

Toward Compact and Structured Visual Representations in VLMs: SSM-Based Vision Encoders as an Alternative to Transformers

Anonymous CVPR submission

Paper ID ****

Abstract

001 *Spatial concept understanding is important for vision-*
002 *language models (VLMs), not only for explicit grounding*
003 *and localization tasks but also for general visual question*
004 *answering. The common solution to this challenge has been*
005 *to scale image resolution or visual token counts, yet this*
006 *conflicts with the quadratic attention cost of ViT-based en-*
007 *coders that nearly all current VLMs rely on. We explore*
008 *an orthogonal question and investigate whether a vision en-*
009 *coder can produce compact and structured visual token rep-*
010 *resentations that carry richer spatial concept information*
011 *under a fixed token budget. We conduct the first controlled*
012 *evaluation of SSM-based vision encoders as frozen visual*
013 *backbones in VLMs and find that VMamba-based VLMs*
014 *achieve stronger spatial concept understanding than ViT-*
015 *family alternatives at matched scale and token budget, with*
016 *improvements that transfer from grounding benchmarks to*
017 *general visual question answering. Token-region similar-*
018 *ity maps computed at intermediate layers of the language*
019 *model show that SSM visual tokens produce sharper and*
020 *more spatially selective concept-region binding during LLM*
021 *reasoning, showing that the language model can use SSM*
022 *representations better than ViT-based vision backbones un-*
023 *der similar settings. These findings suggest that architec-*
024 *tural inductive bias is an underexplored direction for im-*
025 *proving visual concept representations in VLMs without in-*
026 *creasing token budgets.*

027 1. Introduction

028 Spatial concept understanding is a core requirement for
029 vision-language models (VLMs). Beyond explicit ground-
030 ing and localization tasks, spatial understanding also plays
031 an important role in general visual question answering. For
032 example, GQA [10] is built around scene graphs with ex-
033 plicit spatial relationships, and several works show that im-
034 proving visual grounding leads to better performance on
035 both GQA and VQA-v2 [14, 21].

A recurring challenge in VLMs is extracting spatially 036
grounded evidence from images under a fixed token bud- 037
get [4, 9]. The dominant response has been to increase 038
image resolution or the number of visual tokens [15, 039
29]. However, this approach has a fundamental limita- 040
tion: virtually all current VLMs rely on ViT-based vi- 041
sion encoders [12], whose self-attention mechanism scales 042
quadratically with the number of tokens. Increasing reso- 043
lution or token count therefore quickly becomes computa- 044
tionally infeasible. More importantly, this approach treats a 045
representation quality problem as a quantity problem, with- 046
out questioning whether the visual tokens themselves carry 047
sufficient spatial concept information. 048

We argue that the more principled question is whether a 049
vision encoder can produce *compact and structured* visual 050
representations that encode richer spatial concept informa- 051
tion per token, without increasing the token budget. This is 052
also well aligned with the goal of visual concept learning, 053
which seeks compact and structured representations of the 054
visual world. Notably, existing VLMs rely almost exclu- 055
sively on ViT-family encoders, leaving the representation 056
design space largely underexplored. 057

State space model (SSM) based vision encoders are a 058
promising alternative. Unlike ViTs, which rely on global 059
self-attention over a flattened token sequence, SSM vision 060
models such as VMamba [19] build representations through 061
structured 2D state-space scanning, where each token inte- 062
grates spatial context from its 2D neighborhood across mul- 063
tiple scan directions. This mechanism embeds spatial rela- 064
tional structure into every token by construction, rather than 065
relying on positional encodings that may be underutilized 066
under standard classification pretraining. To our knowl- 067
edge, no prior work has conducted a controlled evaluation 068
of SSM encoders as visual backbones in generative VLMs, 069
isolating architectural effects from other confounding fac- 070
tors such as training data, connector design, and pretraining. 071

In this work, we conduct such a controlled evaluation us- 072
ing a LLaVA-style [18] VLM framework, where we swap 073
only the vision backbone while keeping the rest of the 074
pipeline fixed. Our main findings are: (i) VLMs built on 075

