

Reimplementation and Fine-Grained Analysis of Visual Token Sparsification Methods

Anonymous ACL submission

Abstract

Visual token sparsification represents a critical optimization technique for reducing computational overhead in Vision-Language Models (VLMs) while maintaining performance. This work presents a comprehensive reimplementation and fine-grained analysis of two state-of-the-art methods: VisionZip, a text-agnostic approach utilizing CLIP attention weights, and SparseVLM, a text-aware method employing cross-modal attention. Through systematic evaluation on the POPE dataset with 6,823 samples, we demonstrate that both methods achieve near-equivalent performance (Accuracy: 80.3% vs 80.6%) at 64-token sparsification (88.9% reduction), despite fundamentally different architectural designs. Our analysis reveals that text-aware token selection provides marginal advantages in precision (0.7%) but fails to meaningfully improve recall, suggesting that for object existence tasks, question-adaptive selection may not provide substantial benefits over question-agnostic approaches.

1 Introduction

Modern Vision-Language Models (VLMs) typically process 576 visual tokens per image, representing a significant computational bottleneck. Visual token sparsification methods aim to reduce this number while preserving model performance, with recent approaches demonstrating that 88.9% token reduction (from 576 to 64 tokens) can be achieved with minimal performance degradation. However, the optimal strategy for token selection remains an open question: should tokens be selected based on intrinsic visual importance independent of the query, or should selection adapt to the specific question being asked?

This work addresses this question through a rigorous comparative analysis of two representative approaches: VisionZip (1), which employs text-agnostic token selection based on CLIP attention

patterns, and SparseVLM (3), which utilizes text-aware selection through cross-modal attention. By integrating, debugging, and modifying the official implementations of both methods, we successfully configured them for evaluation under identical experimental conditions. This process allowed us to provide insights into the mechanisms, trade-offs, and limitations of each approach.

Our contributions are threefold: (1) a successful replication of the methods, paying careful attention to implementation details and necessary modifications for compatibility, (2) a comprehensive evaluation on the POPE dataset with 6,823 samples, reporting multiple metrics, and (3) a fine-grained analysis of token selection mechanisms, performance characteristics, and failure modes.

2 Method Overview

2.1 VisionZip: Text-Agnostic Token Selection

VisionZip implements a two-stage token selection mechanism that operates entirely within the vision encoder, independent of the language component. The method is grounded in the observation that CLIP token attention weights in deep layers of the CLIP encoder capture global visual saliency.

VisionZip operates in two stages as described in the original paper (1). In Stage 1, the method selects 54 dominant tokens by computing CLIP token attention weights at layer 22 of the CLIP encoder and selecting the top-k tokens with highest attention scores. The CLIP token attention naturally captures global visual saliency, prioritizing prominent objects and high-contrast regions.

In Stage 2, VisionZip selects 10 contextual tokens from the remaining tokens through similarity-based clustering. The method normalizes the remaining tokens and samples target tokens uniformly, then assigns other tokens to the nearest target based on cosine similarity. This ensures diverse spatial coverage, capturing background con-

081 text and less prominent visual regions that may be
082 missed by the dominant selection.

083 **2.2 SparseVLM: Text-Aware Progressive**
084 **Token Selection**

085 SparseVLM implements a fundamentally different
086 approach: token selection adapts dynamically to
087 the input question by leveraging cross-modal atten-
088 tion between visual and text tokens. The method ap-
089 plies progressive sparsification at multiple decoder
090 layers, allowing the model to gradually refine token
091 selection based on question-specific requirements.

092 SparseVLM computes cross-modal attention be-
093 tween visual tokens and text tokens at decoder
094 layers 2, 6, and 15. The method uses standard
095 scaled dot-product attention (2) to compute atten-
096 tion weights between text token queries and visual
097 token keys. The relevance of each visual token
098 to the question is determined by averaging atten-
099 tion scores from all text tokens to that visual token,
100 effectively measuring how much the question “at-
101 tends” to each visual region.

102 The method applies progressive sparsification:
103 for 64 total tokens, it retains 66 tokens at layer
104 2, 30 tokens at layer 6, and 17 tokens at layer 15.
105 At each layer, tokens with highest text-visual rele-
106 vance scores are selected via top-k selection, and
107 remaining tokens are merged into selected tokens
108 using attention-weighted pooling. This progres-
109 sive approach allows the model to maintain rich
110 visual information in early layers while gradually
111 focusing on question-relevant tokens.

