

# OpenVision 2 : A Family of Generative Pretrained Visual Encoders for Multimodal Learning

Yanqing Liu<sup>1</sup> Xianhang Li<sup>1</sup> Letian Zhang<sup>1</sup> Zirui Wang<sup>2</sup> Zeyu Zheng<sup>3</sup> Yuyin Zhou<sup>1</sup> Cihang Xie<sup>1</sup>

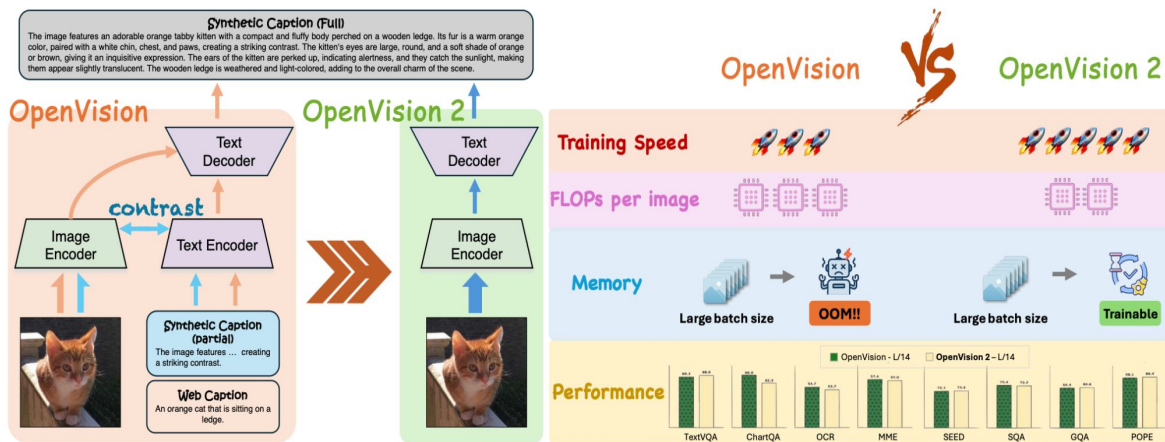
<sup>1</sup>University of California Santa Cruz    <sup>2</sup>Apple    <sup>3</sup>University of California Berkeley

Project Page: <https://ucsc-vlaa.github.io/OpenVision2>

Model Training: <https://github.com/UCSC-VLAA/OpenVision>

# Contributions

- Challenges the belief that **CLIP-style contrastive learning is essential** for vision encoders.
- OpenVision 2 shows a **caption-only generative objective can match multimodal performance**.
- Approach **reduces computation and memory costs** compared to contrastive methods.
- Full training suite and pretrained checkpoints of **OpenVision 2 are publicly released**.



# OpenVision

## Fully-Open Vision Encoders

- Open release of datasets, training recipes, and model checkpoints for transparency and reproducibility.

## Wide Range of Model Scales

- A family of encoders ranging from *Tiny* (~5.9 million params) to *Huge* (~632.1 million).
- Flexibility for deployment across a spectrum from edge devices to big compute servers.

## Superior Multimodal Performance

- Matches or exceeds performance of proprietary encoders (e.g. OpenAI's CLIP, Google's SigLIP) on several multimodal benchmarks, particularly in frameworks like LLaVA-1.5 and Open-LLaVA-Next.

## Efficiency: Progressive & Resolution Training

- Use of progressive resolution training (start with lower resolution images, move to higher) to speed up training and save compute.
- Significant reductions in training time and memory usage in comparison with existing large models and proprietary CLIP models.

Model	Data	Training	Evaluation	Model Num.	Training Time
OpenAI's CLIP	Closed	Closed	Open	4	
Google's SigLIP	Closed	Open	Open	10	
OpenVision	Open	Open	Open	>25	

# OpenVision

## Efficiency (CLIPA @UCSC)

- OpenVision adopts this two-stage curriculum by training on low-resolution images( $84^2$ )and conducting a fine-tuning at full resolution( $224^2$ ).

## Data Quality (Recap @UCSC)

- a LLaMA-3-powered LLaVA model recaptions the entire DataComp-1B collection; this high-quality synthetic set serves as the training corpus for OpenVision.