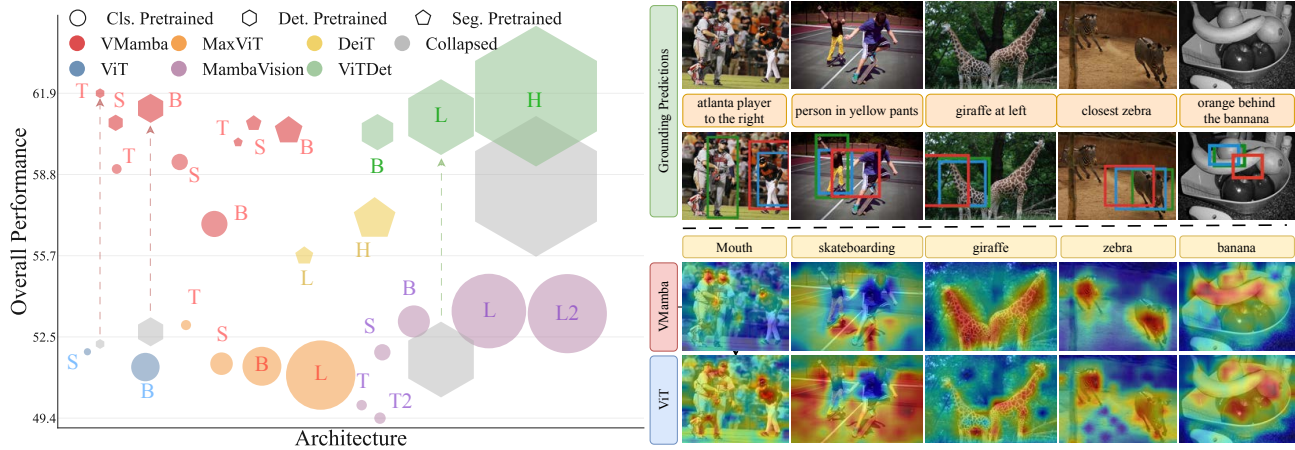


Figure 1. **Left: Overall VLM performance across backbone families.** Each marker represents one backbone configuration plugged into the same VLM pipeline. Colors denote backbone family; marker shapes denote pretraining objective; marker size reflects encoder scale. Gray markers indicate collapsed configurations, with arrows pointing to their stabilized variants. VMamba (red) consistently achieves stronger overall performance than ViT-family alternatives at matched or smaller model scale, despite similar ImageNet accuracy. **Right: Grounding predictions and token-region similarity maps.** Top: predicted bounding boxes from VMamba-T (blue) and ViT-S (red) against ground truth (green). Bottom: similarity maps between visual tokens and text concept tokens computed at an intermediate layer of the language model, showing that VMamba-T tokens produce sharper and more spatially selective concept-region binding than ViT-S tokens. This reveals that the language model can retrieve spatial concept information more effectively from SSM visual representations.

076 SSM visual representations achieve stronger spatial concept
 077 understanding than ViT-family alternatives at matched token
 078 budget, with improvements that transfer from ground-
 079 ing benchmarks to general VQA, leading to strong overall
 080 performance especially at small model scale (Figure 1, left)
 081 (ii) token-region similarity maps computed at intermediate
 082 layers of the language model show that SSM visual tokens
 083 produce sharper and more spatially selective concept-region
 084 binding, providing direct interpretability evidence that the
 085 language model can use SSM representations more effective-
 086 ly; (iii) classification pretraining is misaligned with spatial
 087 concept retention, and dense pretraining objectives partial-
 088 ly recover spatial structure; and (iv) grounding failures
 089 in some configurations stem from the vision-language inter-
 090 face rather than the visual representation itself, and can
 091 be fixed with simple architecture-agnostic stabilization. In
 092 our experiments, we also observe that localization bench-
 093 mark scores correlate strongly with general VQA scores
 094 further supporting this connection. This suggests that im-
 095 proving the spatial concept quality of visual representations
 096 is a general strategy for improving VLMs, not only a fix for
 097 grounding-specific tasks. Our contributions are:

