

On the Limits of Token Reduction for Efficient Unified Vision Language Training

Anonymous CVPR submission

Paper ID ****

Abstract

001 *Unified vision-language models (VLMs) integrate visual*
002 *understanding and visual generation within a single au-*
003 *to-regressive backbone, but their joint training is computa-*
004 *tionally expensive and largely overlooked from an efficiency*
005 *perspective. In this work, we study the feasibility and lim-*
006 *its of token-reduction-based acceleration for unified VLM*
007 *training. Through a systematic analysis of layerwise atten-*
008 *tion allocation, we uncover a fundamental asymmetry: vi-*
009 *sual understanding exhibits substantial late-layer visual re-*
010 *dundancy, whereas visual generation maintains persistent*
011 *dependence on image tokens across depth. Guided by this*
012 *observation, we design task-specific accelerators that selec-*
013 *tively reduce image-token computation for each objective.*
014 *While these methods achieve significant efficiency gains in*
015 *isolated settings, we observe a consistent synergy loss un-*
016 *der unified training—task-specific token dropping necessi-*
017 *tates divergent parameter pathways and eliminates the mu-*
018 *tual performance gains typically observed in joint optimiza-*
019 *tion. Our findings suggest that efficient unified modeling re-*
020 *quires preserving shared cross-task structures, highlighting*
021 *the need for synergy-aware acceleration strategies.*

022 1. Introduction

023 Unified Vision-Language Models (VLMs) [11, 25, 26, 36,
024 38, 39] integrate visual generation [6, 7, 29, 33] and un-
025 derstanding [4, 21, 22, 27] within a single model and have
026 demonstrated remarkable scalability and cross-task poten-
027 tial [32, 37, 39]. However, the training of these models
028 is prohibitively expensive; for instance, VILA-U [39] re-
029 quires approximately 20K A100 GPU hours. While many
030 prior methods propose to reduce inference-time computa-
031 tion in understanding-only VLMs via token pruning or spe-
032 cial attention masks [1, 3, 12, 24, 28, 30, 44], these strate-
033 gies do not directly translate to improve training-time ef-
034 ficiency. Furthermore, existing acceleration techniques for
035 visual understanding do not account for the distinct struc-
036 tural requirements of visual generation, nor do they study
037 the complexities inherent in unifying generative and dis-
038 criminative objectives within a single VLM.

In this paper, we investigate the feasibility and limits of
accelerating the training of unified vision language models.
We adopt the pure autoregressive framework as our testbed,
as it represents one of the most prevalent architectures for
integrating multimodal capabilities [9, 14, 23, 36, 39, 41–
43]. Through an analysis of the attention dynamics within
this framework (in Figure 2), we reveal a critical asymme-
try in task-specific redundancy: while visual understanding
tasks exhibit high token redundancy in the deeper layers,
visual generation depends heavily on the context of pre-
viously generated image tokens within many deep layers.
Building on these insights, we develop task-specific strate-
gies to accelerate training by selectively dropping image to-
kens tailored to the unique requirements of each objective.

Furthermore, we reveal a critical “synergy loss” phe-
nomenon that occurs when task-specific token reduction
methods are applied to the joint training of unified mod-
els. We find that task-specific token dropping disrupts the
inherent synergy between understanding and generation by:
(1) necessitating divergent sets of image-related model pa-
rameters, and (2) eliminating the mutual performance gains
typically observed when both tasks are trained concurrently.
Our diagnostic analysis suggests that aggressive token drop-
ping amplifies task conflicts, offering a cautionary lesson
and a new perspective for future research in efficient unified
modeling. Our contributions are summarized as follows:

- **Unified Redundancy Analysis:** We characterize task-specific attention patterns in unified VLMs, identifying distinct redundancy zones.
- **Task-Specific Accelerators:** We design and implement training-time acceleration for isolated tasks.
- **Discovery of Synergy Loss:** We discover that task-specific optimization strategies fail in unified settings, revealing that forced token reduction disrupts mutual improvements of discriminative and generative objectives.
- **Lessons for Unified Acceleration:** Our results suggest that effective acceleration methods may benefit from preserving shared cross-task structures and carefully accounting for impact on cross-task learning dynamics.

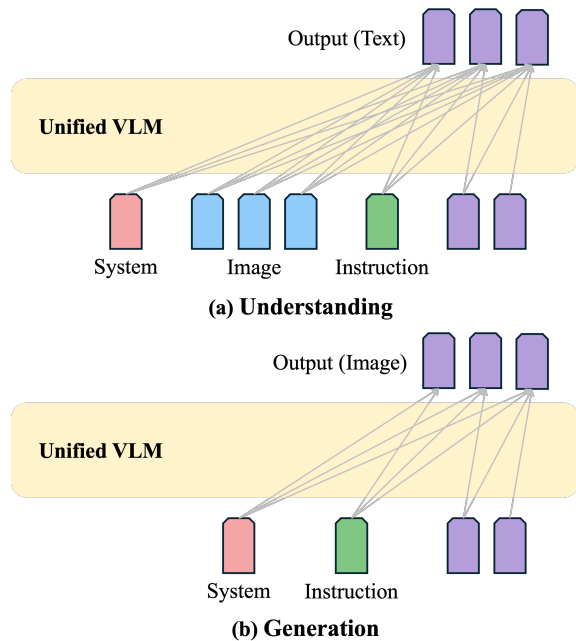


Figure 1. **Unified autoregressive VLM.** A single Transformer backbone processes multimodal sequences under a unified next-token prediction objective. (a) In visual understanding, the model predicts text tokens conditioned on image and textual context. (b) In visual generation, the model autoregressively predicts image tokens conditioned on preceding text and image tokens.

078

2. Related Works

079

080

081

082

083

084

085

086

087

088

089

090

091

092

093

094

095

096

097

098

099

100

101

102

103

104

Unified Vision-Language Models. Recent advancements have shifted toward unifying perception and generation within a single framework. Models like VILA-U [39], Janus [37], and Chameleon [32] utilize discrete visual tokenizers (e.g., VQVAE [35]) to treat images as a “foreign language.” While these models simplify the pipeline by using a single next-token prediction objective, their joint training is computationally demanding. Many other hybrid models that append diffusion heads [29] to a transformer also require fine-tuning the entire backbone across multiple modalities, creating the need for efficient training.

Efficiency in Vision-Language Models. Efficiency research in VLMs has primarily focused on visual understanding during inference-time [3, 12, 17, 20, 24, 30, 45]. For instance, LLaVA-PruMerge [30] and LLaMA-VID [19] reduce the number of visual tokens by identifying spatial redundancy and merging tokens. Other works explore efficient attention mechanisms or special masks to skip redundant computations during inference [44, 46]. However, these methods are often designed for “understanding-only” tasks where the model’s output is limited to text, and a complete set of image tokens is treated as input, thus having difficulty applying to visual generation, and how to reduce training-time computation remains a challenging problem.