112 **2.3 Key Differences**

113 The fundamental difference between VisionZip and
114 SparseVLM lies in their token selection basis. Vi-
115 sionZip is text-agnostic: it selects tokens based
116 purely on visual features (CLS attention weights)
117 without considering the input question. This means
118 the same set of tokens is selected for an image re-
119 gardless of what question is asked. SparseVLM
120 is text-aware: it adapts token selection to each
121 specific question by using cross-modal attention
122 between visual and text tokens. This allows the
123 method to focus on different image regions depend-
124 ing on what is being asked, potentially providing
125 advantages for questions requiring attention to less
126 prominent visual regions.

127 **3 Experimental Setup**

128 **3.1 Dataset: POPE**

129 The Polling-based Object Probing Evaluation
130 (POPE) dataset (4) is specifically designed to eval-
131 uate object hallucination in VLMs. The dataset
132 consists of yes/no questions about object exist-
133 ence, where positive cases refer to objects that ex-
134 ist in the image, and negative cases refer to objects
135 that do not exist.

136 The dataset contains 6,823 total questions di-
137 vided into three categories: Adversarial (2,274 sam-
138 ples), Popular (2,274 samples), and Random (2,275
139 samples). Images are sourced from the COCO
140 val2014 dataset. The LLaVA model zoo reports
141 that Vanilla LLaVA-1.5-7B achieves 85.9% ac-
142 curacy on POPE with full token retention (576 to-
143 kens), which we use as a reference baseline (not
144 re-evaluated in this work).

145 **3.2 Model Configuration**

146 We use LLaVA-1.5-7B (5)
147 (liuhaojian/llava-v1.5-7b) as the base
148 model for both methods. For VisionZip, we
149 configure dominant tokens as 54 and contextual
150 tokens as 10, totaling 64 tokens as specified in the
151 paper (1). Token selection occurs at CLIP encoder
152 layer 22. For SparseVLM, we configure retained
153 tokens as 64 total, with progressive sparsification
154 applying token counts of [66, 30, 17] at decoder
155 layers [2, 6, 15] respectively, matching the
156 method’s design.

157 **3.3 Hardware and Software Environment**

158 All experiments were conducted on the CS de-
159 partment’s instructional GPU cluster, specifically
160 on the instgpu-01.cs.wisc.edu machine. This
161 server is equipped with eight NVIDIA GeForce
162 RTX 2080 Ti GPUs, each with approximately 11
163 GB of VRAM. Our experiments were primarily
164 executed on GPUs 6 and 7 of this shared machine.

165 The software environment was built on a Linux
166 operating system with CUDA version 12.4 and
167 NVIDIA driver version 550.67. We used Python
168 3.10 with PyTorch 2.1.2. Key libraries included
169 transformers, bitsandbytes for 4-bit quantiza-
170 tion attempts, and accelerate for model handling.
171 Due to the 11 GB memory limitation of the RTX
172 2080 Ti GPUs when loading the 7B parameter
173 LLaVA model (which requires 14 GB in float16),
174 we utilized the device_map="auto" feature from
175 the accelerate library. This configuration auto-

176 matically offloaded some model layers to the CPU
177 to prevent out-of-memory errors, a necessary deviation
178 for running these large models on the available
179 hardware.

180 3.4 Evaluation Metrics

181 We report four standard classification metrics: Accuracy,
182 Precision, Recall, and F1 Score. Metrics are
183 calculated separately for each category (adversarial,
184 popular, random) and then averaged to obtain
185 overall performance. Since the VisionZip paper
186 reports accuracy as the primary metric for POPE,
187 we use accuracy as our primary comparison metric
188 to ensure fair comparison with published results.

189 4 Results

190 4.1 Overall Performance

191 Table 1 presents comprehensive results on the
192 POPE dataset.

193 4.2 Comparison with Published Results

194 The VisionZip paper reports “raw benchmark accu-
195 racy” for the POPE dataset at a 64-token configura-
196 tion (1). Their findings show VisionZip achieving
197 77.0% accuracy and SparseVLM achieving 75.1%
198 accuracy. This establishes a performance hierarchy
199 where VisionZip outperforms SparseVLM by 1.9
200 percentage points.

201 Our reimplementation yields different results in
202 two key aspects. First, our measured accuracies
203 (VisionZip: 80.3%, SparseVLM: 80.6%) are
204 notably higher than those reported in the paper. This
205 discrepancy may stem from differences in the eval-
206 uation protocol, minor implementation details, or
207 the specific hardware environment. Second, the per-
208 formance ranking is reversed: in our experiments,
209 SparseVLM surpasses VisionZip by a margin of 0.3
210 percentage points. While this margin is minimal
211 and likely within the bounds of experimental varia-
212 nce, it contradicts the paper’s conclusion. The key
213 takeaway is that our reimplementation finds the two
214 methods to have virtually equivalent performance
215 on this task.