- 098 • We identify compact and structured visual token repre-
 099 sentations as a key axis for improving spatial concept un-
 100 derstanding in VLMs, and argue that architectural induc-
 101 tive bias is an underexplored direction compared to scal-
 102 ing token counts or resolution.
- 103 • We provide the first controlled evaluation of SSM-based
 104 vision encoders as frozen visual backbones in VLMs,

showing that VMamba-based VLMs achieve stronger spatial
 concept understanding than ViT-family alternatives at matched
 scale and token budget.

- We provide interpretability evidence through token-
 region similarity maps at intermediate LLM layers, showing
 that the language model can leverage SSM visual tokens
 more effectively for spatial concept reasoning.
- We identify and diagnose two failure modes: classifica-
 tion pretraining eroding spatial concept structure, and
 vision-language interface bottlenecks causing localiza-
 tion collapse, and propose simple stabilizations for both.

2. Preliminaries

We adopt a LLaVA-style [18] VLM architecture consisting
 of a frozen vision encoder, a lightweight two-layer MLP
 connector, and a decoder-only language model (Vicuna-
 7B). The vision encoder is kept frozen during training; only
 the connector and language model are updated. This design
 allows us to swap vision backbones while keeping all other
 components fixed, enabling a controlled comparison of ar-
 chitectural effects. We evaluate six vision encoder families:
 VMamba [19], ViT [6], MaxViT [26], MambaVision [8],
 ViTDet [16], and DeiT [25]. A brief description of each
 architecture is provided in the supplementary material.

We follow the training recipe of [12], using a one-stage
 instruction tuning setup on 665K multimodal instruction-
 tuning examples [18]. Full training details are provided in
 the supplementary material.

We evaluate on two groups of benchmarks. For general VQA, we use VQA-v2 [7], GQA [10], VizWiz [2], TextVQA [23], POPE [17], and TallyQA [1]. For spatial grounding and localization, we use RefCOCO, RefCOCO+, RefCOCOg [13], and OCID-Ref [28]. We report weighted average scores across each group, as well as an overall weighted average across all benchmarks.

3. SSM Visual Representations Improve VLMs

3.1. Matched Token Budget Comparison

We compare vision encoders from four architecture families under a strictly matched setting: all backbones are pre-trained on ImageNet-1K at 224×224 resolution and produce 196 visual tokens. The four families are VMamba [19] (pure SSM), ViT [6] (plain transformer), MaxViT [26] (hybrid convolution and attention), and MambaVision [8] (hybrid SSM and transformer). This setting isolates the effect of backbone architecture from all other factors.

Table 1 reports the weighted average VQA and localization scores for each backbone. VMamba-T/S achieves the strongest performance across both VQA and localization, with a particularly large margin on localization over all other families. Notably, higher ImageNet accuracy does not predict better VLM performance: backbones with the highest classification accuracy consistently underperform VMamba, suggesting that classification-task accuracy is orthogonal to spatial concept quality.

Table 1. **Matched IN1K/224 backbone swaps.** All backbones use 224×224 inputs and $L=196$ visual tokens. IN1K denotes ImageNet-1K top-1 accuracy. We report weighted average VQA, localization (Loc.), and overall scores. Best results in bold, second best underlined.

Visual Encoder	Size	IN1K	VQA	Loc.	Overall
ViT-S	22M	78.8	57.25	17.82	51.95
ViT-B	87M	81.1	57.17	13.92	51.36
MaxViT-T	31M	83.4	58.75	15.79	52.98
MaxViT-S	69M	84.4	57.42	13.32	51.49
MaxViT-B	119M	84.9	57.60	11.41	51.39
MaxViT-L	212M	84.9	57.30	10.81	51.05
MambaVision-T	32M	82.3	54.38	20.98	49.89
MambaVision-S	50M	83.3	55.61	28.22	51.93
MambaVision-B	98M	84.2	56.53	31.17	53.12
MambaVision-L	228M	85.0	56.94	31.51	53.52
VMamba-T	30M	82.6	<u>62.07</u>	39.20	<u>59.00</u>
VMamba-S	50M	83.6	62.39	<u>39.17</u>	59.27
VMamba-B	80M	83.9	61.38	27.89	56.88