Token Reduction and Attention Redundancy. The concept of “token reduction” or “pruning” originates from the

Vision Transformer (ViT) and NLP literature to handle long-sequence data [1, 10, 28, 40]. These methods typically use attention weights or activation statistics as proxies for token importance. In the multimodal domain, recent studies have analyzed attention sinks [40] and sparsity to prune background patches. While effective for single-task models, these importance metrics are not directly transferable to unified models where tokens must serve dual roles in discriminative perception and generative synthesis.

Multi-task Synergy in VLMs. The relationship between understanding and generation has been a subject of ongoing debate. While some studies suggest that generative pre-training provides a stronger world model for perception [36, 41], others have noted the difficulty of balancing these disparate objectives during joint optimization [37]. We build upon this line of inquiry by investigating how structural constraints—specifically, token dropping—affect the stability and synergy of this multi-task learning process.

3. Problem Setup

3.1. Unified Autoregressive Vision-Language Model

We study a unified vision-language model (VLM) that jointly performs visual understanding and visual generation within a single autoregressive Transformer backbone, following the unified next-token prediction paradigm of VILA-U (7B) [39]. Let $x = (x_1, \dots, x_T)$ denote text tokens and $v = (v_1, \dots, v_M)$ denote discrete image tokens obtained from a visual tokenizer (e.g., VQ-based). We construct a multimodal sequence

$$z = (z_1, \dots, z_{|z|}) = (\text{system}, x, v). \quad 134$$

The model, parameterized by θ , is trained using autoregressive next-token prediction:

$$P_\theta(z) = \prod_{t=1}^{|z|} P_\theta(z_t | z_{<t}). \quad 137$$

Under this formulation, both text and image tokens are treated uniformly as discrete tokens in a single sequence, and a shared next-token objective is applied across modalities. We visualize the generation process in Figure 1.

3.2. Training Objective

Unified training mixes data from visual understanding and visual generation tasks under a single objective. For a multimodal training sample z , the loss is defined as the negative log-likelihood:

$$\mathcal{L}_{\text{unified}} = - \sum_{t=1}^{|z|} \log P_\theta(z_t | z_{<t}). \quad (1) \quad 147$$

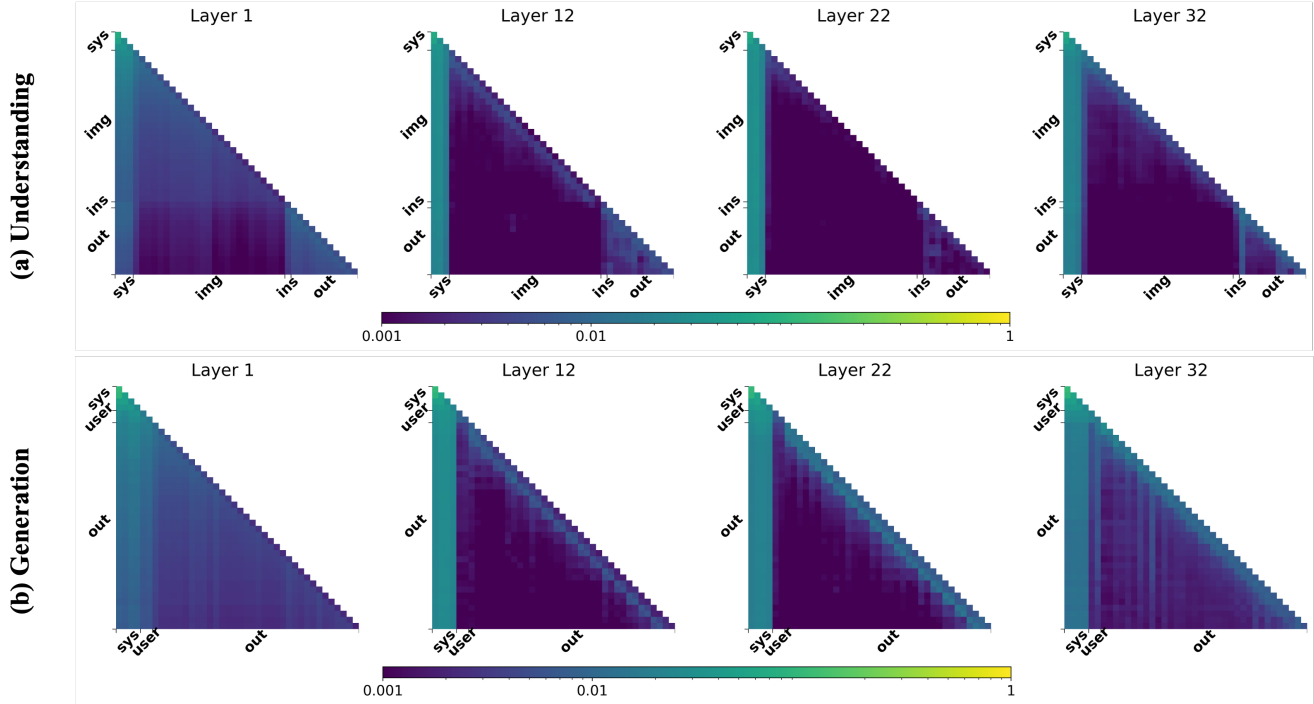


Figure 2. **Asymmetric depth-wise attention patterns in unified VLMs.** Visualization of self-attention heatmaps across layers for understanding (a) and generation (b). Understanding exhibits strong early cross-modal interactions followed by a sharp decay in image-token attention. Generation, however, preserves substantial image-token attention throughout depth, highlighting a fundamental asymmetry in token utilization.

148 This unified objective simultaneously optimizes: **Visual**
 149 **Understanding:** predicting text tokens conditioned on image
 150 and textual context; **Visual Generation:** predicting image
 151 tokens conditioned on preceding text and image tokens.

152 3.3. Computational Bottleneck

153 Let $N = |z|$ denote the total sequence length. A stan-
 154 dard Transformer layer incurs quadratic self-attention cost:
 155 FLOPs per layer $\propto N^2$. Since $N = T + M$, where T and
 156 M are the numbers of text and image tokens respectively,

$$157 (T + M)^2 = T^2 + 2TM + M^2.$$

158 In unified VLM training, image tokens typically dominate
 159 the sequence ($M \gg T$), making the M^2 term the primary
 160 computational bottleneck.

161 This naturally motivates token reduction strategies that
 162 limit the effective participation of image tokens in atten-
 163 tion. In this work, we investigate the effectiveness and lim-
 164 itations of *token-reduction-based training acceleration* for
 165 visual understanding, visual generation, and their unifica-
 166 tion. We begin by analyzing task-specific redundancy pat-
 167 terns through attention statistics.

168 4. Redundancy Analysis

169 To guide the design of our acceleration strategies, we ana-
 170 lyze the layerwise attention behavior of a pre-trained uni-
 171 fied VLM. Our goal is to analyze task-specific redundancy
 172 in visual tokens that inspires method design.