216 5 Fine-Grained Analysis

217 5.1 Token Selection Mechanism Analysis

218 VisionZip’s Design:

219 VisionZip’s two-stage selection mechanism is
220 designed to balance global saliency with spatial
221 diversity. The method uses CLS token attention

222 weights at layer 22 of the CLIP encoder, operat-
223 ing on the principle that deep-layer CLS attention
224 aggregates global visual information. By design,
225 the dominant token selection (54 tokens) should
226 capture visually salient regions—typically objects,
227 high-contrast areas, and text—while the context-
228 ual token selection (10 tokens) through similarity-
229 based clustering should provide spatial coverage of
230 less prominent regions.

231 A key characteristic of this approach is that to-
232 ken selection is question-independent: the same 64
233 tokens are selected for an image regardless of what
234 is being asked. This design choice prioritizes effi-
235 ciency but theoretically limits the method’s ability
236 to adapt when questions focus on background or
237 non-salient objects.

238 SparseVLM’s Design:

239 SparseVLM employs a fundamentally different
240 strategy: progressive, question-aware token selec-
241 tion across multiple decoder layers. The method
242 computes cross-modal attention between text and
243 visual tokens, allowing token selection to adapt
244 dynamically based on the specific question. The
245 progressive schedule (66 tokens at layer 2, 30 at
246 layer 6, 17 at layer 15) is designed to maintain
247 rich visual information early in processing while
248 gradually focusing on question-relevant regions.

249 Theoretically, this text-aware approach should
250 provide advantages when questions require atten-
251 tion to specific regions that differ from global
252 visual saliency patterns. However, our empiri-
253 cal results show that both methods achieve nearly
254 identical performance (80.3% vs 80.6% accuracy),
255 suggesting that for POPE’s object existence ques-
256 tions—which primarily query about prominent,
257 visually salient objects—the adaptability of text-
258 aware selection does not translate to measurable
259 performance gains.

260 5.2 Precision-Recall Trade-off Analysis

261 Both methods exhibit a **high precision, moderate**
262 **recall** profile, indicating conservative prediction
263 strategies. This is particularly important for hallu-
264 cination detection, where false positives are costly.

265 **Precision Analysis:** VisionZip achieves 94.5%
266 precision, while SparseVLM achieves 95.2% (0.7%
267 higher). The text-aware selection provides more
268 relevant visual tokens, leading to slightly higher
269 confidence in predictions and reducing false pos-
270 itives.

271 **Recall Analysis:** Both methods achieve nearly
272 identical recall (64.9% vs 65.0%), indicating that at

Method	Accuracy	Precision	Recall	F1 Score
VisionZip	0.803 (80.3%)	0.945 (94.5%)	0.649 (64.9%)	0.769 (76.9%)
SparseVLM	0.806 (80.6%)	0.952 (95.2%)	0.650 (65.0%)	0.772 (77.2%)
Difference	-0.003 (-0.3%)	-0.007 (-0.7%)	-0.001 (-0.1%)	-0.003 (-0.3%)

Table 1: Performance comparison on POPE dataset. Both methods achieve near-equivalent performance with minimal differences across all metrics.

64 tokens (88.9% reduction), both struggle equally with detecting all relevant objects. This suggests that 64 tokens may be insufficient to capture all object information needed for high recall, and neither selection strategy can overcome the fundamental information loss from aggressive sparsification.

5.3 Theoretical Limitations and Failure Modes

Based on the architectural designs and our quantitative results, we can identify theoretical limitations for each method.

VisionZip’s Potential Limitations: The question-independent selection means that when POPE asks about objects in background regions, the CLS attention-based selection may prioritize foreground tokens over the queried object. Questions about small or non-salient objects may fail if those objects are not among the 64 selected tokens. The high precision (94.5%) but moderate recall (64.9%) we measured is consistent with this conservative selection strategy.

SparseVLM’s Potential Limitations: The progressive reduction schedule (66→30→17 tokens) means that if critical visual information is not captured at layer 2, later layers cannot recover it. The method’s reliance on cross-modal attention means that poor question representations could lead to suboptimal token selection. Despite the theoretical advantage of question-aware selection, SparseVLM achieves only marginally higher precision (95.2%) and nearly identical recall (65.0%) compared to VisionZip, suggesting these theoretical advantages may not materialize for simple object existence queries.