3.2. Interpretability: How LLMs Use Visual Tokens

Figure 1 (right, bottom) shows token-region similarity maps computed between visual tokens and text concept tokens at

an intermediate layer of the language model, reflecting both the spatial structure encoded in the visual tokens and the degree to which it is preserved through the language model’s intermediate representations.

VMamba tokens produce tight and spatially selective concept-region binding: for example, the concept “Mouth” activates a small, precise facial region, and “zebra” activates the correct animal with minimal spread to surrounding regions. In contrast, ViT tokens produce more diffuse responses, with high similarity spread across multiple regions or objects. This pattern is consistent across all examples in , and aligns with the grounding predictions in Figure 1 (right), where VMamba-based VLMs predict tighter and more accurate bounding boxes than their ViT counterparts. Together, these results suggest that SSM visual tokens carry denser spatial concept structure per token, and that the language model is better able to ground referenced concepts when using SSM representations.

Furthermore, across all evaluated configurations, localization benchmark scores correlate strongly with general VQA scores (Pearson r between 0.65 and 0.80 for VQA-v2, GQA, POPE, and TallyQA), confirming that improvements in spatial concept structure propagate broadly beyond explicit grounding tasks.

4. What Determines Spatial Concept Quality?

4.1. Architectural Inductive Bias

Why does VMamba produce better-structured visual tokens under identical training conditions? VMamba’s SS2D mechanism applies state-space aggregation in four scan directions across the 2D token grid, so each token is built by integrating spatial context from its neighbors along rows and columns. Spatial relational structure is therefore embedded into every token by construction. In contrast, ViT’s self-attention is permutation-invariant: spatial structure is carried only by positional encodings, which standard classification pretraining does not explicitly reward the model to utilize [11, 27]. As a result, SSM tokens are intrinsically more compact spatial concept carriers per unit of token budget, which explains the sharper concept-region binding observed in Figure 1 (right, bottom).

4.2. Effect of Pretraining Objectives

Classification pretraining optimizes for global category identity and creates a representation bottleneck: models are rewarded for retaining information that distinguishes between object categories, while spatial layout information that is irrelevant for predicting the image label is discarded. This explains why higher ImageNet accuracy does not predict better VLM performance, as seen in Table 1. Adapting backbones with dense pretraining objectives such as detection or segmentation partially recovers spatial concept

211 fidelity by directly supervising spatial layouts and region-
 212 level discrimination. As shown in Table 2, segmentation-
 213 adapted VMamba variants consistently outperform DeiT
 214 baselines of comparable size on both VQA and localization.
 215 For detection-adapted backbones, naive application of these
 216 checkpoints causes localization collapse for both ViTDet
 217 and VMamba variants (highlighted in blue), which we ad-
 218 dress in Section 4.3; after stabilization, VMamba achieves
 219 the strongest overall performance at a much smaller model
 220 scale than ViTDet. The improvement from dense objectives
 221 is larger for ViT-family backbones, which lack a built-in
 222 spatial prior, than for VMamba, which already preserves
 223 spatial structure architecturally, confirming that inductive
 224 bias and pretraining objective are complementary levers.

Table 2. **Effect of dense pretraining objectives and interface stabilizations.** We report weighted average VQA, localization (Loc.), and overall scores. Blue rows are collapsed configurations. (f) denotes a stronger 3-layer MLP connector. † denotes square 512×512 input geometry.

