Early VLMS

Philipp Krähenbühl, UT Austin

Computer vision

"Traditional" Computer Vision 2014 - now



. . .

Unified Computer Vision 2020 - now

Task



Unified Model Output

CLIP as a VLM

- Clip maps
 - Images to text
 - Text to images
- Primitive image and text models
- No dialogue

(1) Contrastive pre-training

(2) Create dataset classifier from label text





-	
YΤ	
-	

Unified-IO

- Tokenize everything
 - Image, Text, Sparse (few numbers),
 Dense (depth maps etc)
 - Image, Dense: VQ-GAN
 - Text: SentencePiece
 - Sparse: 1000 special tokens (coordinates)

Unified-IO: A Unified Model for Vision, Language, and Multi-Modal Tasks, Lu etal 2022



Tasks

Image Classification Object Detection Semantic Segmentation Depth Estimation Surface Normal Estimation Segment-based Image Generation Image Inpainting Pose Estimation Relationship Detection Image Caption ng Visual QA Referring Expressions Situation Recognition Text-based Image Generation Visual Commonsense Classification ir context Region Captioning GLUE Benchmark tasks Reading comprehension Natural Language Inference



Unified-IO

- Everything is sequence prediction (with 1D or 2D positional embeddings)
- Unsupervised pre-training of vision and language inputs (masked prediction)
- Fine-tune (train) on 95 unified datasets

Unified-IO: A Unified Model for Vision, Language, and Multi-Modal Tasks, Lu etal 2022

	Example		Siz	e		Input Modalities				Output Modalities			
	Source	Datasets	Size	Percent	Rate	Text	Image	Sparse	Dense	Text	Imaga	Sparse	Dense
Image Synthesis		14	56m	43.0	18.7	1	1	~	1	-	1		-
Image Synthesis from Text	RedCaps	9	55m	41.9	16.7	1	-	-	-	-	1	-	-
Image Inpainting	VO	3	1,2m	0.9	1.5	1	1	1	-	-	1	-	-
Image Synthesis from Seg.	LVIS	2	220k	0.2	0.6	1	-	-	×	-	1		-
Sparse Labelling		10	8.2m	6.3	12.5	1	1	 Image: A second s	-	-	-	<	-
Object Detection	Open Images	3	1.9m	1.5	3.6	-	1	-	-	-	-	\checkmark	-
Object Localization	VG	3	6m	4.6	7.1	1	1	-	-	-	-	1	_
Keypoint Estimation	coco	1	140k	0.1	0.7	-	e.	5	-	-	-	\checkmark	-
Referring Expression	RefCoco	3	130k	0.1	1.1	1	1	-	-	-	-	1	-
Dense Labelling		6	2.4m	1.8	6.2	1	1	-	-	-	-		1
Depth Estimation	NYU Depth	1	48k	0.1	0.4	-	1	-	-	-	-		4
Surface Normal Estimation	Framenet	2	210k	0.2	1,1	-	1	-	-	-	-	-	4
Object Segmentation	LVIS	3	$2.1\mathrm{m}$	1.6	4.7	1	1	-	-	-	-	-	
Image Classification		9	22m	16.8	12.5	-	1	1	-	1	-	-	-
Image Classification	ImageNet	6	16m	12.2	8.1	1	1	-	-	√	-	-	-
Object Categorization	COCO	3	бm	4.6	4.4	-	×.	V	-		-		-
Image Captioning		7	31m	23.7	12.5	-	1	 Image: A second s	-	√	-		-
Webly Supervised Captioning	CC12M	3	26m	19.7	8.8	•	<	-	-	4	-	•	-
Supervised Captioning	VizWiz	3	1.4m	1.1	1.7	-	1	-	-	4	-	-	-
Region Captioning	VG	1	3.8m	2.9	2.0	-	1	5	-	1	-	-	-
Vision & Language		16	4m	3.0	12.5	1	1	5	-	- 4	-	-	1
Visual Question Answering	VQA 2.0	13	3.3m	2.5	10.4	√	<	1	-	- √	-	-	-
Relationship Detection	VG	2	640k	0.5	1.9	-	×	× .	-	*	-	-	-
Grounded VQA	VizWiz	1	6.5k	0.1	0.1	1	8	-	-		-	-	1
NLP		31	7.1m	5.4	12.5	1	-	-	-	- √	-	-	-
Text Classification	MALL	17	1.6m	1.2	4.8	1	-	-	-	1	-	-	-
Question Answering	SQHAD	13	1.7m	1.3	5.2		-	-	-	- √	-	-	-
Text Summarization	Gigaword	1	3.8m	2.9	2.5	1		-	-	4	-		-
Language Modelling		2	-	-	12.5	1	-	-	-	4	-		-
Masked Language Modelling	C4	2	-	-	12.5		-	-	-	*	-		-
All Tasks		95	130m	100	100	1	1	 Image: A second s	1	4	1	1	1