173 4.1. Analysis Setup

174 **Model and Data.** We analyze the VILA-U model and
 175 collect attention statistics on both visual understanding
 176 (with ShareGPT-4v dataset) and visual generation (with
 177 JournyDB dataset). For each task, we record attention maps
 178 across all transformer layers.

179 **Attention Allocation.** Following prior attention decom-
 180 position analysis [3], we measure how attention mass is dis-
 181 tributed across token segments. Let $A_{i,j}^{(\ell,h)}$ denote the atten-
 182 tion weight at layer ℓ and head h from query token i to key
 183 token j , with

$$184 \sum_j A_{i,j}^{(\ell,h)} = 1.$$

185 Given a partition of tokens into segments (e.g., system,
 186 image, instruction, output), the *attention allocation* of segment S
 187 at layer ℓ is defined as:

$$188 \alpha_S^{(\ell)} = \frac{1}{H} \sum_{h=1}^H \sum_i \sum_{j \in S} A_{i,j}^{(\ell,h)}. \quad (2)$$

189 This metric captures the fraction of total attention mass
 190 directed to each token segment at a given layer. We use $\alpha_S^{(\ell)}$
 191 to quantify redundancy patterns across depth.

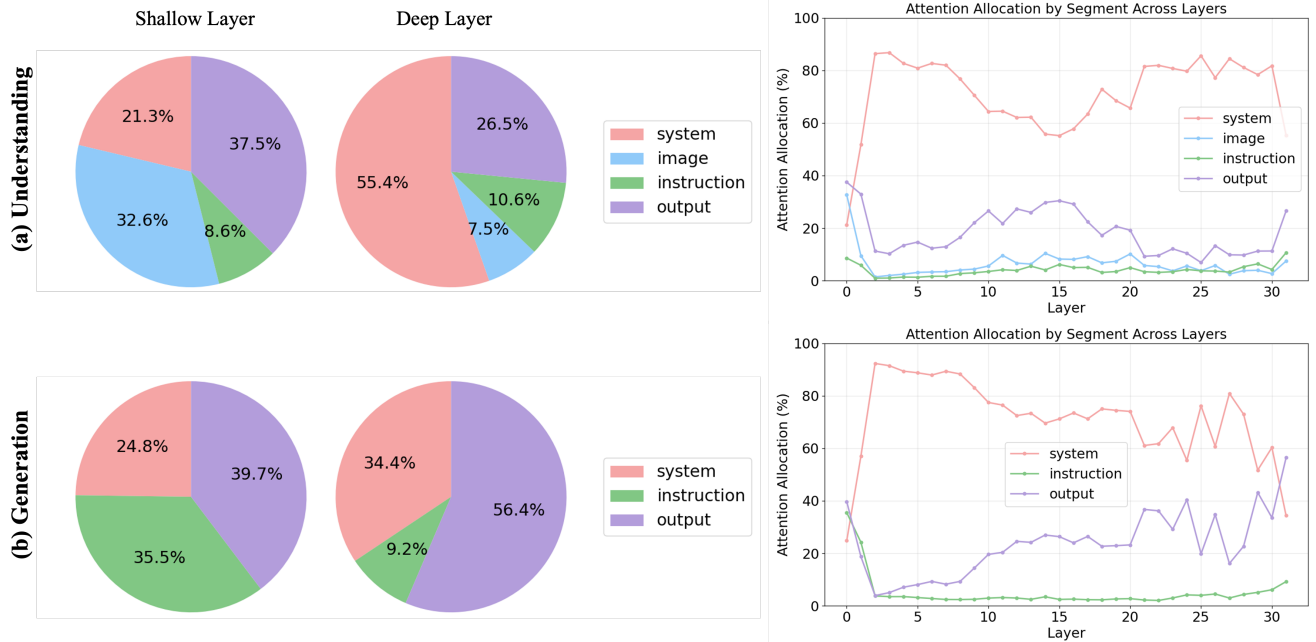


Figure 3. **Quantitative attention allocation reveals depth-dependent visual redundancy asymmetry.** Left: Attention mass distribution over token segments (system, image, instruction, output) at representative shallow and deep layers. Right: Layerwise attention allocation across the full transformer depth. For visual understanding (top), attention to image tokens sharply decreases in deeper layers, while instruction and output tokens dominate, indicating substantial late-layer visual redundancy. In contrast, visual generation (bottom) maintains consistently high attention to image tokens across layers, with increasing allocation to output image tokens in deep layers, reflecting persistent autoregressive dependence on generated image tokens.

192 4.2. Task-Specific Attention Patterns

193 We visualize (1) attention allocation (Figure 3) across token
194 segments over layers and (2) attention heatmaps (Figure 2)
195 at representative layers for both visual understanding (U)
196 and visual generation (G). The results reveal a clear asym-
197 metry in visual token redundancy patterns in different tasks.

198 **Visual Understanding (U).** For perception tasks (e.g.,
199 VQA), visual tokens exhibit clear depth-dependent redun-
200 dancy. As shown in Figure 2 and Figure 3, attention rapidly
201 shifts away from image tokens as depth increases. Image
202 tokens account for roughly $\sim 30\%$ of attention in the first
203 layer, but this drops below 10% in middle and late lay-
204 ers. Instead, attention becomes dominated by instruction
205 and output tokens, which together exceed 80% of the total
206 attention mass in deeper layers. Across layers, we observe
207 a consistent transition:

- 208 • **Early layers:** Strong cross-modal interactions between
209 image and text tokens, indicating visual grounding and
210 alignment.
- 211 • **Middle layers:** Attention increasingly concentrates on
212 text tokens, with diminishing image-to-image and image-
213 to-text interactions.
- 214 • **Late layers:** Attention is almost entirely confined to tex-
215 tual tokens, suggesting that high-level reasoning becomes
216 predominantly linguistic.

Visual Generation (G). In contrast to understanding, im-
age generation exhibits a persistent and structured depen-
dence on image tokens. Output (image) tokens receive a
substantial fraction of attention across early and late lay-
ers, typically ranging from 30% to 60%. Unlike the rapid
decay observed in understanding tasks, attention to image
tokens exhibits a consistent increase of attention allocation
in deeper layers. Across depth, the attention pattern follows
a hierarchical structure:

- 226 • **Early layers:** Broad attention over both textual prompts
227 and previously generated image tokens, establishing
228 global conditioning.
- 229 • **Middle layers:** Attention concentrates on recent image
230 tokens and specific prefix positions, reflecting localized
231 autoregressive dependencies.
- 232 • **Late layers:** Image-token attention becomes increasingly
233 significant, ensuring consistency in token prediction.

Robustness Across Scales. We repeat the same analy-
sis on a smaller-scale VILA-U model trained by ourselves
(LLaMA-3-3B backbone [5]). The qualitative and quanti-
tative patterns remain consistent: late-layer visual redun-
dancy emerges for understanding, while generation pre-
serves significant image-token attention. This suggests the
observed asymmetry is not scale-specific.