Visual Encoder	Size	Obj.	VQA	Loc.	Overall
<i>Detection-adapted (IN1K → COCO)</i>					
ViTDet-B	111M	Det	63.00	43.74	60.42
ViTDet-L	331M	Det	57.64	13.05	51.65
+ (f)	331M	Det	63.55	44.58	61.00
VMamba-S	50M	Det	62.78	47.94	60.78
VMamba-T	30M	Det	58.05	14.86	52.25
+ †	30M	Det	64.22	44.52	61.57
+ (f) + †	30M	Det	64.28	46.75	61.92
VMamba-B	89M	Det	58.57	15.02	52.72
+ †	89M	Det	63.58	45.63	61.17
+ (f) + †	89M	Det	63.58	46.90	61.34
<i>Segmentation-adapted (IN1K → ADE20K)</i>					
DeiT-S	58M	Seg	59.36	31.78	55.65
DeiT-B	134M	Seg	60.69	33.77	57.07
VMamba-T	30M	Seg	62.60	43.47	60.03
VMamba-S	50M	Seg	63.21	44.98	60.76
VMamba-B	89M	Seg	62.87	45.08	60.48

225 4.3. Vision-Language Interface Bottlenecks

226 Even when a backbone produces high-quality spatial repre-
 227 sentations, the language model may still fail to use them
 228 if the vision-language interface does not transmit them
 229 faithfully. Table 2 shows that naively applying detection-
 230 pretrained checkpoints at their original pretraining reso-
 231 lution (1024×1024 for ViTDet, 1333×800 for VMamba-
 232 Det) leads to localization collapse in several configurations:
 233 ViTDet-L drops to 13.05 and VMamba-T/B drop to around
 234 15, far below their stable counterparts ViTDet-B (43.74)
 235 and VMamba-S (47.94). Since the spatial concept informa-
 236 tion is clearly present in these backbone families, the fail-

237 ure may lie at the interface rather than in the representation
 238 itself. We identify two separable bottlenecks: a transmis-
 239 sion bottleneck, where connector capacity is insufficient to
 240 preserve spatial structure when projecting into the language
 241 model embedding space, and a utilization bottleneck, where
 242 the language model cannot reliably interpret spatial cues
 243 from non-square input geometries. Switching to a square
 244 512×512 input geometry eliminates collapse for VMamba-
 245 T and VMamba-B, recovering localization scores from be-
 246 low 15 to above 44. Increasing connector depth provides ad-
 247 ditional gains and is complementary to the geometry fix, as
 248 seen in the VMamba-T(f)† and VMamba-B(f)† rows. Both
 249 fixes are architecture-agnostic and add no additional param-
 250 eters to the vision encoder.

251 Together, these results show that spatial concept quality
 252 is jointly determined by three separable and independently
 253 addressable factors: architectural inductive bias, pretraining
 254 objective alignment, and vision-language interface fidelity.

255 5. Related Work

256 **Vision backbone comparisons in VLMs.** Several works
 257 study the effect of vision encoder choice in VLMs [5, 12,
 258 24]. However, these comparisons often change multiple fac-
 259 tors simultaneously, including pretraining data, connector
 260 design, and training recipe, making it difficult to isolate the
 261 contribution of backbone architecture. We address this by
 262 conducting a controlled backbone swap under a fixed train-
 263 ing recipe.

264 **SSMs in vision-language models.** Prior work applying
 265 SSMs to vision-language modeling has focused on the lan-
 266 guage and fusion side, using Mamba-style layers to re-
 267 place transformer-based text backbones or multimodal fu-
 268 sion blocks [20, 30], or exploring fully SSM-based con-
 269 trastive pretraining [22]. None of these works evaluate SSM
 270 architectures as frozen vision encoders in generative VLMs
 271 under controlled conditions.

272 6. Conclusion

273 We have shown that spatial concept understanding in vision-
 274 language models is not only a question of how many visual
 275 tokens to use, but also of how much spatial concept struc-
 276 ture each token carries. SSM-based vision encoders pro-
 277 duce visual tokens with richer spatial concept structure per
 278 token than ViT-family alternatives at matched scale and to-
 279 ken budget, with improvements that transfer from ground-
 280 ing to general VQA. Classification pretraining is misaligned
 281 with spatial concept retention, dense objectives partially re-
 282 cover this structure, and grounding failures in some config-
 283 urations are localized to the vision-language interface and
 284 are correctable with simple fixes. We hope this work en-
 285 courages further exploration of SSM-based vision encoders
 286 as a compact and structured alternative to transformer-based
 287 backbones.