5.4 Theoretical Implications

Our experimental findings suggest that for object existence tasks like POPE, which primarily focus on prominent objects, text-aware token selection does not provide substantial advantages over text-

agnostic selection. This is because POPE questions typically query about visually salient objects that are naturally captured by both selection strategies. However, we hypothesize that for more complex tasks requiring spatial reasoning, attribute understanding, or multi-object relationships, text-aware selection may provide advantages by allowing the model to focus on different image regions depending on the question.

At 64 tokens (88.9% reduction from 576), both methods achieve approximately 80% accuracy, compared to the baseline performance. The absolute performance drop suggests that 64 tokens capture most but not all critical information needed for accurate object existence detection. Further token reduction would likely degrade performance more significantly, suggesting that 64 tokens may be near the lower bound for maintaining reasonable performance on this task.

6 Limitations

This study has several limitations: (1) evaluation on POPE only—results may not generalize to other benchmarks, (2) single sparsification level—only 64 tokens evaluated, (3) accuracy values higher than paper—our accuracy (80.3%, 80.6%) exceeds paper’s (77.0, 75.1%), suggesting potential differences in evaluation protocol, (4) implementation variance—subtle differences may affect results, and (5) hardware constraints—the 7B LLaVA model requires approximately 14GB in float16, exceeding the 11GB capacity of our RTX 2080 Ti GPUs. Both methods used `device_map="auto"` which likely offloaded some model layers to CPU. This memory constraint may have affected absolute performance, though both methods were evaluated under identical conditions.

348 7 Conclusion

349 This work presents a comprehensive reimplementa-
350 tion and fine-grained analysis of two visual token
351 sparsification methods: VisionZip (text-agnostic)
352 and SparseVLM (text-aware). Through rigorous
353 evaluation on the POPE dataset with 6,823 sam-
354 ples, we demonstrate that both methods achieve
355 near-equivalent performance (Accuracy: 80.3% vs
356 80.6%) at 64-token sparsification, despite funda-
357 mentally different architectural designs.

358 Our analysis reveals that for object exis-
359 tence tasks, text-agnostic and text-aware selection
360 achieve essentially identical performance, suggest-
361 ing that question-adaptive selection may not pro-
362 vide substantial benefits for simple binary questions
363 about prominent objects. However, SparseVLM’s
364 text-aware selection provides marginal precision
365 advantage (0.7%), indicating better alignment be-
366 tween visual tokens and question requirements.
367 Both methods struggle equally with recall at high
368 sparsification levels, suggesting that 64 tokens may
369 be insufficient regardless of selection strategy.

370 These findings suggest that for object existence
371 tasks, the choice between text-agnostic and text-
372 aware selection may not significantly impact per-
373 formance. However, for more diverse question
374 types requiring attention to different image regions,
375 text-aware selection may provide advantages. Fu-
376 ture work should evaluate both methods on more
377 diverse benchmarks to test this hypothesis.

378 Acknowledgments

379 We thank the authors of VisionZip and SparseVLM
380 for making their code publicly available, enabling
381 this reimplementation study.

382 References

- 383 [1] Yang, S., Chen, Y., Tian, Z., Wang, C., Li, J., Yu,
384 B., & Jia, J. (2024). VisionZip: Longer is Better but
385 Not Necessary in Vision Language Models. *arXiv*
386 preprint arXiv:2412.04467.
- 387 [2] Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J.,
388 Jones, L., Gomez, A. N., ... Polosukhin, I. (2017).
389 Attention is all you need. *Advances in Neural Infor-*
390 *mation Processing Systems (NeurIPS)*.
- 391 [3] Zhang, Y., et al. (2024). SparseVLM: Visual Token
392 Sparsification for Efficient Vision-Language Model
393 Inference. *arXiv preprint arXiv:2410.04417*.
- 394 [4] Li, Y., Du, Y., Zhou, K., Wang, J., Zhao, W. X., &
395 Wen, J. R. (2023). Evaluating Object Hallucination in
Large Vision-Language Models. *Proceedings of the
2023 Conference on Empirical Methods in Natural
Language Processing (EMNLP)*.
396 [5] Liu, H., Li, C., Wu, Q., & Lee, Y. J. (2023). Visual
397 Instruction Tuning. *Advances in Neural Information
398 Processing Systems (NeurIPS)*.
399 [6] ...
400 [7] ...
401 [8] ...