## Optimization (CLIPS @UCSC)

- To better leverage synthetic captions, CLIPS introduces two additional objectives: (i) **a dual contrastive loss** that pairs each image with both web-crawled and generated captions, and (ii) **a caption loss** that asks the model to predict the synthetic caption given the image and its web caption. OpenVision integrates both losses to enhance training.

# OpenVision2

## OpenVision Challenges

- the text encoder must process **two captions per image for the dual contrastive objective**
- **an additional text decoder is required** to autoregressively predict the synthetic caption.

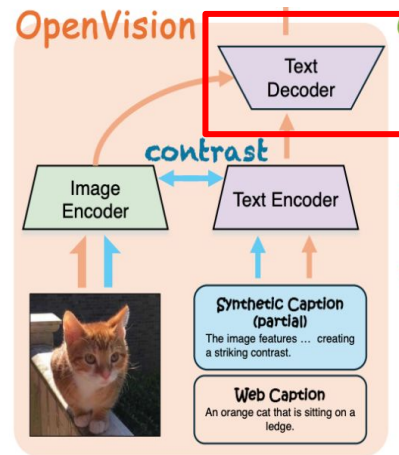
→ Together, these two components substantially increase FLOPs and GPU memory in training.

## OpenVision 2 approaches

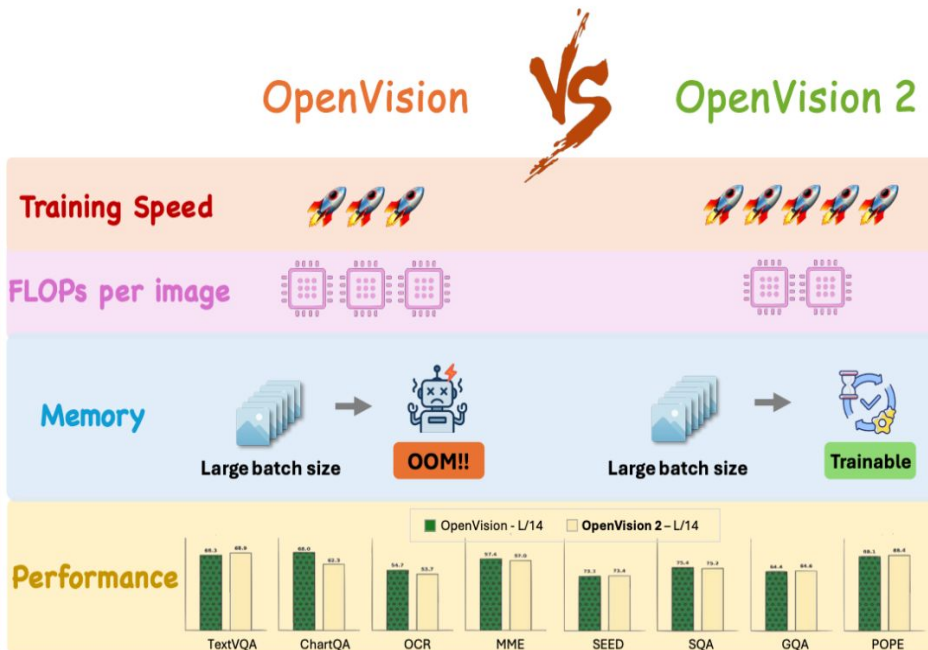
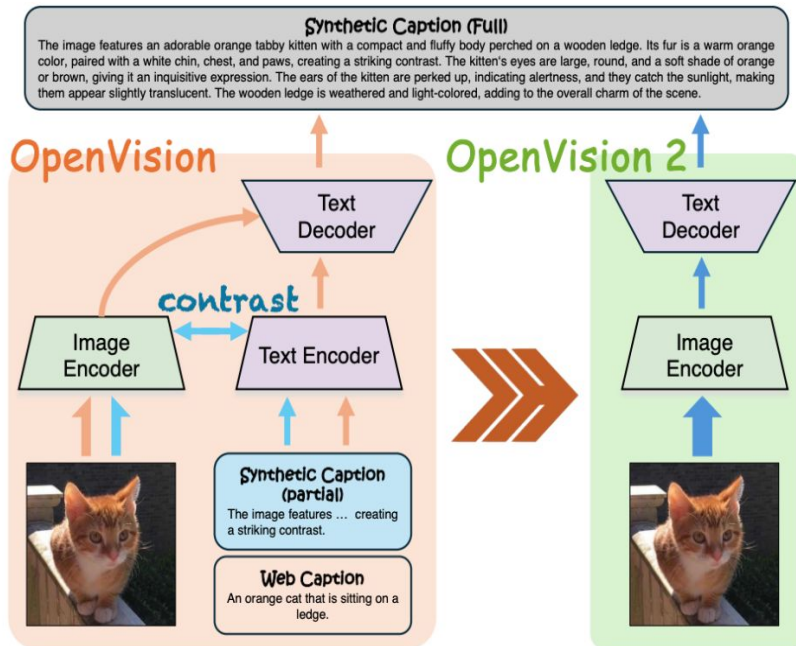
- Discarding the text encoder and contrastive loss.
- Simplifying training to:
  1. Vision encoder → visual tokens.
  2. Text decoder → synthetic caption.

