aggregate token cost scales with the number of interacting nodes. Enforcing strict token budgets on inter-agent exchange compels information-dense communication. This sparse protocol enhances scalability by reducing redundant chatter and mitigating hallucination loops in collaborative workflows [163, 174].

5.9 Towards AI for Broader ML and Scientific Domains

Token reduction methods can also offer powerful opportunities to reshape broader machine learning and scientific applications. In particular, domains such as medicine, biology, chemistry, and temporal data analysis frequently encounter complex data structures, heterogeneous data sources, and intricate domain-specific relationships. Informed tokenization approaches promise to address these challenges by transforming complex and rich scientific data into concise, informative, and flexible representations, significantly enhancing the utility of transformer-based foundation models across these domains.

Building Biomedical Tokenizers. Recent works exemplify the transformative potential of advanced tokenization methods in the biomedical domain, including protein [33, 119, 155], genomic [27], and chemical structure [147] tokenizers. Collectively, these methods illustrate how informed reduction and condensation of input tokens can lead to more effective and interpretable scientific models. For example, traditional tokenizers in EHR foundation models typically treat medical codes as isolated textual units, neglecting their inherent structured and relational context, such as hierarchical relationships, disease co-occurrences, and drug-treatment associations found within biomedical ontologies. To solve this issue, MedTok [114] integrates textual descriptions and graph-based relational data into a unified tokenization framework. It first uses a language model encoder to extract embeddings from medical code descriptions and employs a graph encoder to capture relational structures from biomedical ontologies. These embeddings are combined into a compact token space through vector quantization, preserving both modality-specific and cross-modality information.

To enhance informativeness and reduce redundancy, MedTok employs a token packing mechanism. It optimizes shared tokens and modality-specific tokens, ensuring that the final tokens encode both shared semantic meaning and modality-specific structure. This process drastically reduces effective vocabulary size, addressing the scalability challenge of 600,000+ medical codes by collapsing redundant representations while preserving critical clinical context. Inspired by adaptive tokenization methods for vision [26, 60, 148], future EHR tokenization would be adaptive, enabling the dynamic representation of patients’ medical histories, where the length of the token series for each patient’s history would be directly correlated with the length and complexity. Such adaptive tokenization can significantly improve training and inference efficiency across diverse healthcare systems.

Time-Series Data and Clinical Reasoning. Temporal dynamics form an essential component of clinical reasoning, particularly through longitudinal patient data like lab tests and vital signs. However, current large language models struggle to effectively incorporate time-series inputs due to challenges in temporal tokenization [112, 6, 29, 112, 92]. Future tokenization methods should not only dynamically adjust the number of tokens according to temporal complexity but also selectively focus on time segments most relevant to the clinical context, prompt, or task at hand [120]. This could enhance training effectiveness and inference accuracy, helping create the next generation of EHR foundation models, which are flexible not only over different tasks or prompts, but also over different data sources, patients, and populations. The complexity and richness of EHR data offer opportunities for AI-driven advancements in patient health outcomes. Future EHR models should support comprehensive reasoning capabilities, encompassing complete patient histories, such as vitals, lab results, diagnoses, and procedures over time. They could facilitate timely disease predictions, accurately forecast chronic disease trajectories, and anticipate patient responses to treatments.

6 Conclusion

In this position paper, we have argued that token reduction must evolve beyond a mere efficiency optimization to become a core design principle in generative modeling. We have shown how principled token reduction can address key challenges such as enhancing semantic fidelity in vision-language alignment, curbing verbose reasoning trajectories, preserving long-range coherence, and stabilizing learning dynamics. Looking forward, the roadmap we outlined points to a broad landscape of opportunities, ranging from algorithmic innovations and hardware-algorithm co-design to specialized applications in scientific domains. We anticipate that future work will increasingly focus on constructive compression and reinforcement learning-guided selection, enabling models to autonomously optimize their information bandwidth. Ultimately, by treating token reduction as a holistic and task-aware mechanism, the community can develop next-generation systems that effectively balance scalability with effectiveness, interpretability, and performance.

7 Limitations

While token reduction offers significant benefits, our review identifies critical limitations and trade-offs that must be considered to avoid indiscriminate application.

Information Loss in Dense Prediction Token reduction methods, particularly pruning, inherently discard information. While this is acceptable for semantic classification or generation, it poses severe risks for dense prediction tasks (e.g., segmentation, object detection) or medical analysis where fine-grained spatial details are crucial. Merging strategies like ToMe [13] mitigate this better than pruning, but artifacts often remain at high compression ratios. In scenarios requiring pixel-perfect reconstruction, the trade-off between reduction and precision often favors preserving the full token set. Furthermore, the lack of specialized designs for token recovery often leads to performance mismatches in real-world applications like autonomous driving and robotic control. Future research must therefore explore custom reconstruction mechanisms to ensure these methods can meet the rigorous demands of dense tasks.

Overhead vs. Gain Dynamic token reduction introduces computational overhead (e.g., scoring networks, predictors). For short sequences or small batch sizes, the cost of computing token importance may outweigh the savings from processing fewer tokens. Furthermore, unstructured pruning can lead to irregular memory access patterns that are inefficient on standard hardware (GPUs/TPUs), potentially negating theoretical FLOPs reductions.

Alternatives: Reduction vs. Retrieval Critics may argue that techniques like Retrieval-Augmented Generation (RAG) or simply scaling context windows render token reduction unnecessary. However, we argue these are complementary. While RAG selects *documents*, token reduction operates at the *sub-document* level, filtering noise within relevant chunks. Similarly, while larger context windows allow for more input, they exacerbate the "lost-in-the-middle" phenomenon; token reduction acts as an attention-sharpening mechanism that helps models utilize these long contexts more effectively.

References

- [1] Griffin Adams, Alexander R Fabbri, Faisal Ladhak, Eric Lehman, and Noémie Elhadad. From sparse to dense: Gpt-4 summarization with chain of density prompting. *EMNLP, 4th New Frontier Summarization Workshop*, 2023.
- [2] Pranjal Aggarwal and Sean Welleck. L1: Controlling how long a reasoning model thinks with reinforcement learning. *COLM*, 2025.
- [3] Yash Akhauri, Ahmed F AbouElhamayed, Yifei Gao, Chi-Chih Chang, Nilesh Jain, and Mohamed S Abdelfattah. Tokenbutler: Token importance is predictable. *arXiv preprint arXiv:2503.07518*, 2025.
- [4] Saeed Ranjbar Alvar, Gursimran Singh, Mohammad Akbari, and Yong Zhang. Divprune: Diversity-based visual token pruning for large multimodal models. *CVPR*, 2025.
- [5] Sotiris Anagnostidis, Dario Pavllo, Luca Biggio, Lorenzo Noci, Aurelien Lucchi, and Thomas Hofmann. Dynamic context pruning for efficient and interpretable autoregressive transformers. *NeurIPS*, 2023.
- [6] Md Fahim Anjum. Lipcot: Linear predictive coding based tokenizer for self-supervised learning of time series data via language models. *arXiv preprint arXiv:2408.07292*, 2024.
- [7] Anthropic. Effective context engineering for ai agents. <https://www.anthropic.com/engineering/effective-context-engineering-for-ai-agents>, 2025.
- [8] Kazi Hasan Ibn Arif, JinYi Yoon, Dimitrios S Nikolopoulos, Hans Vandierendonck, Deepu John, and Bo Ji. Hired: Attention-guided token dropping for efficient inference of high-resolution vision-language models. In *AAAI*, 2025.
- [9] Daman Arora and Andrea Zanette. Training language models to reason efficiently. *NeurIPS*, 2025.
- [10] Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibao Song, et al. Qwen2.5-vl technical report. *arXiv preprint arXiv:2502.13923*, 2025.
- [11] Benjamin Bergner, Christoph Lippert, and Aravindh Mahendran. Token crop: Faster vits for quite a few tasks. *CVPR*, 2025.