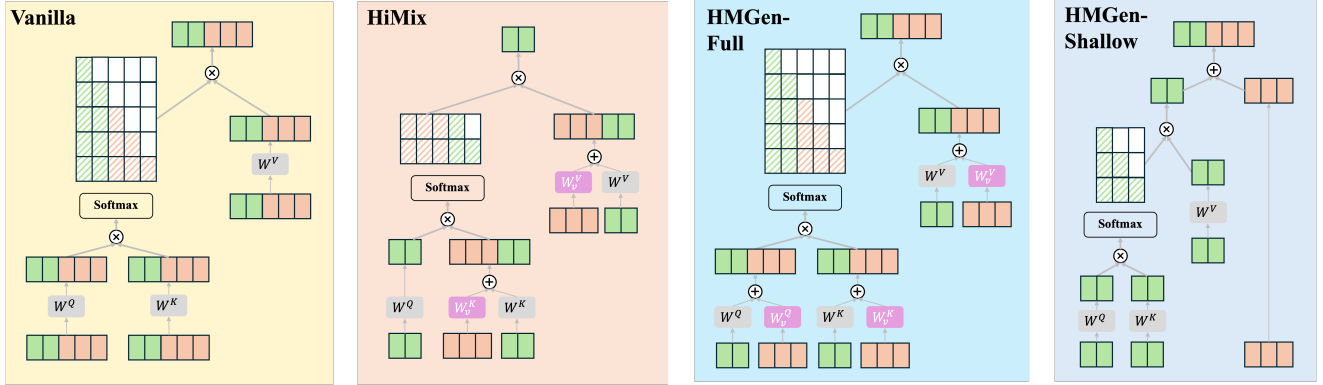


Figure 4. **Task-specific token-reduction-based acceleration mechanisms for unified VLM training.** From left to right: (1) Vanilla Transformer layer, where both text and image tokens participate fully in self-attention and feed-forward computation. (2) HiMix (Understanding) reduce image tokens from the query stream while retaining them in key/value projections, eliminating quadratic image-to-image attention while preserving text-to-image interactions. (3) HMGen-Full layer (Generation) maintains full autoregressive attention but separates image- and text-related projections for stable hierarchical conditioning. (4) HMGen-Shallow layer (Generation) skips image-token attention and feed-forward updates, forwarding their hidden states to reduce computation while preserving autoregressive structure.

241 4.3. Implications for Acceleration

242 This implies token reduction must be task-aware:

- 243 1. **Understanding:** Visual tokens are redundant after the
244 first few layers. It is possible to reduce image tokens
245 in some way and significantly reduce quadratic attention
246 cost with minimal performance impact.
- 247 2. **Generation:** Visual tokens are autoregressively gener-
248 ated during inference, and removing training computa-
249 tion on them must still enable the same next-token
250 prediction for inference. Deep layers in visual genera-
251 tion also have limited bandwidth to reduce image token-
252 related computation.

253 Therefore, a unified model cannot rely on a single token-
254 dropping rule. The structural roles of visual tokens differ
255 fundamentally between discriminative and generative objec-
256 tives. We introduce the task-specific training accelera-
257 tion methods below.

258 5. Proposed Task-Specific Accelerators

259 Motivated by the task-specific redundancy revealed in
260 Sec. 4, we investigate whether token-reduction-based accel-
261 eration can be done separately for visual understanding and
262 generation. We first evaluate these strategies in isolation
263 before analyzing their behavior under unified training.

264 5.1. Understanding (U)

265 **Method.** We adopt HiMix [46] as the baseline accelerator
266 for visual understanding. Unlike token merging/dropping
267 given complete image token sets based on inference-time
268 analysis, HiMix modifies the attention computation in a
269 manner compatible with both training and inference, mak-
270 ing it suitable for unified autoregressive VLMs. The key
271 idea is to *reduce tokens in queries*.

272 Specifically, as illustrated in Figure 4, image tokens are

273 removed from the *query* projections while retained in the
274 *key* and *value* projections. This eliminates quadratic image-
275 to-image attention while preserving text-to-image interac-
276 tions. As shown in Sec. 4, visual tokens become increas-
277 ingly redundant in deeper layers for understanding tasks;
278 thus, removing them from queries reduces computation
279 with minimal impact on prediction. Moreover, noticing that
280 the final output of each layer only includes text tokens as
281 image tokens are removed from queries, this strategy re-
282 quires the input original image tokens to each layer of the
283 transformer.

Theoretical Efficiency. For a sequence with T text tokens and M image tokens (total length $T + M$) and hidden size d , the per-layer complexity of a vanilla Transformer can be decomposed into: (i) self-attention, dominated by the QK^\top operation, $\mathcal{O}((T + M)^2d)$; and (ii) the feed-forward network (two linear layers), $\mathcal{O}(8(T + M)d^2)$. Thus,

$$\text{Cost}_{\text{base}} = \mathcal{O}((T + M)^2d + 8(T + M)d^2). \quad 290$$

291 With HiMix, image tokens are removed from the *query*
292 stream, so attention is computed only for T text queries
293 over $(T + M)$ keys/values, reducing the attention term to
294 $\mathcal{O}(T(T + M)d)$ while keeping the FFN term unchanged:

$$\text{Cost}_{\text{HiMix}} = \mathcal{O}(T(T + M)d + 8Td^2). \quad 295$$

296 When $M \gg T$, the dominant $\mathcal{O}(M^2d)$ attention term is
297 removed. In practice, this leads to substantial FLOPs re-
298 duction while preserving the cross-modal interactions nec-
299 essary for visual understanding.

Experimental Results. We evaluate HiMix in an
300 understanding-only setting of VILA-U (LLaMA-3-3B
301 backbone [5, 34]), with the ShareGPT-4v dataset [2]. We
302 follow VILA-U to conduct pretraining and finetuning each
303

Table 1. **HiMix for visual understanding.** Performance and computational cost comparison between the understanding-only VILA-U baseline and HiMix. HiMix reduces training FLOPs to 0.24× (76% reduction) by removing image-token queries, while incurring only moderate performance degradation across GQA, MME, POPE, and SeedBench benchmarks. The relatively small accuracy drop compared to the substantial computational savings confirms significant late-layer visual redundancy in understanding tasks.

Method	GQA	MME-C	MME-P	POPE-A	POPE-P	POPE-R	POPE-F1	SeedBench-Img	FLOPs
VILA-U (U-only)	52.86	258.21	1054.88	81.30	84.67	74.76	79.40	46.05	1×
HiMix (U-only)	49.92	224.64	983.30	78.56	78.03	79.49	78.75	40.88	0.24×

304 for one epoch, and evaluate on several visual understanding
305 benchmarks [8, 13, 15, 18]. Results are in Table 1.