→ This makes pretraining **purely generative**, aligning better with downstream fine-tuning (e.g., LLaVA).

→ Efficiency tweak: randomly mask  $\sim \frac{2}{3}$  of visual tokens, which still allows good captioning while reducing computation.



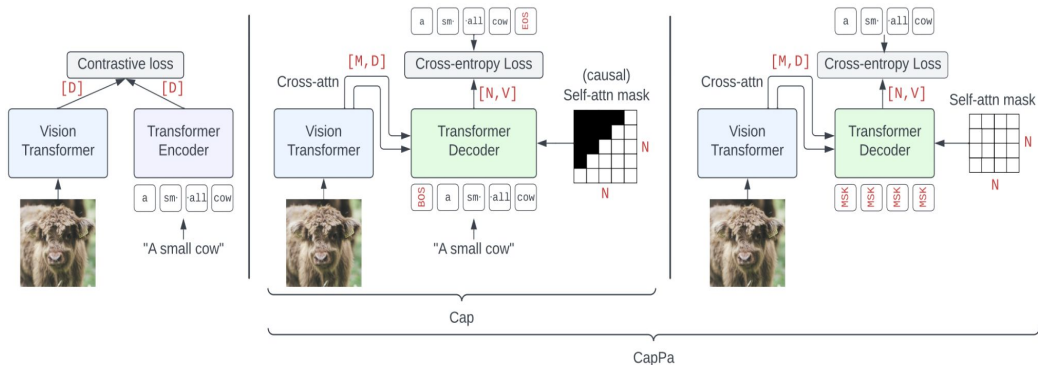
# OpenVision2



# CapPa

## Background

- **Contrastive Pretraining Dominance**
- **Generative Captioning Considered Inferior**
- **Lack of Fair Comparisons**



## Trade-offs

- Zero-shot classification: **contrastive wins in many standard benchmarks.**
- Fine-grained tasks, compositionality, ordering, relations: contrastive models tend to ignore word order or relational structure and treat text more like a “bag of words.” **Captioning models are potentially better at capturing these finer structures.**
- Efficiency / inference cost: captioning models (encoder-decoder) require decoding (autoregressive or parallel), which is costlier in some settings compared to just encoding text/image separately (as in CLIP).