- [12] Daniel Bolya, Cheng-Yang Fu, Xiaoliang Dai, Peizhao Zhang, Christoph Feichtenhofer, and Judy Hoffman. Token merging: Your vit but faster. *ICLR*, 2023.
- [13] Daniel Bolya and Judy Hoffman. Token merging for fast stable diffusion. In *CVPR*, 2023.
- [14] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *NeurIPS*, 2020.
- [15] Mu Cai, Jianwei Yang, Jianfeng Gao, and Yong Jae Lee. Matryoshka multimodal models. In *ICLR*, 2025.
- [16] Jiajun Cao, Qizhe Zhang, Peidong Jia, Xuhui Zhao, Bo Lan, Xiaohan Zhang, Zhuo Li, Xiaobao Wei, Sixiang Chen, Liyun Li, et al. Fastdrivevla: Efficient end-to-end driving via plug-and-play reconstruction-based token pruning. *arXiv preprint arXiv:2507.23318*, 2025.
- [17] Jianjian Cao, Peng Ye, Shengze Li, Chong Yu, Yansong Tang, Jiwen Lu, and Tao Chen. Madtp: Multimodal alignment-guided dynamic token pruning for accelerating vision-language transformer. In *CVPR*, 2024.
- [18] Qingqing Cao, Bhargavi Paranjape, and Hannaneh Hajishirzi. Pumer: Pruning and merging tokens for efficient vision language models. *ACL*, 2023.
- [19] Junbum Cha, Wooyoung Kang, Jonghwan Mun, and Byungseok Roh. Honeybee: Locality-enhanced projector for multimodal llm. In *CVPR*, 2024.
- [20] Feng Chen, Yefei He, Lequan Lin, Jing Liu, Bohan Zhuang, and Qi Wu. Zipr1: Reinforcing token sparsity in mllms. *arXiv preprint arXiv:2504.18579*, 2025.
- [21] Lei Chen, Zhan Tong, Yibing Song, Gangshan Wu, and Limin Wang. Efficient video action detection with token dropout and context refinement. In *ICCV*, 2023.
- [22] Liang Chen, Haozhe Zhao, Tianyu Liu, Shuai Bai, Junyang Lin, Chang Zhou, and Baobao Chang. An image is worth 1/2 tokens after layer 2: Plug-and-play inference acceleration for large vision-language models. In *ECCV*. Springer, 2024.
- [23] Alexis Chevalier, Alexander Wettig, Anirudh Ajith, and Danqi Chen. Adapting language models to compress contexts. *EMNLP*, 2023.
- [24] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 conference of the North American chapter of the association for computational linguistics: human language technologies, volume 1 (long and short papers)*, 2019.
- [25] Peiyan Dong, Mengshu Sun, Alec Lu, Yanyue Xie, Kenneth Liu, Zhenglun Kong, Xin Meng, Zhengang Li, Xue Lin, Zhenman Fang, et al. Heatvit: Hardware-efficient adaptive token pruning for vision transformers. In *2023 IEEE International Symposium on High-Performance Computer Architecture (HPCA)*. IEEE, 2023.
- [26] Shivam Duggal, Phillip Isola, Antonio Torralba, and William T Freeman. Adaptive length image tokenization via recurrent allocation. In *First Workshop on Scalable Optimization for Efficient and Adaptive Foundation Models*, 2024.
- [27] Bell Raj Eapen. Genomic tokenizer: Toward a biology-driven tokenization in transformer models for dna sequences. *bioRxiv*, 2025.
- [28] Haipeng Fang, Sheng Tang, Juan Cao, Enshuo Zhang, Fan Tang, and Tong-Yee Lee. Attend to not attended: Structure-then-detail token merging for post-training dit acceleration. *CVPR*, 2025.
- [29] Liri Fang, Yuncong Chen, Wenchao Yu, Yanchi Liu, Lu-an Tang, Vette I Torvik, and Haifeng Chen. Tsia: A multi-task time series language model. In *ICASSP 2025-2025 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2025.
- [30] Qichen Fu, Minsik Cho, Thomas Merth, Sachin Mehta, Mohammad Rastegari, and Mahyar Najibi. Lazyllm: Dynamic token pruning for efficient long context llm inference. *arXiv preprint arXiv:2407.14057*, 2024.
- [31] Shanghua Gao, Pan Zhou, Ming-Ming Cheng, and Shuicheng Yan. Masked diffusion transformer is a strong image synthesizer. In *ICCV*, 2023.

- [32] Shanghua Gao, Richard Zhu, Zhenglun Kong, Ayush Noori, Xiaorui Su, Curtis Ginder, Theodoros Tsiligkaridis, and Marinka Zitnik. Txagent: An ai agent for therapeutic reasoning across a universe of tools. *arXiv preprint arXiv:2503.10970*, 2025.
- [33] Zhangyang Gao, Cheng Tan, Jue Wang, Yufei Huang, Lirong Wu, and Stan Z Li. Foldtoken: Learning protein language via vector quantization and beyond. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2025.
- [34] Tao Ge, Jing Hu, Lei Wang, Xun Wang, Si-Qing Chen, and Furu Wei. In-context autoencoder for context compression in a large language model. *ICLR*, 2024.
- [35] Fabian Gloeckle, Badr Youbi Idrissi, Baptiste Rozière, David Lopez-Paz, and Gabriel Synnaeve. Better & faster large language models via multi-token prediction. *arXiv preprint arXiv:2404.19737*, 2024.
- [36] Saurabh Goyal, Anamitra Roy Choudhury, Saurabh Raje, Venkatesan Chakaravarthy, Yogish Sabharwal, and Ashish Verma. Power-bert: Accelerating bert inference via progressive word-vector elimination. In *ICML*. PMLR, 2020.
- [37] Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
- [38] Yiju Guo, Wenkai Yang, Zexu Sun, Ning Ding, Zhiyuan Liu, and Yankai Lin. Learning to focus: Causal attention distillation via gradient-guided token pruning. *NeurIPS*, 2025.
- [39] Kai Han, An Xiao, Enhua Wu, Jianyuan Guo, Chunjing XU, and Yunhe Wang. Transformer in transformer. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan, editors, *NeurIPS*, 2021.
- [40] Tingxu Han, Zhenting Wang, Chunrong Fang, Shiyu Zhao, Shiqing Ma, and Zhenyu Chen. Token-budget-aware llm reasoning. *ACL*, 2025.
- [41] Yuhang Han, Xuyang Liu, Zihan Zhang, Pengxiang Ding, Donglin Wang, Honggang Chen, Qingsen Yan, and Siteng Huang. Filter, correlate, compress: Training-free token reduction for mllm acceleration. *AAAI*, 2026.
- [42] Joakim Bruslund Haurum, Sergio Escalera, Graham W Taylor, and Thomas B Moeslund. Agglomerative token clustering. In *European Conference on Computer Vision*. Springer, 2024.
- [43] Lukas Helff, Ruben HÅd’rle, Wolfgang Stammer, Felix Friedrich, Manuel Brack, Antonia WÅijst, Hikaru Shindo, Patrick Schramowski, and Kristian Kersting. Activationreasoning: Logical reasoning in latent activation spaces. *arXiv preprint arXiv:2510.18184*, 2025.
- [44] Cheng-Yao Hong and Tyng-Luh Liu. Multimodal promptable token merging for diffusion models. *AAAI*, 2025.
- [45] Xiangyu Hong, Che Jiang, Biqing Qi, Fandong Meng, Mo Yu, Bowen Zhou, and Jie Zhou. On the token distance modeling ability of higher rope attention dimension. *EMNLP Findings*, 2024.
- [46] Bairu Hou, Yang Zhang, Jiabao Ji, Yujian Liu, Kaizhi Qian, Jacob Andreas, and Shiyu Chang. Thinkprune: Pruning long chain-of-thought of llms via reinforcement learning. *arXiv preprint arXiv:2504.01296*, 2025.
- [47] Wenxuan Huang, Zijie Zhai, Yunhang Shen, Shaosheng Cao, Fei Zhao, Xiangfeng Xu, Zheyu Ye, Yao Hu, and Shaohui Lin. Dynamic-llava: Efficient multimodal large language models via dynamic vision-language context sparsification. *ICLR*, 2025.
- [48] Xiaohu Huang, Hao Zhou, and Kai Han. Prunevid: Visual token pruning for efficient video large language models. *ACL*, 2024.
- [49] Xijie Huang, Li Lina Zhang, Kwang-Ting Cheng, Fan Yang, and Mao Yang. Fewer is more: Boosting llm reasoning with reinforced context pruning. *EMNLP*, 2024.
- [50] Xin Huang, Ashish Khetan, Rene Bidart, and Zohar Karnin. Pyramid-bert: Reducing complexity via successive core-set based token selection. *ACL*, 2022.
- [51] Jeongseok Hyun, Sukjun Hwang, Su Ho Han, Taeh Kim, Inwoong Lee, Dongyoon Wee, Joon-Young Lee, Seon Joo Kim, and Minho Shim. Multi-granular spatio-temporal token merging for training-free acceleration of video llms. *ICCV*, 2025.