306 HiMix reduces FLOPs to 0.24× of the baseline, corre-
307 sponding to a 76% reduction in computation. Despite this
308 substantial saving, performance degradation remains moder-
309 ate. For example, GQA accuracy decreases from 52.86 to
310 49.92, while POPE F1 drops slightly from 79.40 to 78.75.
311 Notably, the performance drop is significantly smaller than
312 the reduction in computational cost, indicating substantial
313 redundancy in late-layer visual processing for understand-
314 ing tasks. Overall, these results confirm that visual under-
315 standing exhibits considerable late-layer image redundancy.
316 Structured removal of image-token queries yields large ef-
317 ficiency gains while largely preserving cross-modal reason-
318 ing capability.

319 5.2. Generation (G)

320 **Design Constraints from Autoregressive Image Genera-**
321 **tion.** Unlike visual understanding, visual generation fol-
322 lows a strict autoregressive process: each predicted image
323 token is appended to the sequence and must serve as a
324 valid query for predicting subsequent tokens. Therefore,
325 image tokens *must remain in the query stream*. Remov-
326 ing them from queries would break the autoregressive de-
327 pendency chain and make inference inconsistent with train-
328 ing. This constraint fundamentally differentiates generation
329 from understanding and prevents directly applying HiMix-
330 style query removal.

331 One might instead consider removing image tokens from
332 key/value projections while keeping them in queries. Al-
333 though this reduces part of the attention computation, two
334 major issues arise. (1) **Limited FLOPs Reduction.** Even if
335 image-to-image attention is partially suppressed, the feed-
336 forward network (FFN) still processes all image tokens.
337 When $M \gg T$, the dominant $\mathcal{O}(8Md^2)$ FFN term remains
338 intact, resulting in minimal overall computational savings.
339 (2) **Severe Performance Degradation.** Image generation
340 exhibits persistent image-token dependence across depth
341 (Sec. 4). Suppressing key/value participation disrupts hier-
342 archical autoregressive conditioning, leading to substantial
343 quality degradation in practice. Empirically, we observe
344 that this naive modification yields both limited efficiency
345 gains and large drops in generative performance. For ex-
346 ample, applying this design to one middle layer leads to a
347 significant drop (-3.52) on MJHQ-30K [16].



Figure 5. **Inference-time-only HMGen.** Qualitative comparison between the original model (bottom) and inference-time-only HMGen (top), where image-token computation in shallow layers is skipped without retraining. Visual quality and semantic consistency are largely preserved despite reduced computation.

348 **HMGen: Hierarchical Mixture for Generation.** Moti-
349 vated by the hierarchical attention structure observed in
350 Sec. 4, we instead propose **HMGen**, which is composed
351 of two kinds of layers illustrated in Figure 4. HMGen pre-
352 serves the autoregressive structure with image in query *from*
353 *model level* while reducing tokens in specific layers.

354 We introduce K designated *shallow layers* in the middle
355 portion of the transformer (out of L total layers) while other
356 layers remain as *full layers*. This is because early full layers
357 are required to preserve global conditioning, while late full
358 layers are required to ensure high-fidelity final token predic-
359 tion. We empirically find that *alternating shallow and full*
360 *layers in the middle layers* yields the best trade-off between
361 efficiency and generation quality.

362 In *shallow layers*, image-token attention computation is
363 skipped, and the feed-forward network is applied only to
364 text tokens. The image-token hidden states are directly for-
365 forwarded from the previous layer to the next without partici-
366 pating in self-attention or FFN updates.

367 In *full layers*, we further introduce separate projection
368 parameters for image and text tokens. Although the back-
369 bone remains unified, decoupling image-related projections
370 stabilizes training and improves generation quality. This
371 separation allows image-token representations to maintain
372 dedicated pathways even when their participation in atten-
373 tion is selectively reduced. Empirically, we observe im-
374 proved performance compared to fully shared parameteriza-
375 tion under the same FLOPs budget.

376 **Theoretical Efficiency.** HMGen maintains the autore-
377 gressive dependency chain while reducing computation in
378 K designated middle “shallow” layers (out of L total) by
379 skipping image-token attention/MLP computation and for-
380 warding their hidden states. Using the same decomposition

381 as above, a vanilla layer costs

$$382 \quad \text{Cost}_{\text{base}} = \mathcal{O}((T + M)^2d + 8(T + M)d^2).$$

383 In a shallow layer, attention is computed only for T text
384 queries, giving $\mathcal{O}(T^2d)$, and the FFN is applied only to text
385 tokens, giving $\mathcal{O}(8Td^2)$:

$$386 \quad \text{Cost}_{\text{shallow}} = \mathcal{O}(T^2d + 8Td^2).$$

387 The total complexity across L layers is therefore

$$388 \quad \text{Cost}_{\text{HMGen}} = (L - K) \mathcal{O}((T + M)^2d + 8(T + M)d^2) \\ 389 \quad + K \mathcal{O}(T^2d + 8Td^2).$$

390 When $M \gg T$, each shallow layer removes the domi-
391 nant image-related costs in both attention and FFN, i.e., the
392 $\mathcal{O}(M^2d)$ and $\mathcal{O}(8Md^2)$ terms. Consequently, the overall
393 FLOPs reduction scales with the fraction of layers made
394 shallow (K/L), and is upper-bounded by the compute in
395 the remaining $(L - K)$ full layers. In the idealized regime
396 where shallow layers contribute negligible cost relative to
397 full layers, the relative cost approaches $1 - K/L$, yielding
398 an approximate speedup of $1/(1 - K/L)$.

399 **Experimental Results.** We first evaluate HMGen in a
400 generation-only setting of VILA-U using the JourneyDB
401 [31] dataset, and evaluate visual generation on MJHQ-30K
402 [16]. Quantitative results are shown in Table 2, and quali-
403 tative inference-time only results (without training, just di-
404 rectly skipping image-related computations in middle lay-
405 ers) are visualized in Figure 5.

406 **(1) Inference-Time Applicability.** Figure 5 demonstrates
407 that HMGen can be directly applied at inference time with-
408 out architectural modification. By design, shallow layers
409 preserve the autoregressive query structure, allowing image
410 tokens to be appended and used as subsequent queries dur-
411 ing generation. This confirms that HMGen is not merely
412 a training-time approximation but a structurally consistent
413 acceleration mechanism.

414 **(2) Reasonable FLOPs Reduction.** As shown in Table 2,
415 introducing K shallow layers yields significant computa-
416 tional savings. With $K = 3$, FLOPs are reduced to $0.85\times$
417 of the baseline, and with $K = 5$, to $0.75\times$. Since each
418 shallow layer removes both the dominant $\mathcal{O}(M^2d)$ atten-
419 tion term and the $\mathcal{O}(8Md^2)$ FFN term, the efficiency gain
420 scales approximately with the fraction of shallow layers.

421 **(3) Improved Generation Quality.** Notably, HMGen
422 achieves substantially better MJHQ-30K scores compared
423 to the generation-only VILA-U baseline ($17.45 \rightarrow 12.16$
424 with $K = 3$). This improvement arises from our separa-
425 tion of image and text projection parameters within the full
426 layers. By decoupling image-specific transformations, the
427 model maintains more stable hierarchical image representa-
428 tions even when computation is selectively reduced.