# CapPa

## Architecture

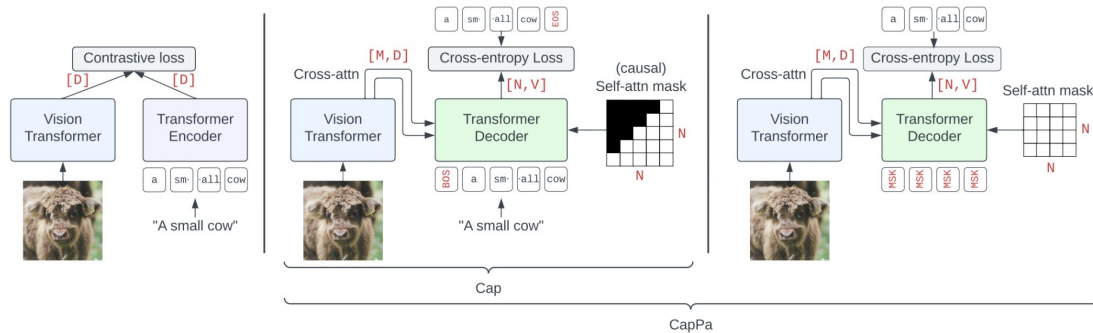
- Vision Transformer (ViT) as the image encoder.
- Standard Transformer decoder that takes the encoder's output via cross-attention to generate captions.
- The decoder has fewer layers (half the depth) than the encoder in their setups, but matches width & attention heads.

## Captioner (Cap) Variant

- Pure image captioning: autoregressive decoding (teacher forcing) predicting next token given previous text tokens + image encoding.

## Parallel Prediction / Mixed Mode (“CapPa”)

- CapPa uses a *parallel prediction* mode for a fraction of training data: the decoder's input text tokens are masked (all [MASK]), attention mask is changed so there's no causal masking. The decoder must predict all tokens at once (positions matter) given only the image (not previous text tokens).





# Difference from CapPa

## Higher-quality captions

- Uses *ReCap-DataComp-1B* (Llama-3 recaptioned dataset) with improved captioning strategy → produces longer, more grounded captions for stronger generative supervision.

## Fusion simplification

- Replaces CapPa's cross-attention with simple concatenation of visual tokens in the text decoder; randomly drops tokens during training to regularize and reduce cost.

## Scale & evaluation

- Scales vision encoder to **1.01B parameters** trained on **12.8B image–caption pairs**; evaluates on advanced benchmarks (MME, ChartQA), beyond classification/QA.

## Decoding strategy

- Uses **standard autoregressive decoding** only, instead of CapPa's hybrid approach.

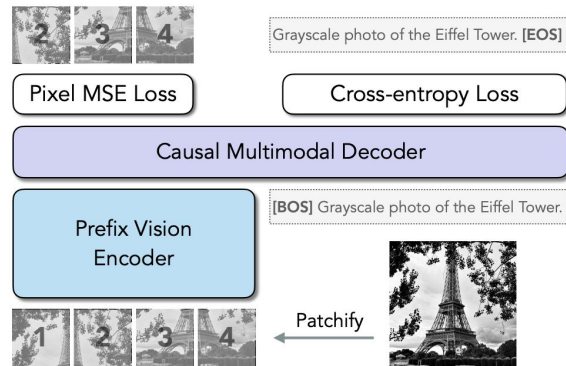
# AIMv2

## Background

- 대표적인 modal align 방법은 generative vs discriminative(contrastive)
- generative 직관적인 사전학습 방법, 그러나 높은 capacity
- discriminative 방법은 parameter efficient, 그러나 학습이 까다로움  
→ generative 사전 학습의 간단함과 확장성 그리고 discriminative 방법의 parameter-efficient 방법을 고려

## Architecture

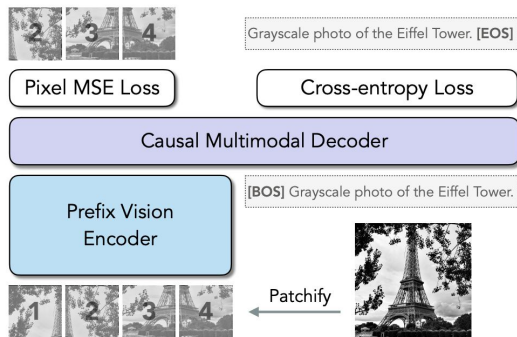
- 구성: vision encoder + multimodal decoder (autoregressive로 다음 이미지 패치와 텍스트 토큰을 예측할 수 있도록)
- ViT architecture (300M~3B)
- Prefix attention (to facilitate the use of bidirectional attention)
- SwiGLU + RMSNorm
- Multimodal decoder
  - The outputs of the decoder are processed through two separate linear heads to predict the next token in each modality respectively. (image token head, text token head)