288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343

References

- [1] Manoj Acharya, Kushal Kafle, and Christopher Kanan. Tal-lyqa: Answering complex counting questions. In *Proceedings of the AAAI conference on artificial intelligence*, pages 8076–8084, 2019. 3
- [2] Jeffrey P Bigham, Chandrika Jayant, Hanjie Ji, Greg Little, Andrew Miller, Robert C Miller, Robin Miller, Aubrey Tatarowicz, Brandyn White, Samuel White, et al. Vizwiz: nearly real-time answers to visual questions. In *Proceedings of the 23rd annual ACM symposium on User interface software and technology*, pages 333–342, 2010. 3
- [3] Zhe Chen, Yuchen Duan, Wenhai Wang, Junjun He, Tong Lu, Jifeng Dai, and Yu Qiao. Vision transformer adapter for dense predictions. *arXiv preprint arXiv:2205.08534*, 2022. 1
- [4] An-Chieh Cheng, Hongxu Yin, Yang Fu, Qiushan Guo, Ruihan Yang, Jan Kautz, Xiaolong Wang, and Sifei Liu. Spatialrgpt: Grounded spatial reasoning in vision-language models. *Advances in Neural Information Processing Systems*, 37:135062–135093, 2024. 1
- [5] Federico Cocchi, Nicholas Moratelli, Davide Caffagni, Sara Sarto, Lorenzo Baraldi, Marcella Cornia, and Rita Cucchiara. Llava-more: A comparative study of llms and visual backbones for enhanced visual instruction tuning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4278–4288, 2025. 4
- [6] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020. 2, 3, 1
- [7] Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. Making the v in vqa matter: Elevating the role of image understanding in visual question answering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6904–6913, 2017. 3
- [8] Ali Hatamizadeh and Jan Kautz. Mambavision: A hybrid mamba-transformer vision backbone. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 25261–25270, 2025. 2, 3, 1
- [9] Yihong Huang, Fei Ma, Yihua Shao, Jingcai Guo, Zitong Yu, Laizhong Cui, and Qi Tian. N\” uwa: Mending the spatial integrity torn by vlm token pruning. *arXiv preprint arXiv:2602.02951*, 2026. 1
- [10] Drew A Hudson and Christopher D Manning. Gqa: A new dataset for real-world visual reasoning and compositional question answering. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 6700–6709, 2019. 1, 3
- [11] Samy Jelassi, Michael Sander, and Yuanzhi Li. Vision transformers provably learn spatial structure. *Advances in Neural Information Processing Systems*, 35:37822–37836, 2022. 3
- [12] Siddharth Karamcheti, Suraj Nair, Ashwin Balakrishna, Percy Liang, Thomas Kollar, and Dorsa Sadigh. Prismatic vlms: Investigating the design space of visually-conditioned language models. In *International Conference on Machine Learning (ICML)*, 2024. 1, 2, 4
- [13] Sahar Kazemzadeh, Vicente Ordonez, Mark Matten, and Tamara Berg. Referitgame: Referring to objects in photographs of natural scenes. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, pages 787–798, 2014. 3
- [14] Aisha Urooj Khan, Hilde Kuehne, Chuang Gan, Niels Da Victoria Lobo, and Mubarak Shah. Weakly supervised grounding for vqa in vision-language transformers. In *European Conference on Computer Vision*, pages 652–670. Springer, 2022. 1
- [15] Kevin Li, Sachin Goyal, João D. Semedo, and J Zico Kolter. Inference optimal VLMS need fewer visual tokens and more parameters. In *The Thirteenth International Conference on Learning Representations*, 2025. 1
- [16] Yanghao Li, Hanzi Mao, Ross B. Girshick, and Kaiming He. Exploring plain vision transformer backbones for object detection. *ArXiv*, abs/2203.16527, 2022. 2, 1
- [17] Yifan Li, Yifan Du, Kun Zhou, Jinpeng Wang, Wayne Xin Zhao, and Ji-Rong Wen. Evaluating object hallucination in large vision-language models. In *Proceedings of the 2023 conference on empirical methods in natural language processing*, pages 292–305, 2023. 3
- [18] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *Advances in neural information processing systems*, 36:34892–34916, 2023. 1, 2
- [19] Yue Liu, Yunjie Tian, Yuzhong Zhao, Hongtian Yu, Lingxi Xie, Yaowei Wang, Qixiang Ye, Jianbin Jiao, and Yunfan Liu. Vmamba: Visual state space model. *Advances in neural information processing systems*, 37:103031–103063, 2024. 1, 2, 3
- [20] Yanyuan Qiao, Zheng Yu, Longteng Guo, Sihan Chen, Zijia Zhao, Mingzhen Sun, Qi Wu, and Jing Liu. VI-mamba: Exploring state space models for multimodal learning. *arXiv preprint arXiv:2403.13600*, 2024. 4
- [21] Kanchana Ranasinghe, Satya Narayan Shukla, Omid Pour-saeed, Michael S Ryoo, and Tsung-Yu Lin. Learning to localize objects improves spatial reasoning in visual-llms. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12977–12987, 2024. 1
- [22] Nimrod Shabtay, Itamar Zimmerman, Eli Schwartz, and Raja Giryes. Climp: Contrastive language-image mamba pretraining. *arXiv preprint arXiv:2601.06891*, 2026. 4
- [23] Amanpreet Singh, Vivek Natarajan, Meet Shah, Yu Jiang, Xinlei Chen, Dhruv Batra, Devi Parikh, and Marcus Rohrbach. Towards vqa models that can read. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8317–8326, 2019. 3
- [24] Peter Tong, Ellis Brown, Penghao Wu, Sanghyun Woo, Adithya Jairam Vedagiri Iyer, Sai Charitha Akula, Shusheng Yang, Jihan Yang, Manoj Middepogu, Ziteng Wang, et al. Cambrian-1: A fully open, vision-centric exploration of multimodal llms. *Advances in Neural Information Processing Systems*, 37:87310–87356, 2024. 4
- [25] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou. Training