- [52] JetBrains Research. Cutting through the noise: Smarter context management for llm-powered agents. <http://blog.jetbrains.com/research/2025/12/efficient-context-management/>, December 2025. Published Dec 2025; Accessed: 2026-01-05.
- [53] Ke Ji, Jiahao Xu, Tian Liang, Qiuzhi Liu, Zhiwei He, Xingyu Chen, Xiaoyuan Liu, Zhijie Wang, Junying Chen, Benyou Wang, et al. The first few tokens are all you need: An efficient and effective unsupervised prefix fine-tuning method for reasoning models. *arXiv preprint arXiv:2503.02875*, 2025.
- [54] Huiqiang Jiang, Qianhui Wu, Chin-Yew Lin, Yuqing Yang, and Lili Qiu. Llmllingua: Compressing prompts for accelerated inference of large language models. *EMNLP*, 2023.
- [55] Jindong Jiang, Xiuyu Li, Zhijian Liu, Muyang Li, Guo Chen, Zhiqi Li, De-An Huang, Guilin Liu, Zhiding Yu, Kurt Keutzer, et al. Token-efficient long video understanding for multimodal llms. *arXiv preprint arXiv:2503.04130*, 2025.
- [56] Pengfei Jiang, Hanjun Li, Linglan Zhao, Fei Chao, Ke Yan, Shouhong Ding, and Rongrong Ji. Visa: Group-wise visual token selection and aggregation via graph summarization for efficient mllms inference. *ACM MM*, 2025.
- [57] Titong Jiang, Xuefeng Jiang, Yuan Ma, Xin Wen, Bailin Li, Kun Zhan, Peng Jia, Yahui Liu, Sheng Sun, and Xianpeng Lang. The better you learn, the smarter you prune: Towards efficient vision-language-action models via differentiable token pruning. *arXiv preprint arXiv:2509.12594*, 2025.
- [58] Chen Ju, Haicheng Wang, Haozhe Cheng, Xu Chen, Zhonghua Zhai, Weilin Huang, Jinsong Lan, Shuai Xiao, and Bo Zheng. Turbo: Informativity-driven acceleration plug-in for vision-language large models. In *ECCV*. Springer, 2024.
- [59] Minki Kang, Wei-Ning Chen, Dongge Han, Huseyin A Inan, Lukas Wutschitz, Yanzhi Chen, Robert Sim, and Saravan Rajmohan. Acon: Optimizing context compression for long-horizon llm agents. *arXiv preprint arXiv:2510.00615*, 2025.
- [60] Minseo Kang and Byunghan Lee. Tictok: Time-series anomaly detection with contrastive tokenization. *IEEE Access*, 2023.
- [61] Oğuzhan Fatih Kar, Alessio Tonioni, Petra Poklukur, Achin Kulshrestha, Amir Zamir, and Federico Tombari. Brave: Broadening the visual encoding of vision-language models. In *ECCV*. Springer, 2024.
- [62] Gyuwan Kim and Kyunghyun Cho. Length-adaptive transformer: Train once with length drop, use anytime with search. *ACL*, 2021.
- [63] Kwonyoung Kim, Jungin Park, Jin Kim, Hyeongjun Kwon, and Kwanghoon Sohn. Faster parameter-efficient tuning with token redundancy reduction. *CVPR*, 2025.
- [64] Minchul Kim, Shangqian Gao, Yen-Chang Hsu, Yilin Shen, and Hongxia Jin. Token fusion: Bridging the gap between token pruning and token merging. In *WACV*, 2024.
- [65] Sehoon Kim, Sheng Shen, David Thorsley, Amir Gholami, Woosuk Kwon, Joseph Hassoun, and Kurt Keutzer. Learned token pruning for transformers. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, 2022.
- [66] Yeachan Kim, Junho Kim, Jun-Hyung Park, Mingyu Lee, and SangKeun Lee. Leap-of-thought: Accelerating transformers via dynamic token routing. In *EMNLP*, 2023.
- [67] Zhenglun Kong, Peiyan Dong, Xiaolong Ma, Xin Meng, Wei Niu, Mengshu Sun, Xuan Shen, Geng Yuan, Bin Ren, Hao Tang, et al. Spvit: Enabling faster vision transformers via latency-aware soft token pruning. In *ECCV*. Springer, 2022.
- [68] Zhenglun Kong, Haoyu Ma, Geng Yuan, Mengshu Sun, Yanyue Xie, Peiyan Dong, Xin Meng, Xuan Shen, Hao Tang, Minghai Qin, et al. Peeling the onion: Hierarchical reduction of data redundancy for efficient vision transformer training. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2023.
- [69] Seon-Ho Lee, Jue Wang, Zhikang Zhang, David Fan, and Xinyu Li. Video token merging for long video understanding. *NeurIPS*, 2024.
- [70] Cheng Lei, Ao Li, Hu Yao, Ce Zhu, and Le Zhang. Rethinking token reduction with parameter-efficient fine-tuning in vit for pixel-level tasks. *CVPR*, 2025.