# AIMv2

## Contributions

- 쉽고 직관적인 사전학습 방법 (이미지 패치와 텍스트 토큰을 같이 **auto-regressively** 예측하는 **causal multimodal encoder** 사용, 이때 **contrastive** 방법과 달리 배치 사이즈, 배치 간 고려는 하지 않기 때문에 학습하기 쉬움)
- 확장성 (다양한 모달 지원)
- **localization, grounding, classification** 포함한 비전 벤치마크, 멀티-모달 벤치마크를 포함 우수한 성능
- 멀티모달 이해에서도 CLIP SigLIP과 같은 **sota** 모델보다 우수함.



모델	특징 및 강점	단점
AIMv2	멀티모달 및 이미지 이해에서 높은 성능, 다양한 해상도 지원	학습 및 사용을 위해 고사양 필요
CLIP	텍스트-이미지 align 작업에서 높은 성능	멀티모달 확장성 제한적
DINOv2	객체탐지에서 우수성능	멀티모달성능은 제한적

# Difference from AIMv2

## Training signal

- AIMv2 → combines image patch reconstruction + text generation.
- OpenVision 2 → caption-only supervision (no image reconstruction).

## Token masking

- OpenVision 2 masks  $\sim\frac{2}{3}$  of visual tokens → improves efficiency & performance.

## Data composition

- AIMv2 → mix of human (67%) + synthetic (33%) captions.
- OpenVision 2 → fully synthetic captions from ReCap-DataComp-1B (richer & consistent).

## Vision encoder

- AIMv2 → prefixViT with special attention mask.
- OpenVision 2 → standard ViT backbone (simple & efficient).

# Results

Under LLaVA 1.5 Framework

Method	Vision Encoder	Params	# Res.	Text VQA	Chart QA	OCR.	MME	SEED	SQA	GQA	POPE
OpenAI-CLIP [44]	L/14	304M	224	56.1	13.2	177	1443/306	66.0	73.4	60.8	85.0
LAION-2B-CLIP [19]	L/14	304M	224	54.2	12.8	165	1434/298	65.5	76.0	59.0	84.5
DataComp-1B-CLIP [16]	L/14	304M	224	53.0	12.3	131	1382/312	62.4	74.2	57.8	83.0
DFN-2B-CLIP [12]	L/14	304M	224	53.2	12.4	246	1447/306	65.6	76.3	59.1	85.0
MetaCLIP-5B [59]	L/14	304M	224	55.6	12.8	313	1552/315	67.4	78.0	61.3	85.4
OpenVision [30]	L/14	304M	224	57.7	13.9	315	1487/317	69.5	73.6	62.9	86.4
<b>OpenVision 2</b>	<b>L/14</b>	<b>304M</b>	<b>224</b>	<b>59.0</b>	<b>13.7</b>	<b>327</b>	<b>1460/312</b>	<b>69.3</b>	<b>76.5</b>	<b>62.6</b>	<b>87.1</b>
OpenAI-CLIP [44]	L/14	304M	336	59.1	13.8	201	1475/288	67.5	73.1	61.1	85.7
OpenVision [30]	L/14	304M	336	61.2	15.7	339	1525/315	70.5	75.1	63.7	87.2
<b>OpenVision 2</b>	<b>L/14</b>	<b>304M</b>	<b>336</b>	<b>63.0</b>	<b>14.5</b>	<b>357</b>	<b>1486/321</b>	<b>70.1</b>	<b>77.5</b>	<b>63.0</b>	<b>87.7</b>
SigLIP [62]	SoViT-400M/14	400M	384	62.6	14.5	338	1481/347	69.4	76.7	63.3	87.0
OpenVision [30]	SoViT-400M/14	400M	384	62.4	16.1	357	1493/320	70.4	72.4	63.8	88.0
<b>OpenVision 2</b>	<b>SoViT-400M/14</b>	<b>400M</b>	<b>384</b>	<b>64.3</b>	<b>15.0</b>	<b>387</b>	<b>1472/310</b>	<b>70.7</b>	<b>74.9</b>	<b>63.5</b>	<b>87.5</b>
<b>OpenVision 2</b>	<b>H/14</b>	<b>632M</b>	<b>224</b>	<b>60.2</b>	<b>13.5</b>	<b>340</b>	<b>1470/305</b>	<b>69.3</b>	<b>75.4</b>	<b>62.5</b>	<b>87.2</b>
<b>OpenVision 2</b>	<b>H/14</b>	<b>632M</b>	<b>336</b>	<b>63.4</b>	<b>16.3</b>	<b>391</b>	<b>1470/311</b>	<b>70.6</b>	<b>76.4</b>	<b>63.1</b>	<b>88.4</b>
<b>OpenVision 2</b>	<b>H/14</b>	<b>632M</b>	<b>448</b>	<b>65.6</b>	<b>18.1</b>	<b>416</b>	<b>1499/331</b>	<b>70.6</b>	<b>75.6</b>	<b>63.1</b>	<b>88.7</b>
<b>OpenVision 2</b>	<b>g/14</b>	<b>1.01B</b>	<b>224</b>	<b>60.2</b>	<b>13.7</b>	<b>338</b>	<b>1469/290</b>	<b>69.3</b>	<b>75.0</b>	<b>62.6</b>	<b>86.9</b>