- 401 data-efficient image transformers & distillation through at-
402 tention. *arXiv preprint arXiv:2012.12877*, 2(3),
403 2020. 2, 1
- 404 [26] Zhengzhong Tu, Hossein Talebi, Han Zhang, Feng Yang,
405 Peyman Milanfar, Alan Bovik, and Yinxiao Li. Maxvit:
406 Multi-axis vision transformer. In *European conference on*
407 *computer vision*, pages 459–479. Springer, 2022. 2, 3, 1
- 408 [27] Haochen Wang, Junsong Fan, Yuxi Wang, Kaiyou Song,
409 Tong Wang, and ZHAO-XIANG ZHANG. Droppos: Pre-
410 training vision transformers by reconstructing dropped posi-
411 tions. *Advances in Neural Information Processing Systems*,
412 36:46134–46151, 2023. 3
- 413 [28] Ke-Jyun Wang, Yun-Hsuan Liu, Hung-Ting Su, Jen-Wei
414 Wang, Yu-Siang Wang, Winston Hsu, and Wen-Chin Chen.
415 Ocid-ref: A 3d robotic dataset with embodied language for
416 clutter scene grounding. In *Proceedings of the 2021 Confer-*
417 *ence of the North American Chapter of the Association for*
418 *Computational Linguistics: Human Language Technologies*,
419 pages 5333–5338, 2021. 3
- 420 [29] Weiyun Wang, Zhangwei Gao, Lixin Gu, Hengjun Pu, Long
421 Cui, Xingguang Wei, Zhaoyang Liu, Linglin Jing, Shenglong
422 Ye, Jie Shao, et al. Internv13. 5: Advancing open-source
423 multimodal models in versatility, reasoning, and efficiency.
424 *arXiv preprint arXiv:2508.18265*, 2025. 1
- 425 [30] Han Zhao, Min Zhang, Wei Zhao, Pengxiang Ding, Siteng
426 Huang, and Donglin Wang. Cobra: Extending mamba to
427 multi-modal large language model for efficient inference.
428 In *Proceedings of the AAAI Conference on Artificial Intel-*
429 *ligence*, pages 10421–10429, 2025. 4