- [71] Sheng Li, Qitao Tan, Yue Dai, Zhenglun Kong, Tianyu Wang, Jun Liu, Ao Li, Ninghao Liu, Yufei Ding, Xulong Tang, et al. Mutual effort for efficiency: A similarity-based token pruning for vision transformers in self-supervised learning. *ICLR*, 2025.
- [72] Siyuan Li, Luyuan Zhang, Zedong Wang, Juanxi Tian, Cheng Tan, Zicheng Liu, Chang Yu, Qingsong Xie, Haonan Lu, Haoqian Wang, et al. Mergevq: A unified framework for visual generation and representation with disentangled token merging and quantization. *arXiv preprint arXiv:2504.00999*, 2025.
- [73] Wentong Li, Yuqian Yuan, Jian Liu, Dongqi Tang, Song Wang, Jie Qin, Jianke Zhu, and Lei Zhang. Tokenpacker: Efficient visual projector for multimodal llm. *arXiv preprint arXiv:2407.02392*, 2024.
- [74] Xirui Li, Chao Ma, Xiaokang Yang, and Ming-Hsuan Yang. Vidtoe: Video token merging for zero-shot video editing. In *CVPR*, 2024.
- [75] Yanwei Li, Chengyao Wang, and Jiaya Jia. Llama-vid: An image is worth 2 tokens in large language models. In *ECCV*. Springer, 2024.
- [76] Yanwei Li, Yuechen Zhang, Chengyao Wang, Zhisheng Zhong, Yixin Chen, Ruihang Chu, Shaoteng Liu, and Jiaya Jia. Mini-gemini: Mining the potential of multi-modality vision language models. *arXiv preprint arXiv:2403.18814*, 2024.
- [77] Yize Li, Yihua Zhang, Sijia Liu, and Xue Lin. Pruning then reweighting: Towards data-efficient training of diffusion models. In *IEEE ICASSP*, 2025.
- [78] Weicong Liang, Yuhui Yuan, Henghui Ding, Xiao Luo, Weihong Lin, Ding Jia, Zheng Zhang, Chao Zhang, and Han Hu. Expediting large-scale vision transformer for dense prediction without fine-tuning. *NeurIPS*, 2022.
- [79] Youwei Liang, Chongjian Ge, Zhan Tong, Yibing Song, Jue Wang, and Pengtao Xie. Not all patches are what you need: Expediting vision transformers via token reorganizations. *ICLR*, 2022.
- [80] Zhenghao Lin, Zhibin Gou, Yeyun Gong, Xiao Liu, Ruochen Xu, Chen Lin, Yujiu Yang, Jian Jiao, Nan Duan, Weizhu Chen, et al. Not all tokens are what you need for pretraining. *NeurIPS*, 2024.
- [81] Zhihang Lin, Mingbao Lin, Luxi Lin, and Rongrong Ji. Boosting multimodal large language models with visual tokens withdrawal for rapid inference. In *AAAI*, 2025.
- [82] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *NeurIPS*, 2023.
- [83] Xiangrui Liu, Yan Shu, Zheng Liu, Ao Li, Yang Tian, and Bo Zhao. Video-xl-pro: Reconstructive token compression for extremely long video understanding. *arXiv preprint arXiv:2503.18478*, 2025.
- [84] Xin Liu, Jie Liu, Jie Tang, and Gangshan Wu. Catanet: Efficient content-aware token aggregation for lightweight image super-resolution. *CVPR*, 2025.
- [85] Xuyang Liu, Yiyu Wang, Junpeng Ma, and Linfeng Zhang. Video compression commander: Plug-and-play inference acceleration for video large language models. *EMNLP*, 2025.
- [86] Xuyang Liu, Zichen Wen, Shaobo Wang, Junjie Chen, Zhishan Tao, Yubo Wang, Tailai Chen, Xiangqi Jin, Chang Zou, Yiyu Wang, et al. Shifting ai efficiency from model-centric to data-centric compression. *arXiv preprint arXiv:2505.19147*, 2025.
- [87] Yuqi Liu, Pengfei Xiong, Luhui Xu, Shengming Cao, and Qin Jin. Ts2-net: Token shift and selection transformer for text-video retrieval. In *ECCV*. Springer, 2022.
- [88] Zhihang Liu, Chen-Wei Xie, Pandeng Li, Liming Zhao, Longxiang Tang, Yun Zheng, Chuanbin Liu, and Hongtao Xie. Hybrid-level instruction injection for video token compression in multi-modal large language models. *CVPR*, 2025.
- [89] Lei Lu, Yize Li, Yanzhi Wang, Wei Wang, and Wei Jiang. Hdcompression: Hybrid-diffusion image compression for ultra-low bitrates. *PRICAI*, 2025.
- [90] Haotian Luo, Li Shen, Haiying He, Yibo Wang, Shiwei Liu, Wei Li, Naiqiang Tan, Xiaochun Cao, and Dacheng Tao. O1-pruner: Length-harmonizing fine-tuning for o1-like reasoning pruning. *arXiv preprint arXiv:2501.12570*, 2025.
- [91] Haoyu Ma, Zhe Wang, Yifei Chen, Deying Kong, Liangjian Chen, Xingwei Liu, Xiangyi Yan, Hao Tang, and Xiaohui Xie. Ppt: token-pruned pose transformer for monocular and multi-view human pose estimation. In *European Conference on Computer Vision*. Springer, 2022.

- [92] Luca Masserano, Abdul Fatir Ansari, Boran Han, Xiyuan Zhang, Christos Faloutsos, Michael W Mahoney, Andrew Gordon Wilson, Youngsuk Park, Syama Rangapuram, Danielle C Maddix, et al. Enhancing foundation models for time series forecasting via wavelet-based tokenization. *arXiv preprint arXiv:2412.05244*, 2024.
- [93] Jesse Mu, Xiang Li, and Noah Goodman. Learning to compress prompts with gist tokens. *NeurIPS*, 2023.
- [94] Piotr Nawrot, Jan Chorowski, Adrian Łańcucki, and Edoardo M Ponti. Efficient transformers with dynamic token pooling. *ACL*, 2023.
- [95] NVIDIA. Llama-nemotron: Efficient reasoning models. *arXiv preprint arXiv:2505.00949*, 2025.
- [96] OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2024.
- [97] Yue Pan, Minxuan Zhou, Chonghan Lee, Zheyu Li, Rishika Kushwah, Vijaykrishnan Narayanan, and Tajana Rosing. Primate: Processing in memory acceleration for dynamic token-pruning transformers. In *Proceedings of the 29th Asia and South Pacific Design Automation Conference, ASPDAC '24*. IEEE Press, 2024.
- [98] Dhruv Parikh, Shouyi Li, Bingyi Zhang, Rajgopal Kannan, Carl Busart, and Viktor Prasanna. Accelerating vit inference on fpga through static and dynamic pruning. In *2024 IEEE 32nd Annual International Symposium on Field-Programmable Custom Computing Machines (FCCM)*. IEEE, 2024.
- [99] William Peebles and Saining Xie. Scalable diffusion models with transformers. *ICCV*, 2023.
- [100] Xiaohuan Pei, Yuxing Chen, Siyu Xu, Yunke Wang, Yuheng Shi, and Chang Xu. Action-aware dynamic pruning for efficient vision-language-action manipulation. *arXiv preprint arXiv:2509.22093*, 2025.
- [101] Qwen. Qwq-32b: Embracing the power of reinforcement learning, 2025.
- [102] Yongming Rao, Wenliang Zhao, Benlin Liu, Jiwen Lu, Jie Zhou, and Cho-Jui Hsieh. Dynamicvit: Efficient vision transformers with dynamic token sparsification. *NeurIPS*, 2021.
- [103] Shuhuai Ren, Sishuo Chen, Shicheng Li, Xu Sun, and Lu Hou. Testa: Temporal-spatial token aggregation for long-form video-language understanding. In *EMNLP*, 2023.
- [104] Michael Ryoo, AJ Piergiovanni, Anurag Arnab, Mostafa Dehghani, and Anelia Angelova. Tokenlearner: Adaptive space-time tokenization for videos. *NeurIPS*, 34, 2021.
- [105] Yuzhang Shang, Mu Cai, Bingxin Xu, Yong Jae Lee, and Yan Yan. Llava-prumerge: Adaptive token reduction for efficient large multimodal models. *arXiv preprint arXiv:2403.15388*, 2024.
- [106] Kele Shao, Keda Tao, Kejia Zhang, Sicheng Feng, Mu Cai, Yuzhang Shang, Haoxuan You, Can Qin, Yang Sui, and Huan Wang. When tokens talk too much: A survey of multimodal long-context token compression across images, videos, and audios. *arXiv preprint arXiv:2507.20198*, 2025.
- [107] Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Y Wu, et al. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv:2402.03300*, 2024.
- [108] Leqi Shen, Tianxiang Hao, Tao He, Sicheng Zhao, Yifeng Zhang, Pengzhang Liu, Yongjun Bao, and Guiguang Ding. Tempme: Video temporal token merging for efficient text-video retrieval. *ICLR*, 2025.
- [109] Xuan Shen, Peiyan Dong, Lei Lu, Zhenglun Kong, Zhengang Li, Ming Lin, Chao Wu, and Yanzhi Wang. Agile-quant: Activation-guided quantization for faster inference of llms on the edge. In *AAAI*, volume 38, 2024.
- [110] Xuan Shen, Yizhou Wang, Xiangxi Shi, Yanzhi Wang, Pu Zhao, and Jiuxiang Gu. Efficient reasoning with hidden thinking. *arXiv preprint arXiv:2501.19201*, 2025.
- [111] Dingjie Song, Wenjun Wang, Shunian Chen, Xidong Wang, Michael Guan, and Benyou Wang. Less is more: A simple yet effective token reduction method for efficient multi-modal llms. *COLING*, 2025.
- [112] Dimitris Spathis and Fahim Kawsar. The first step is the hardest: Pitfalls of representing and tokenizing temporal data for large language models. *Journal of the American Medical Informatics Association*, 2024.
- [113] DiJia Su, Sainbayar Sukhbaatar, Michael Rabbat, Yuandong Tian, and Qinqing Zheng. Dualformer: Controllable fast and slow thinking by learning with randomized reasoning traces. *ICLR*, 2025.