Table 2. **HMGen for visual generation.** Introducing shallow lay-
ers reduces FLOPs (to $0.85\times$ and $0.75\times$) while improving genera-
tive quality compared to the VILA-U baseline, demonstrating ef-
ficient and structure-aware acceleration.

Method	#Shallow Layers	MJHQ-30K	FLOPs
VILA-U (G-only)	0	17.45	$1\times$
HMGen	3	12.16	$0.85\times$
HMGen	5	12.55	$0.75\times$

Overall, HMGen not only reduces computation but also
enhances generative quality, demonstrating that hierarchi-
cal, structure-aware acceleration is better aligned with the
intrinsic dependencies of visual generation.

6. The Limits of Unified Efficiency

While the task-specific token-reduction-based accelerators
in Sec. 4 demonstrate substantial efficiency gains when ap-
plied to understanding or generation in isolation, our pri-
mary objective is to evaluate their behavior under unified
training. In unified VLMs, both objectives are optimized
jointly under a shared backbone, and improvements in one
task often influence the other through shared representa-
tions. Efficiency modifications may therefore interact with
cross-task learning dynamics in subtle ways. In this section,
we examine whether task-specific acceleration strategies re-
main effective in a unified setting, and identify structural
barriers that emerge during joint optimization.

6.1. Synergy Breakage: The Cost of Efficiency

Positive Cross-Task Synergy Baseline We first examine
the unified baseline without token reduction. From Table 3,
joint training improves both tasks relative to their single-
task counterparts. **Understanding improves under unified
training:** GQA increases from 52.86 (U-only) to 56 (Uni-
fied), POPE F1 from 79.40 to 82.3, and SeedBench from
46.05 to 47.88. **Generation also improves:** MJHQ-30K
improves from 17.45 (G-only) to 15.78 (Unified), indicat-
ing better generative quality. Formally, let performance on
understanding and generation be $\mathcal{U}(\theta)$ and $\mathcal{G}(\theta)$. For the
unified baseline,

$$\mathcal{U}(\theta_{\text{unified}}) > \mathcal{U}(\theta_{\text{U-only}}), \quad \mathcal{G}(\theta_{\text{unified}}) > \mathcal{G}(\theta_{\text{G-only}}).$$

This mutual improvement confirms the presence of positive
cross-task transfer, which motivates unified modeling.

**Severe Collapse with Fully Shared Parameters in
HiMix-HMGen.** The row *HiMix-HMGen (Share All)* in
Table 3 shows substantial degradation than HiMix (U-only)
and HMGen (g-only): GQA drops from 56 to 33, POPE
F1 from 82.3 to 67.59, SeedBench from 47.88 to 31.38,
while MJHQ-30K worsens from 12.16 to 12.53. Although
FLOPs are reduced to $0.56\times$, both tasks suffer significant
performance collapse. Notably, unified performance be-
comes worse than the single-task baseline in some metrics,

Table 3. **Unified training performance and efficiency.** The unified baseline improves both understanding (e.g., GQA 0.5600 vs. 0.5286 U-only) and generation (MJHQ 15.78 vs. 17.45 G-only), demonstrating positive cross-task synergy. However, combining HiMix and HMGen under joint training substantially reduces FLOPs (0.55–0.56 \times) but degrades performance on both objectives (e.g., GQA drops to 0.4705/0.3300 and MJHQ worsens to 14.54/12.53), indicating that task-specific token reduction disrupts mutual gains.

Method	GQA	MME-C	MME-P	POPE-A	POPE-P	POPE-R	POPE-F1	SeedBench-Img	MJHQ	FLOPs
VILA-U (U-only)	52.86	258.21	1054.88	81.30	84.67	74.76	79.40	46.05	–	1 \times
VILA-U (G-only)	–	–	–	–	–	–	–	–	17.45	1 \times
VILA-U (Unified)	56.00	250.00	1135.91	83.37	87.95	77.33	82.30	47.88	15.78	1 \times
HiMix (U-only)	49.92	224.64	983.30	78.56	78.03	79.49	78.75	40.88	–	0.24 \times
HMGen (G-only)	–	–	–	–	–	–	–	–	12.16	0.85 \times
HiMix-HMGen (Share All Params)	33.00	233.21	662.26	60.84	57.66	81.64	67.59	31.38	12.53	0.56 \times
HiMix-HMGen (Share Partial Params)	47.05	255.00	847.82	76.43	76.10	77.10	76.58	34.50	14.54	0.55 \times

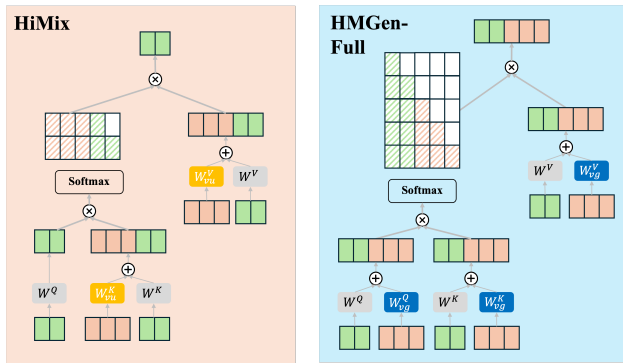


Figure 6. **Separate image projection strategy.** To reduce interference between task-specific routing, we decouple image-related projection parameters (e.g., W_v^Q , W_v^K , W_v^V) for HiMix (highlighted by yellow) and HMGen-Full (highlighted by blue), while keeping the backbone shared. This design aims to stabilize hierarchical image representations when token participation differs across tasks.

458 indicating negative transfer. Thus, naively combining task-
459 specific accelerators destroys cross-task synergy.

460 6.2. Separate Image Projection Strategy

To mitigate this issue, we introduce partial decoupling of image-related projections, as illustrated in Figure 6. Instead of fully shared projections, we decompose:

$$W_v^Q = \{W_{vu}^Q, W_{vx}^Q\},$$

461 and similarly for W_v^K and W_v^V . This creates semi-
462 independent image pathways while preserving the unified
463 backbone. From Table 3, *HiMix-HMGen (Share Partial)*
464 improves over the fully shared variant significantly: GQA
465 increases from 33 to 47.05, POPE F1 from 67.59 to 76.58,
466 and MJHQ-30K from 12.53 to 14.54. However, perfor-
467 mance still falls short of the unified baseline per understand-
468 ing, while better than the unified baseline on generation
469 (56 GQA and 15.78 MJHQ-30K), indicating that param-
470 eter separation partially restores synergy.