# Results

Under Open-LLaVA next Framework

Method	Vision Encoder	Params	# Res.	Text VQA	Chart QA	OCR.	MME	SEED	SQA	GQA	POPE
OpenAI-CLIP [44]	L/14	304M	224	62.8	60.7	459	1600/334	70.6	75.0	62.8	86.9
LAION-2B-CLIP [19]	L/14	304M	224	59.4	50.8	396	1533/323	70.0	72.9	62.7	86.4
DataComp-1B-CLIP [16]	L/14	304M	224	58.1	48.5	373	1524/348	70.2	75.6	62.3	86.2
DFN-2B-CLIP [12]	L/14	304M	224	57.0	42.7	303	1486/328	68.3	70.6	61.7	86.0
MetaCLIP-5B [59]	L/14	304M	224	63.0	62.9	493	1590/335	72.3	77.1	64.0	86.8
OpenVision	L/14	304M	224	65.7	61.5	503	1567/332	73.1	73.1	64.7	87.8
<b>OpenVision 2</b>	<b>L/14</b>	<b>304M</b>	<b>224</b>	<b>66.1</b>	<b>60.4</b>	<b>501</b>	<b>1577/297</b>	<b>73.1</b>	<b>68.4</b>	<b>64.6</b>	<b>87.6</b>
OpenAI-CLIP [44]	L/14	304M	336	69.4	70.0	535	1591/351	73.3	76.9	64.5	87.6
OpenVision	L/14	304M	336	68.3	68.0	547	1520/310	73.3	75.4	64.4	88.1
<b>OpenVision 2</b>	<b>L/14</b>	<b>304M</b>	<b>336</b>	<b>68.9</b>	<b>62.3</b>	<b>537</b>	<b>1585/278</b>	<b>73.4</b>	<b>75.2</b>	<b>64.6</b>	<b>88.4</b>
SigLIP [62]	SoViT-400M/14	400M	384	68.2	61.3	494	1539/325	72.9	74.7	62.9	86.8
OpenVision	SoViT-400M/14	400M	384	67.4	63.1	540	1500/353	72.2	73.5	63.4	87.8
<b>OpenVision 2</b>	<b>SoViT-400M/14</b>	<b>400M</b>	<b>384</b>	<b>69.0</b>	<b>63.4</b>	<b>549</b>	<b>1521/319</b>	<b>72.2</b>	<b>72.7</b>	<b>63.1</b>	<b>87.7</b>
<b>OpenVision 2</b>	<b>H/14</b>	<b>632M</b>	<b>224</b>	<b>66.4</b>	<b>60.2</b>	<b>514</b>	<b>1597/314</b>	<b>73.3</b>	<b>76.2</b>	<b>64.7</b>	<b>88.4</b>
<b>OpenVision 2</b>	<b>H/14</b>	<b>632M</b>	<b>336</b>	<b>69.9</b>	<b>64.8</b>	<b>573</b>	<b>1572/337</b>	<b>73.8</b>	<b>74.5</b>	<b>64.4</b>	<b>87.8</b>
<b>OpenVision 2</b>	<b>H/14</b>	<b>632M</b>	<b>448</b>	<b>71.9</b>	<b>64.9</b>	<b>590</b>	<b>1542/324</b>	<b>74.1</b>	<b>75.6</b>	<b>64.4</b>	<b>88.8</b>
<b>OpenVision 2</b>	<b>g/14</b>	<b>1.01B</b>	<b>224</b>	<b>67.3</b>	<b>62.4</b>	<b>514</b>	<b>1558/323</b>	<b>73.4</b>	<b>74.4</b>	<b>64.7</b>	<b>88.0</b>