- [114] Xiaorui Su, Shvat Messica, Yepeng Huang, Ruth Johnson, Lukas Fesser, Shanghua Gao, Faryad Sahnesh, and Marinka Zitnik. Multimodal medical code tokenizer. *arXiv preprint arXiv:2502.04397*, 2025.
- [115] Yi Su, Yuechi Zhou, Quantong Qiu, Juntao Li, Qingrong Xia, Ping Li, Xinyu Duan, Zhefeng Wang, and Min Zhang. Accurate kv cache quantization with outlier tokens tracing. *ACL*, 2025.
- [116] Zhaochen Su, Peng Xia, Hangyu Guo, Zhenhua Liu, Yan Ma, Xiaoye Qu, Jiaqi Liu, Yanshu Li, Kaide Zeng, Zhengyuan Yang, et al. Thinking with images for multimodal reasoning: Foundations, methods, and future frontiers. *arXiv preprint arXiv:2506.23918*, 2025.
- [117] Yang Sui, Yu-Neng Chuang, Guanchu Wang, Jiamu Zhang, Tianyi Zhang, Jiayi Yuan, Hongyi Liu, Andrew Wen, Shaochen Zhong, Hanjie Chen, et al. Stop overthinking: A survey on efficient reasoning for large language models. *arXiv preprint arXiv:2503.16419*, 2025.
- [118] Boyuan Sun, Jiaxing Zhao, Xihan Wei, and Qibin Hou. Llava-scissor: Token compression with semantic connected components for video llms, 2025.
- [119] Burak Suyunu, Özdeniz Dolu, and Arzucan Özgür. evobpe: Evolutionary protein sequence tokenization. *arXiv preprint arXiv:2503.08838*, 2025.
- [120] Sabera Talukder, Yisong Yue, and Georgia Gkioxari. Totem: Tokenized time series embeddings for general time series analysis. *arXiv preprint arXiv:2402.16412*, 2024.
- [121] Hao Tang, Chenwei Xie, Haiyang Wang, Xiaoyi Bao, Tingyu Weng, Pandeng Li, Yun Zheng, and Liwei Wang. Ufo: A unified approach to fine-grained visual perception via open-ended language interface. *NeurIPS*, 2025.
- [122] Yehui Tang, Kai Han, Yunhe Wang, Chang Xu, Jianyuan Guo, Chao Xu, and Dacheng Tao. Patch slimming for efficient vision transformers. In *CVPR*, 2022.
- [123] Yao Tao, Yehui Tang, Yun Wang, Mingjian Zhu, Hailin Hu, and Yunhe Wang. Saliency-driven dynamic token pruning for large language models. *arXiv preprint arXiv:2504.04514*, 2025.
- [124] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *NeurIPS*, 30, 2017.
- [125] Hanrui Wang, Zhekai Zhang, and Song Han. Spatten: Efficient sparse attention architecture with cascade token and head pruning. In *2021 IEEE International Symposium on High-Performance Computer Architecture (HPCA)*. IEEE, 2021.
- [126] Hongjie Wang, Difan Liu, Yan Kang, Yijun Li, Zhe Lin, Niraj K Jha, and Yuchen Liu. Attention-driven training-free efficiency enhancement of diffusion models. In *CVPR*, 2024.
- [127] Junke Wang, Xitong Yang, Hengduo Li, Li Liu, Zuxuan Wu, and Yu-Gang Jiang. Efficient video transformers with spatial-temporal token selection. In *ECCV*. Springer, 2022.
- [128] Renxi Wang, Xudong Han, Lei Ji, Shu Wang, Timothy Baldwin, and Haonan Li. Toolgen: Unified tool retrieval and calling via generation. *arXiv preprint arXiv:2410.03439*, 2024.
- [129] Rui Wang, Hongru Wang, Boyang Xue, Jianhui Pang, Shudong Liu, Yi Chen, Jiahao Qiu, Derek Fai Wong, Heng Ji, and Kam-Fai Wong. Harnessing the reasoning economy: A survey of efficient reasoning for large language models. *arXiv preprint arXiv:2503.24377*, 2025.
- [130] Shenzhi Wang, Le Yu, Chang Gao, Chujie Zheng, Shixuan Liu, Rui Lu, Kai Dang, Xionghui Chen, Jianxin Yang, Zhenru Zhang, et al. Beyond the 80/20 rule: High-entropy minority tokens drive effective reinforcement learning for llm reasoning. *arXiv preprint arXiv:2506.01939*, 2025.
- [131] Song Wang, Gongfan Fang, Lingdong Kong, Xiangtai Li, Jianyun Xu, Sheng Yang, Qiang Li, Jianke Zhu, and Xinchao Wang. Pixelthink: Towards efficient chain-of-pixel reasoning. *arXiv preprint arXiv:2505.23727*, 2025.
- [132] Yixiao Wang, Yifei Zhang, Mingxiao Huo, Ran Tian, Xiang Zhang, Yichen Xie, Chenfeng Xu, Pengliang Ji, Wei Zhan, Mingyu Ding, et al. Sparse diffusion policy: A sparse, reusable, and flexible policy for robot learning. *arXiv preprint arXiv:2407.01531*, 2024.
- [133] Siyuan Wei, Tianzhu Ye, Shen Zhang, Yao Tang, and Jiajun Liang. Joint token pruning and squeezing towards more aggressive compression of vision transformers. In *CVPR*, 2023.
- [134] Zichen Wen, Yifeng Gao, Weijia Li, Conghui He, and Linfeng Zhang. Token pruning in multimodal large language models: Are we solving the right problem? *arXiv preprint arXiv:2502.11501*, 2025.