6.3. Structural Drivers of Synergy Loss

We discuss possible drivers of the observed synergy break-
472 age below to guide future investigation. Unified training
473 implicitly assumes a shared latent space $\phi(z; \theta)$, where dis-
474 criminative and generative signals co-shape representations.
475 Task-specific token dropping changes which tokens partic-
476 ipate in attention and which parameters receive gradients.
477 Consequently, gradients $\nabla_{\theta} \mathcal{L}_U$ and $\nabla_{\theta} \mathcal{L}_G$ are computed
478 under incompatible masking operators, leading to poten-
479 tially fragmented optimization dynamics. This hypothesis
480 is also supported by the separate image projection strategy.
481

Key Takeaway. Table 3 reveals a consistent pattern:
482 the unified baseline exhibits positive cross-task transfer,
483 whereas task-specific token reduction eliminates or reverses
484 these gains. Efficiency improvements achieved in isolation
485 do not compose under unified optimization. Effective uni-
486 fied acceleration must therefore preserve shared computa-
487 tional pathways that enable cross-task representation align-
488 ment, rather than simply aggregating task-optimal pruning
489 strategies.
490

7. Conclusion

We investigate the feasibility and limits of token-reduction-
492 based acceleration for unified vision-language models and
493 identify a fundamental asymmetry in visual token usage:
494 visual understanding exhibits substantial late-layer redun-
495 dancy, whereas visual generation maintains persistent
496 image-token dependence across depth. Based on this in-
497 sight, we design task-specific accelerators that achieve sig-
498 nificant efficiency gains in isolated settings; however, when
499 combined under unified training, they induce a consistent
500 *synergy loss*, as task-specific token dropping leads to diver-
501 gent parameter usage and removes the mutual performance
502 gains typically observed in joint optimization. Our findings
503 suggest that efficient unified modeling requires preserving
504 shared computational pathways that enable cross-task rep-
505 resentation alignment, rather than simply aggregating task-
506 specific strategies.
507

508

References

509

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

560

561

562

563

564

- [1] Daniel Bolya, Cheng-Yang Fu, Xiaoliang Dai, Peizhao Zhang, Christoph Feichtenhofer, and Judy Hoffman. Token merging: Your vit but faster, 2023. 1, 2
- [2] Lin Chen, Jinsong Li, Xiaoyi Dong, Pan Zhang, Conghui He, Jiaqi Wang, Feng Zhao, and Dahua Lin. Sharegpt4v: Improving large multi-modal models with better captions. In *European Conference on Computer Vision*, pages 370–387. Springer, 2024. 5
- [3] 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, 2024. 1, 2, 3
- [4] Wenliang Dai, Junnan Li, Dongxu Li, Anthony Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale Fung, and Steven Hoi. InstructBLIP: Towards general-purpose vision-language models with instruction tuning. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. 1
- [5] Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv e-prints*, pages arXiv–2407, 2024. 4, 5
- [6] Patrick Esser, Robin Rombach, and Björn Ommer. Taming transformers for high-resolution image synthesis, 2021. 1
- [7] Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, Dustin Podell, Tim Dockhorn, Zion English, Kyle Lacey, Alex Goodwin, Yan-nik Marek, and Robin Rombach. Scaling rectified flow transformers for high-resolution image synthesis, 2024. 1
- [8] Chaoyou Fu, Peixian Chen, Yunhang Shen, Yulei Qin, Mengdan Zhang, Xu Lin, Jinrui Yang, Xiawu Zheng, Ke Li, Xing Sun, Yunsheng Wu, Rongrong Ji, Caifeng Shan, and Ran He. MME: A comprehensive evaluation benchmark for multimodal large language models. In *The Thirty-ninth Annual Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2025. 6
- [9] Yuying Ge, Sijie Zhao, Ziyun Zeng, Yixiao Ge, Chen Li, Xintao Wang, and Ying Shan. Making llama see and draw with seed tokenizer, 2023. 1
- [10] Xiangming Gu, Tianyu Pang, Chao Du, Qian Liu, Fengzhuo Zhang, Cunxiao Du, Ye Wang, and Min Lin. When attention sink emerges in language models: An empirical view. *arXiv preprint arXiv:2410.10781*, 2024. 2
- [11] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models, 2020. 1
- [12] Wenbo Hu, Zi-Yi Dou, Liunian Harold Li, Amita Kamath, Nanyun Peng, and Kai-Wei Chang. Matryoshka query transformer for large vision-language models, 2024. 1, 2
- [13] Drew A. Hudson and Christopher D. Manning. Gqa: A new dataset for real-world visual reasoning and compositional question answering, 2019. 6
- [14] Yang Jin, Kun Xu, Kun Xu, Liwei Chen, Chao Liao, Jianchao Tan, Quzhe Huang, Bin Chen, Chenyi Lei, An Liu, Chengru Song, Xiaoqiang Lei, Di Zhang, Wenwu Ou, Kun Gai, and Yadong Mu. Unified language-vision pretraining in llm with dynamic discrete visual tokenization, 2024. 1
- [15] Bohao Li, Rui Wang, Guangzhi Wang, Yuying Ge, Yixiao Ge, and Ying Shan. Seed-bench: Benchmarking multimodal llms with generative comprehension, 2023. 6
- [16] Daiqing Li, Aleks Kamko, Ehsan Akhgari, Ali Sabet, Linmiao Xu, and Suhail Doshi. Playground v2.5: Three insights towards enhancing aesthetic quality in text-to-image generation, 2024. 6, 7
- [17] Wentong Li, Yuqian Yuan, Jian Liu, Dongqi Tang, Song Wang, Jie Qin, Jianke Zhu, and Lei Zhang. Tokenpacker: Efficient visual projector for multimodal llm, 2024. 2
- [18] Yifan Li, Yifan Du, Kun Zhou, Jinpeng Wang, Wayne Xin Zhao, and Ji-Rong Wen. Evaluating object hallucination in large vision-language models, 2023. 6
- [19] Yanwei Li, Chengyao Wang, and Jiaya Jia. Llama-vid: An image is worth 2 tokens in large language models, 2023. 2
- [20] Zhihang Lin, Mingbao Lin, Luxi Lin, and Rongrong Ji. Boosting multimodal large language models with visual tokens withdrawal for rapid inference, 2025. 2
- [21] Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning, 2023. 1
- [22] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. In *NeurIPS*, 2023. 1
- [23] Hao Liu, Wilson Yan, Matei Zaharia, and Pieter Abbeel. World model on million-length video and language with blockwise ringattention. In *The Thirteenth International Conference on Learning Representations*, 2025. 1
- [24] Gen Luo, Yiyi Zhou, Tianhe Ren, Shengxin Chen, Xiaoshuai Sun, and Rongrong Ji. Cheap and quick: Efficient vision-language instruction tuning for large language models, 2023. 1, 2
- [25] Chuofan Ma, Yi Jiang, Junfeng Wu, Jihan Yang, Xin Yu, Zehuan Yuan, Bingyue Peng, and Xiaojuan Qi. Unitok: A unified tokenizer for visual generation and understanding, 2025. 1
- [26] Yiyang Ma, Xingchao Liu, Xiaokang Chen, Wen Liu, Chengyue Wu, Zhiyu Wu, Zizheng Pan, Zhenda Xie, Haowei Zhang, Xingkai yu, Liang Zhao, Yisong Wang, Jiaying Liu, and Chong Ruan. Janusflow: Harmonizing autoregression and rectified flow for unified multimodal understanding and generation, 2024. 1
- [27] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision, 2021. 1
- [28] Yongming Rao, Wenliang Zhao, Benlin Liu, Jiwen Lu, Jie Zhou, and Cho-Jui Hsieh. Dynamicvit: Efficient vision transformers with dynamic token sparsification. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2021. 1, 2
- [29] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models, 2022. 1, 2
- [30] Yuzhang Shang, Mu Cai, Bingxin Xu, Yong Jae Lee, and Yan Yan. Llva-prumerge: Adaptive token reduction for efficient