Toward Compact and Structured Visual Representations in VLMs: SSM-Based Vision Encoders as an Alternative to Transformers

Supplementary Material

430 7. Training Details

431 **Optimization.** We use Fully Sharded Data Parallel
432 (FSDP) training on $4\times$ NVIDIA H200 GPUs with BF16
433 mixed precision and activation checkpointing for the lan-
434 guage model. Training is run for 1 epoch with global batch
435 size 128 and per-GPU batch size 16, corresponding to 2 gra-
436 dient accumulation steps. We use AdamW with learning
437 rate 2×10^{-5} , weight decay 0.1, max gradient norm 1.0,
438 and a linear-warmup cosine-decay schedule with warmup
439 ratio 0.03. These hyperparameters are fixed across all ex-
440 periments.

441 **Data and preprocessing.** We fine-tune on the LLaVA-
442 v1.5 665K instruction-tuning mixture [18], which includes
443 a combination of visual question answering, captioning, and
444 conversation data. Images are processed with letterbox re-
445 sizing, which preserves the original aspect ratio by scaling
446 and padding to the target resolution. Incomplete batches
447 are retained rather than dropped. For variable-length inputs,
448 text is truncated to the model maximum length and images
449 are padded within each batch as needed.

450 **Implementation.** Our implementation is based on the
451 Prismatic-VLMs codebase [12] and uses a fixed batch order
452 across all experiments to ensure reproducibility. One-stage
453 instruction tuning, where the connector is randomly initial-
454 ized and jointly trained with the language model from the
455 start, is used throughout following [12].

456 8. Vision Encoder Architectures

457 **VMamba.** VMamba [19] is a pure SSM-based vision
458 backbone built around a 2D Selective Scan (SS2D) mecha-
459 nism. Rather than applying self-attention over a flattened
460 token sequence, SS2D aggregates spatial context in four
461 scan directions across the 2D token grid, embedding spa-
462 tial relational structure into every token by construction. We
463 use VMamba as our primary SSM backbone given its strong
464 performance on dense vision tasks.

465 **ViT.** ViT [6] tokenizes an image into fixed-size patches
466 and applies global self-attention over the full patch se-
467 quence. Spatial structure is carried by positional encod-
468 ings rather than the attention mechanism itself, making it
469 permutation-invariant by design.

MaxViT. MaxViT [26] is a hierarchical hybrid backbone
that combines convolutions with multi-axis attention, ap-
plying both blocked local attention and dilated global atten-
tion at each stage to capture local and global spatial interac-
tions simultaneously.

MambaVision. MambaVision [8] is a hybrid Mamba-
Transformer backbone that uses Mamba blocks in early
stages for efficient spatial aggregation and retains self-
attention in the final layers to capture long-range dependen-
cies.

ViTDet. ViTDet [16] adapts a plain ViT backbone for ob-
ject detection by showing that a simple feature pyramid
built from a single-scale feature map suffices for detec-
tion fine-tuning. We use ViTDet checkpoints pretrained
on ImageNet-1K and fine-tuned for detection on COCO
at 1024×1024 resolution as our ViT-based dense-objective
baseline.

DeiT. DeiT [25] is a plain ViT backbone trained with im-
proved data augmentation and distillation. We use DeiT
checkpoints adapted with the ViT-Adapter framework [3]
for segmentation on ADE20K at 512×512 resolution as our
ViT-based segmentation baseline.

470
471
472
473
474475
476
477
478
479480
481
482
483
484
485
486487
488
489
490
491