- [135] David Wingate, Mohammad Shoeybi, and Taylor Sorensen. Prompt compression and contrastive conditioning for controllability and toxicity reduction in language models. *EMNLP Findings*, 2022.
- [136] Penghao Wu, Lewei Lu, and Ziwei Liu. Streamline without sacrifice – squeeze out computation redundancy in lmm. *ICML*, 2025.
- [137] Qiong Wu, Wenhao Lin, Weihao Ye, Yiyi Zhou, Xiaoshuai Sun, and Rongrong Ji. Accelerating multi-modal large language models via dynamic visual-token exit and the empirical findings. *arXiv preprint arXiv:2411.19628*, 2024.
- [138] Shengqiong Wu, Hao Fei, Xiangtai Li, Jiayi Ji, Hanwang Zhang, Tat-Seng Chua, and Shuicheng Yan. Towards semantic equivalence of tokenization in multimodal llm. *ICLR*, 2025.
- [139] Siye Wu, Jian Xie, Yikai Zhang, Aili Chen, Kai Zhang, Yu Su, and Yanghua Xiao. Arm: Adaptive reasoning model. *NeurIPS*, 2025.
- [140] Tong Wu, Junzhe Shen, Zixia Jia, Yuxuan Wang, and Zilong Zheng. From hours to minutes: Lossless acceleration of ultra long sequence generation up to 100k tokens. *arXiv preprint arXiv:2502.18890*, 2025.
- [141] Wei Wu, Zhuoshi Pan, Chao Wang, Liyi Chen, Yunchu Bai, Tianfu Wang, Kun Fu, Zheng Wang, and Hui Xiong. Tokenselect: Efficient long-context inference and length extrapolation for llms via dynamic token-level kv cache selection. *EMNLP*, 2025.
- [142] Xinjian Wu, Fanhu Zeng, Xiudong Wang, and Xinghao Chen. Ppt: Token pruning and pooling for efficient vision transformers. *arXiv preprint arXiv:2310.01812*, 2023.
- [143] Heming Xia, Yongqi Li, Chak Tou Leong, Wenjie Wang, and Wenjie Li. Tokenskip: Controllable chain-of-thought compression in llms. *EMNLP*, 2025.
- [144] Wenyi Xiao, Leilei Gan, Weilong Dai, Wanggui He, Ziwei Huang, Haoyuan Li, Fangxun Shu, Zhelun Yu, Peng Zhang, Hao Jiang, et al. Fast-slow thinking for large vision-language model reasoning. *arXiv preprint arXiv:2504.18458*, 2025.
- [145] Yuan-An Xiao, Pengfei Gao, Chao Peng, and Yingfei Xiong. Improving the efficiency of llm agent systems through trajectory reduction. *arXiv preprint arXiv:2509.23586*, 2025.
- [146] Ling Xing, Alex Jinpeng Wang, Rui Yan, Xiangbo Shu, and Jinhui Tang. Vision-centric token compression in large language model. *NeurIPS*, 2025.
- [147] Keqiang Yan, Xiner Li, Hongyi Ling, Kenna Ashen, Carl Edwards, Raymundo Arróyave, Marinka Zitnik, Heng Ji, Xiaofeng Qian, Xiaoning Qian, et al. Invariant tokenization of crystalline materials for language model enabled generation. *NeurIPS*, 2024.
- [148] Wilson Yan, Volodymyr Mnih, Aleksandra Faust, Matei Zaharia, Pieter Abbeel, and Hao Liu. Elastictok: Adaptive tokenization for image and video. *ICLR*, 2025.
- [149] Cheng Yang, Yang Sui, Jinqi Xiao, Lingyi Huang, Yu Gong, Chendi Li, Jinghua Yan, Yu Bai, Ponnuswamy Sadayappan, Xia Hu, and Bo Yuan. Topv: Compatible token pruning with inference time optimization for fast and low-memory multimodal vision language model. *CVPR*, 2025.
- [150] Senqiao Yang, Junyi Li, Xin Lai, Bei Yu, Hengshuang Zhao, and Jiaya Jia. Visionthink: Smart and efficient vision language model via reinforcement learning. *arXiv preprint arXiv:2507.13348*, 2025.
- [151] Yantai Yang, Yuhao Wang, Zichen Wen, Luo Zhongwei, Chang Zou, Zhipeng Zhang, Chuan Wen, and Linfeng Zhang. Efficientvla: Training-free acceleration and compression for vision-language-action models. *NeurIPS*, 2025.
- [152] Linli Yao, Yicheng Li, Yuancheng Wei, Lei Li, Shuhuai Ren, Yuanxin Liu, Kun Ouyang, Lean Wang, Shicheng Li, Sida Li, et al. Timechat-online: 80% visual tokens are naturally redundant in streaming videos. *ACM MM*, 2025.
- [153] Deming Ye, Yankai Lin, Yufei Huang, and Maosong Sun. Tr-bert: Dynamic token reduction for accelerating bert inference. *NAACL*, 2021.
- [154] Xubing Ye, Yukang Gan, Xiaoke Huang, Yixiao Ge, and Yansong Tang. Voco-llama: Towards vision compression with large language models. *CVPR*, 2025.
- [155] Xinyu Yuan, Zichen Wang, Marcus Collins, and Huzefa Rangwala. Protein structure tokenization: Benchmarking and new recipe. *arXiv preprint arXiv:2503.00089*, 2025.