# Results

Model	Backbone	Resolution	v4-512 Hours	FLOPs / Image
OpenVision [30]	L/14	224	83	271.75
<b>OpenVision 2</b>	<b>L/14</b>	<b>224</b>	<b>57</b>	<b>208.90</b>
OpenVision [30]	SoViT-400M/14	384	241	1636.75
<b>OpenVision 2</b>	<b>SoViT-400M/14</b>	<b>384</b>	<b>121</b>	<b>1017.74</b>

OpenVision 2 achieves faster training and lower computational cost across model sizes.

Model	Resolution	Batch Size	Peak Memory (GB)
OpenVision [30] (L/14)	224	2k	24.5
	224	4k	OOM
<b>OpenVision 2 (L/14)</b>	<b>224</b>	<b>2k</b>	<b>13.8</b>
	<b>224</b>	<b>4k</b>	<b>22.1</b>
	<b>224</b>	<b>8k</b>	<b>28.4</b>
OpenVision [30] (SoViT-400M/14)	384	512	27.4
	384	1k	OOM
<b>OpenVision 2 (SoViT-400M/14)</b>	<b>384</b>	<b>512</b>	<b>14.5</b>
	<b>384</b>	<b>1k</b>	<b>28.8</b>

OpenVision 2 achieves faster training and lower computational cost across model sizes.

# Results

Caption Type	Text VQA	Chart QA	OCR.	MME	SEED	SQA	GQA	POPE
Alt-text	51.8	12.3	238	1306/293	58.6	75.3	55.4	82.2
ReCap-DataComp-1B	56.9	12.9	291	1426/293	67.9	74.5	61.9	86.5
ReCap-DataComp-1B v2	56.5	13.1	303	1451/310	67.8	74.7	61.2	86.6

Keep Ratio	Text VQA	Chart QA	OCR.	MME	SEED	SQA	GQA	POPE
100%	53.8	12.2	254	1409/350	65.9	73.9	60.3	84.7
90%	56.3	12.4	266	1461/335	67.6	74.8	61.1	85.4
75%	55.8	13.1	293	1438/283	68.6	73.9	61.7	86.3
50%	55.4	12.8	299	1429/313	68.5	73.8	61.6	86.5
35%	56.9	12.9	291	1426/293	67.9	74.5	61.9	86.5
25%	56.7	12.5	283	1430/297	67.8	76.3	61.4	86.3
10%	55.6	13.0	276	1412/301	66.1	75.0	61.2	85.4

A higher keep ratio retains more vision tokens as captioning conditions, while a lower keep ratio masks more tokens



# Discussion

## **Loss of Contrastive Signal / Alignment Robustness**

- Dropping contrastive image-text and relying only on generative captions weakens alignment robustness, hurting retrieval, zero-shot discrimination, and fine-grained image-text matching when captions are noisy and biased.

## **Reliance on Synthetic Captions**

- Heavy reliance on synthetic captions bakes in their quality, biases, omissions, and style; because they focus on silent, generic content rather than exhaustive scene detail, coverage for rare objects, fine-grained attributes, and complex relationships can suffer

## **Caption-Only Objective Might Miss Non-Descriptive Visual Features**

- A caption-only objective neglects non-descriptive visual cues (e.g., low-level textures, subtle spatial relations, background details) that contrastive learning can capture, potentially reducing performance on tasks needing these features.