565

566

567

568

569

570

571

572

573

574

575

576

577

578

579

580

581

582

583

584

585

586

587

588

589

590

591

592

593

594

595

596

597

598

599

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

- 622 large multimodal models. *arXiv preprint arXiv:2403.15388*,
623 2024. 1, 2
- 624 [31] Keqiang Sun, Junting Pan, Yuying Ge, Hao Li, Haodong
625 Duan, Xiaoshi Wu, Renrui Zhang, Aojun Zhou, Zipeng Qin,
626 Yi Wang, et al. Journeymb: A benchmark for generative im-
627 age understanding. *Advances in neural information process-*
628 *ing systems*, 36:49659–49678, 2023. 7
- 629 [32] Chameleon Team. Chameleon: Mixed-modal early-fusion
630 foundation models, 2024. 1, 2
- 631 [33] Keyu Tian, Yi Jiang, Zehuan Yuan, Bingyue Peng, and Li-
632 wei Wang. Visual autoregressive modeling: Scalable image
633 generation via next-scale prediction, 2024. 1
- 634 [34] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier
635 Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste
636 Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aure-
637 lien Rodriguez, Armand Joulin, Edouard Grave, and Guil-
638 laume Lample. Llama: Open and efficient foundation lan-
639 guage models, 2023. 5
- 640 [35] Aaron van den Oord, Oriol Vinyals, and Koray
641 Kavukcuoglu. Neural discrete representation learning,
642 2018. 2
- 643 [36] Xinlong Wang, Xiaosong Zhang, Zhengxiong Luo, Quan
644 Sun, Yufeng Cui, Jinsheng Wang, Fan Zhang, Yueze Wang,
645 Zhen Li, Qiyang Yu, Yingli Zhao, Yulong Ao, Xuebin Min,
646 Tao Li, Boya Wu, Bo Zhao, Bowen Zhang, Liangdong Wang,
647 Guang Liu, Zheqi He, Xi Yang, Jingjing Liu, Yonghua Lin,
648 Tiejun Huang, and Zhongyuan Wang. Emu3: Next-token
649 prediction is all you need, 2024. 1, 2
- 650 [37] Chengyue Wu, Xiaokang Chen, Zhiyu Wu, Yiyang Ma,
651 Xingchao Liu, Zizheng Pan, Wen Liu, Zhenda Xie, Xingkai
652 Yu, Chong Ruan, et al. Janus: Decoupling visual encoding
653 for unified multimodal understanding and generation. *arXiv*
654 *preprint arXiv:2410.13848*, 2024. 1, 2
- 655 [38] Junfeng Wu, Yi Jiang, Chuofan Ma, Yuliang Liu, Heng-
656 shuang Zhao, Zehuan Yuan, Song Bai, and Xiang Bai. Liq-
657 uid: Language models are scalable and unified multi-modal
658 generators. *arXiv preprint arXiv:2412.04332*, 2024. 1
- 659 [39] Yecheng Wu, Zhuoyang Zhang, Junyu Chen, Haotian Tang,
660 Dacheng Li, Yunhao Fang, Ligeng Zhu, Enze Xie, Hongxu
661 Yin, Li Yi, et al. Vila-u: a unified foundation model inte-
662 grating visual understanding and generation. *arXiv preprint*
663 *arXiv:2409.04429*, 2024. 1, 2
- 664 [40] Guangxuan Xiao, Yuandong Tian, Beidi Chen, Song Han,
665 and Mike Lewis. Efficient streaming language models with
666 attention sinks. *arXiv*, 2023. 2
- 667 [41] Jinheng Xie, Weijia Mao, Zechen Bai, David Junhao Zhang,
668 Weihao Wang, Kevin Qinghong Lin, Yuchao Gu, Zhijie
669 Chen, Zhenheng Yang, and Mike Zheng Shou. Show-o:
670 One single transformer to unify multimodal understanding
671 and generation. *arXiv preprint arXiv:2408.12528*, 2024. 1,
672 2
- 673 [42] Lili Yu, Bowen Shi, Ramakanth Pasunuru, Benjamin Muller,
674 Olga Golovneva, Tianlu Wang, Arun Babu, Binh Tang, Brian
675 Karrer, Shelly Sheynin, Candace Ross, Adam Polyak, Rus-
676 sell Howes, Vasu Sharma, Puxin Xu, Hovhannes Tamoyan,
677 Oron Ashual, Uriel Singer, Shang-Wen Li, Susan Zhang,
678 Richard James, Gargi Ghosh, Yaniv Taigman, Maryam
Fazel-Zarandi, Asli Celikyilmaz, Luke Zettlemoyer, and Ar-
679 men Aghajanyan. Scaling autoregressive multi-modal mod-
680 els: Pretraining and instruction tuning, 2023. 681
- [43] Jun Zhan, Junqi Dai, Jiasheng Ye, Yunhua Zhou, Dong
682 Zhang, Zhigeng Liu, Xin Zhang, Ruibin Yuan, Ge Zhang,
683 Linyang Li, Hang Yan, Jie Fu, Tao Gui, Tianxiang Sun, Yu-
684 Gang Jiang, and Xipeng Qiu. Anygpt: Unified multimodal
685 llm with discrete sequence modeling, 2025. 1 686
- [44] Junyang Zhang, Mu Yuan, Ruiguang Zhong, Puhan Luo,
687 Huiyou Zhan, Ningkan Zhang, Chengchen Hu, and Xi-
688 angyang Li. A-vl: Adaptive attention for large vision-
689 language models, 2025. 1, 2 690
- [45] Shaolei Zhang, Qingkai Fang, Zhe Yang, and Yang Feng.
691 Llava-mini: Efficient image and video large multimodal
692 models with one vision token, 2025. 2 693
- [46] Xuange Zhang, Dengjie Li, Bo Liu, Zenghao Bao, Yao Zhou,
694 Baisong Yang, Zhongying Liu, Yujie Zhong, Zheng Zhao,
695 and Tongtong Yuan. Himix: Reducing computational com-
696 plexity in large vision-language models, 2025. 2, 5 697