- [156] Fanhu Zeng and Deli Yu. M2m-tag: Training-free many-to-many token aggregation for vision transformer acceleration. In *Workshop on Machine Learning and Compression, NeurIPS*, 2024.
- [157] Fanhu Zeng, Deli Yu, Zhenglun Kong, and Hao Tang. Token transforming: A unified and training-free token compression framework for vision transformer acceleration. *arXiv preprint arXiv:2506.05709*, 2025.
- [158] Yuting Zeng, Weizhe Huang, Lei Jiang, Tongxuan Liu, Xitai Jin, Chen Tianying Tiana, Jing Li, and Xiaohua Xu. S²-mad: Breaking the token barrier to enhance multi-agent debate efficiency. *NAACL*, 2025.
- [159] Zheng Zhan, Zhenglun Kong, Yifan Gong, Yushu Wu, Zichong Meng, Hangyu Zheng, et al. Exploring token pruning in vision state space models. In *NeurIPS*, 2024.
- [160] Zheng Zhan, Yushu Wu, Zhenglun Kong, Changdi Yang, Yifan Gong, Xuan Shen, Xue Lin, Pu Zhao, and Yanzhi Wang. Rethinking token reduction for state space models. *EMNLP*, 2024.
- [161] Diankun Zhang, Guoan Wang, Runwen Zhu, Jianbo Zhao, Xiwu Chen, Siyu Zhang, Jiahao Gong, Qibin Zhou, Wenyan Zhang, Ningzi Wang, et al. Sparsead: Sparse query-centric paradigm for efficient end-to-end autonomous driving. *arXiv preprint arXiv:2404.06892*, 2024.
- [162] Evelyn Zhang, Jiayi Tang, Xuefei Ning, and Linfeng Zhang. Training-free and hardware-friendly acceleration for diffusion models via similarity-based token pruning. *AAAI*, 2025.
- [163] Guibin Zhang, Yanwei Yue, Zhixun Li, Sukwon Yun, Guancheng Wan, Kun Wang, Dawei Cheng, Jeffrey Xu Yu, and Tianlong Chen. Cut the crap: An economical communication pipeline for llm-based multi-agent systems. *arXiv preprint arXiv:2410.02506*, 2024.
- [164] Haichao Zhang and Yun Fu. Vqtokn: Neural discrete token representation learning for extreme token reduction in video large language models. *NeurIPS*, 2025.
- [165] Jintian Zhang, Yuqi Zhu, Mengshu Sun, Yujie Luo, Shuofei Qiao, Lun Du, Da Zheng, Huajun Chen, and Ningyu Zhang. Lightthinker: Thinking step-by-step compression. *EMNLP*, 2025.
- [166] Qizhe Zhang, Aosong Cheng, Ming Lu, Renrui Zhang, Zhiyong Zhuo, Jiajun Cao, Shaobo Guo, Qi She, and Shanghang Zhang. Beyond text-visual attention: Exploiting visual cues for effective token pruning in vlms. *arXiv preprint arXiv:2412.01818*, 2025.
- [167] Qizhe Zhang, Mengzhen Liu, Lichen Li, Ming Lu, Yuan Zhang, Junwen Pan, Qi She, and Shanghang Zhang. Beyond attention or similarity: Maximizing conditional diversity for token pruning in mllms. *NeurIPS*, 2025.
- [168] Renshan Zhang, Rui Shao, Gongwei Chen, Miao Zhang, Kaiwen Zhou, Weili Guan, and Liqiang Nie. Falcon: Resolving visual redundancy and fragmentation in high-resolution multimodal large language models via visual registers. *ICCV*, 2025.
- [169] Shuheng Zhang, Yuqi Liu, Hongbo Zhou, Jun Peng, Yiyi Zhou, Xiaoshuai Sun, and Rongrong Ji. Adaflow: Efficient long video editing via adaptive attention slimming and keyframe selection. *arXiv preprint arXiv:2502.05433*, 2025.
- [170] Yuan Zhang, Chun-Kai Fan, Junpeng Ma, Wenzhao Zheng, Tao Huang, Kuan Cheng, Denis Gudovskiy, Tomoyuki Okuno, Yohei Nakata, Kurt Keutzer, et al. Sparsevlm: Visual token sparsification for efficient vision-language model inference. In *ICML*, 2025.
- [171] Wangbo Zhao, Yizeng Han, Jiasheng Tang, Kai Wang, Yibing Song, Gao Huang, Fan Wang, and Yang You. Dynamic diffusion transformer. *ICLR*, 2025.
- [172] Xianwei Zhuang, Zhihong Zhu, Yuxin Xie, Liming Liang, and Yuexian Zou. Vaspars: Towards efficient visual hallucination mitigation via visual-aware token sparsification. *CVPR*, 2025.
- [173] Chang Zou, Xuyang Liu, Ting Liu, Siteng Huang, and Linfeng Zhang. Accelerating diffusion transformers with token-wise feature caching. In *ICLR*, 2025.
- [174] Jiaru Zou, Xiyuan Yang, Ruizhong Qiu, Gaotang Li, Katherine Tieu, Pan Lu, Ke Shen, Hanghang Tong, Yejin Choi, Jingrui He, et al. Latent collaboration in multi-agent systems. *arXiv preprint arXiv:2511.20639*, 2025.

A Theoretical Formulation

In this section, we provide a unified mathematical framework for token reduction methods. We formalize the token reduction process into two distinct phases: *Compression Criteria* (how to evaluate tokens) and *Compression Strategies* (how to reduce tokens).

A.1 Problem Definition

Given an input sequence $X = [x_1, x_2, \dots, x_N] \in \mathbb{R}^{N \times D}$, where N represents the sequence length and D denotes the feature dimension, the objective of token reduction is to generate a compressed sequence $X' = [x'_1, x'_2, \dots, x'_M] \in \mathbb{R}^{M \times D}$, such that $M \ll N$.

The reduction process can be formally defined as a composite function of a scoring criterion \mathcal{E} and a compression strategy \mathcal{P} :

$$X' = \mathcal{P}(X, \mathcal{E}(X)) \quad (1)$$

where $\mathcal{E}(X)$ outputs importance scores, clustering assignments, or gradient sensitivities, and \mathcal{P} executes the dimensionality reduction. Figure 2 shows the token reduction pipeline.

A.2 Compression Criteria

The scoring function $\mathcal{E} : X \rightarrow \mathcal{S}$ determines the semantic value or redundancy of each token. We broaden the categorization to include gradient and entropy-based metrics alongside standard parametric/non-parametric approaches.

Attention-based Scoring. Utilizing the inherent sparsity of the self-attention mechanism, the importance of a token x_i is quantified by the attention it receives. This can be *global* (averaged across all heads/tokens) or *targeted* (attention from the special [CLS] token or specific query tokens). The score s_i is typically calculated as:

$$s_i = \sum_{j \in \mathcal{Q}} \text{Attn}(x_j, x_i) \quad (2)$$

where \mathcal{Q} is the set of query tokens (e.g., $\mathcal{Q} = \{x_{\text{CLS}}\}$ for classification tasks or $\mathcal{Q} = \{x_{1 \dots N}\}$ for global density).

Similarity-based Scoring. This approach assumes that tokens close in the feature space contain redundant information. The criterion calculates pairwise distances to identify clusters or redundant pairs. For tokens (x_i, x_j) , the metric is typically cosine similarity:

$$\text{Sim}(x_i, x_j) = \frac{x_i \cdot x_j}{\|x_i\|_2 \|x_j\|_2} \quad (3)$$

High similarity scores ($\text{Sim} > \tau$) trigger merging operations. Advanced methods extend this to density-based clustering (e.g., K-Means or DPC-KNN) to identify representative centroids.

Gradient and Entropy-based Scoring. Beyond static feature analysis, recent methods employ dynamic metrics. *Gradient-based* criteria measure a token’s contribution to the loss function, retaining tokens with high gradient norms ($\|\nabla_{x_i} \mathcal{L}\|$). *Entropy-based* criteria evaluate the uncertainty of the model’s prediction; tokens with low information density (low entropy) are candidates for pruning in early-exit or fast-forwarding frameworks.

Parametric Scoring. Parametric methods introduce a lightweight auxiliary module (e.g., a predictor network \mathcal{M}_ϕ) to explicitly predict token utility:

$$S = \mathcal{M}_\phi(X) \quad (4)$$

where $S \in [0, 1]^N$ represents the keep-probability or saliency score. These predictors are trained via Gumbel-Softmax or reinforcement learning (RL) to maximize downstream accuracy while minimizing token count.

A.3 Compression Strategies

Once the relationships or scores are established, the compression strategy \mathcal{P} transforms the sequence. We classify these into four primary mechanisms.

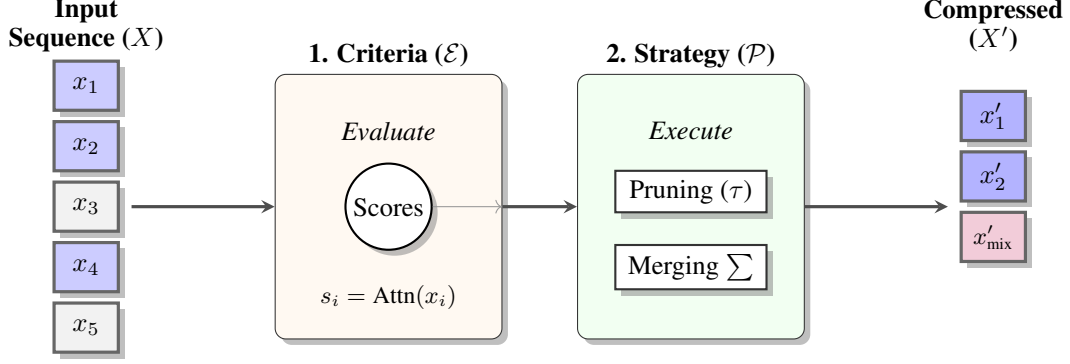


Figure 2: **The token reduction pipeline.** We formulate reduction as a composite of Criteria \mathcal{E} (scoring) and Strategy \mathcal{P} (pruning/merging).

Token Pruning (Hard & Soft). Pruning is a selection process that discards tokens based on the criteria \mathcal{E} .

$$X' = \{x_i \mid s_i \in \text{TopK}(S) \vee s_i > \tau\} \quad (5)$$

While *Hard Pruning* permanently removes tokens, *Soft Pruning* (or "packaging") aggregates the discarded tokens into a single summary token to preserve residual information, preventing total information loss.

Token Merging & Clustering. Merging aggregates information from a set of tokens $\mathcal{C} = \{x_1, \dots, x_k\}$ identified as similar. This ranges from bipartite matching (merging pairs) to density-based clustering (merging large groups). The merged token x'_{cluster} is a weighted average:

$$x'_{\text{cluster}} = \frac{\sum_{x_j \in \mathcal{C}} w_j x_j}{\sum_{x_j \in \mathcal{C}} w_j} \quad (6)$$

where w_j tracks the token's size (number of constituent patches), ensuring proportional representation of fused features.

Transformation-based Compression. These methods reduce sequence length through structural operations, exploiting spatial (image) or temporal (video) priors. Common techniques include:

- **Pooling/Unshuffle:** Non-parametric downsampling: $\text{Pool} : \mathbb{R}^{H \times W} \rightarrow \mathbb{R}^{\frac{H}{r} \times \frac{W}{r}}$.
- **Convolution:** Strided convolution to abstract local neighborhoods: $X' = \text{Conv}_{k \times k}(X, s)$.

Token Distillation (Query-based). Distillation employs a set of learnable latent queries $Q \in \mathbb{R}^{M \times D}$ to extract information from the input X via cross-attention mechanisms (e.g., Perceiver Resampler or Q-Former). This decouples output length M from input length N :

$$X' = \text{Softmax} \left(\frac{Q(XW_K)^T}{\sqrt{D}} \right) (XW_V) \quad (7)$$

This strategy is particularly effective for cross-modal alignment, compressing dense visual features into sparse text-aligned tokens.

We summarize the general workflow of training-free token reduction in Algorithm 1. We also show a visualization of token reduction in Figure 3.

A.4 Controllable Reasoning via Reinforcement Learning

Controllable Reasoning refers to the ability of language models to dynamically adjust the depth and length of their reasoning processes according to user-specified constraints (e.g., exact or maximum token lengths), thereby enabling a tunable trade-off between reasoning efficiency and accuracy. State-of-the-art approaches [90, 2] typically employ reinforcement learning frameworks [107], where a multi-objective reward function is designed to optimize correctness and length compliance jointly. Specifically, the reward function incorporates two key objectives:

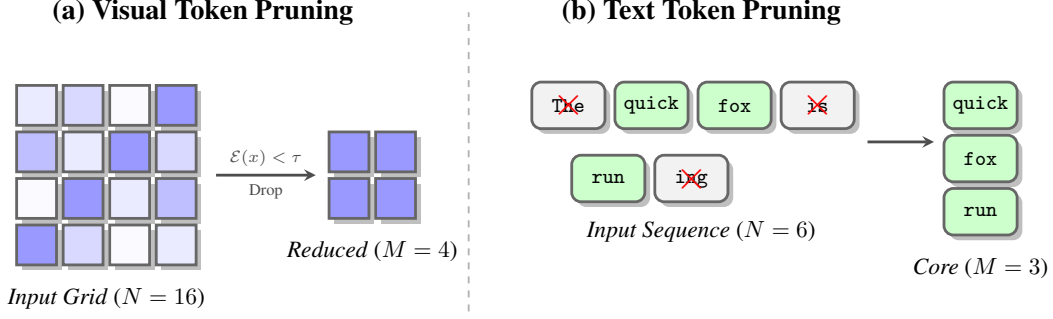


Figure 3: **Visualization of token reduction.** (a) *Image*: Visual tokens are pruned based on saliency, retaining only the most salient patches. (b) *Text*: Low-information stop words (gray) are removed to form a compressed semantic core.

Algorithm 1 General Training-Free Token Reduction Workflow

Require: Input tokens $X \in \mathbb{R}^{N \times D}$, Target ratio r

Ensure: Compressed tokens $X' \in \mathbb{R}^{M \times D}$

```

1:  $M \leftarrow \lfloor N \times (1 - r) \rfloor$ 
2: Phase 1: Criteria Calculation ( $\mathcal{E}$ )
3: if Attention-based then
4:    $S \leftarrow \text{Agg}(\text{AttentionMap}(X))$  ▷ Agg: Sum/Avg over heads
5: else if Similarity-based then
6:    $A \leftarrow XX^T / (\|X\| \|X\|)$  ▷ Cosine Similarity Matrix
7:   Partition  $X$  into sets  $\{C_1, \dots, C_M\}$  via clustering on  $A$ 
8: end if
9: Phase 2: Strategy Execution ( $\mathcal{P}$ )
10: if Pruning then
11:   Indices  $\leftarrow \text{TopK}(S, M)$ 
12:    $X' \leftarrow \text{Gather}(X, \text{Indices})$ 
13: else if Merging then
14:   for  $m = 1$  to  $M$  do
15:      $x'_m \leftarrow \text{WeightedSum}(C_m)$ 
16:   end for
17: end if
18: return  $X'$ 

```

1. **Correctness reward:** awarded if the model’s output matches the ground-truth answer;
2. **Length penalty:** imposed if the generated sequence length deviates from the target length.

Exact Length Control. It requires the model to generate reasoning sequences whose length exactly matches a user-specified target length:

$$R_{\text{exact}}(y, y_t, n_t) = \mathbb{I}(y = y_t) - \alpha \cdot |n_t - n_y|, \quad (8)$$

where y is the generated sequence, y_t is the ground truth answer, n_t is the target token length, n_y is the actual token length of the generated sequence, $\mathbb{I}(\cdot)$ is the indicator function (1 if correct, 0 otherwise), and α is a penalty weight balancing correctness and length.

Maximum Length Control. It controls the model to generate reasoning sequences no longer than a specified upper limit, encouraging efficient reasoning within a token budget:

$$R_{\text{max}}(y, y_t, n_t) = \mathbb{I}(y = y_t) \cdot \text{clip}(\alpha \cdot (n_t - n_y) + \delta, 0, 1), \quad (9)$$

where $\text{clip}(\cdot, 0, 1)$ clamps reward to the range $[0, 1]$ and δ is an offset term to avoid zero reward (typically set to 0.5).

Length Efficiency Optimization. It encourages the model to shorten reasoning length while maintaining correctness, particularly useful for reducing redundancy in long-reasoning models:

$$R_{\text{eff}}(y, y_t, x) = \frac{\bar{L}_{\text{ref}}(x)}{L(y)} - 1 + \lambda (A(y, y_t) - \bar{A}_{\text{ref}}(x)), \quad (10)$$

where $\bar{L}_{\text{ref}}(x)$ is the average length of reference model outputs for problem x , $A(y, y_t)$ is the accuracy function (1 if correct, 0 otherwise), $\bar{A}_{\text{ref}}(x)$ is the average accuracy of the reference model on problem x , and λ is an accuracy penalty weight to prevent performance degradation from over-compression.

By adjusting hyperparameters (e.g., α and λ) in the reward function, one can flexibly control the model’s tendency to prioritize correctness versus length, which not only achieves precise length control but also maintains or even improves model performance while significantly reducing reasoning overhead.