Flamingo

- One of first big VLM
- Based on Chinchilla
 - Frozen vision model spliced in (Gated X-Attn)
- Multi-image support
 - Masked attention only to preceding image



Flamingo: a Visual Language Model for Few-Shot Learning, Alayrac etal 2022



Flamingo

- Started zero-shot eval trend
- Good results on VQA
 - Under VQA metrics



Performance relative

N>

Image

- One network for 3 tasks
 - CLIP-style embeddings (ITC)
 - Image-language matching (ITM)
 - Language modeling / Captioning (LM)
- Different structure, shared weights (except SA)

BLIP: Bootstrapping Language-Image Pre-training for Unified Vision-Language Understanding and Generation, Li etal 2022





N×

Image

- Pre-Training
 - Image captioning data (COCO etal)
- CapFilt
 - Use ITM, LM fine-tuned separately on COCO
 - Clean up web-scale image-text data
 - Create new captions (LM)
 - Filter new + original captions (ITM)

BLIP: Bootstrapping Language-Image Pre-training for Unified Vision-Lar



	Pre-train	NoCaps validation									Caption
Method		in-doi	main	near-domain		out-domain		overall		Karpathy test	
	#IIIages	C C	S	С	S	С	S	С	S	B@4	С
Enc-Dec (Changpinyo et al., 2021)	15M	92.6	12.5	88.3	12.1	94.5	11.9	90.2	12.1	-	110.9
VinVL† (Zhang et al., 2021)	5.7M	103.1	14.2	96.1	13.8	88.3	12.1	95.5	13.5	38.2	129.3
$LEMON_{base}$ † (Hu et al., 2021)	12M	104.5	14.6	100.7	14.0	96.7	12.4	100.4	13.8	-	-
$LEMON_{base}$ † (Hu et al., 2021)	200M	107.7	14.7	106.2	14.3	107.9	13.1	106.8	14.1	40.3	133.3
BLIP	14M	111.3	15.1	104.5	14.4	102.4	13.7	105.1	14.4	38.6	129.7
BLIP	129M	109.1	14.8	105.8	14.4	105.7	13.7	106.3	14.3	39.4	131.4
BLIP _{CapFilt-L}	129M	111.8	14.9	108.6	14.8	111.5	14.2	109.6	14.7	39.7	133.3
LEMON _{large} † (Hu et al., 2021)	200M	116.9	15.8	113.3	15.1	111.3	14.0	113.4	15.0	40.6	135.7
$SimVLM_{huge}$ (Wang et al., 2021)	1.8 B	113.7	_	110.9	_	115.2	_	112.2	_	40.6	143.3
BLIP _{ViT-L}	129M	114.9	15.2	112.1	14.9	115.3	14.4	113.2	14.8	40.4	136.7

Table 7. Comparison with state-of-the-art image captioning methods on NoCaps and COCO Caption. All methods optimize the crossentropy loss during finetuning. C: CIDEr, S: SPICE, B@4: BLEU@4. BLIP_{CapFilt-L} is pre-trained on a dataset bootstrapped by captioner and filter with ViT-L. VinVL[†] and LEMON[†] require an object detector pre-trained on 2.5M images with human-annotated bounding boxes and high resolution (800×1333) input images. SimVLM_{huge} uses $13 \times$ more training data and a larger vision backbone than ViT-L.

> answer 1



nd Generation, Li etal 2022

B||PV/

- Captioning -> Language modeling
- Q-Former
 - Image-text embedder / tokenizer
 - Fed as context to LLM





BIJPV2

- Great results
 - Captioning
 - VQA
 - Retrieval
- Early visual dialogues

BLIP-2: Bootstrapping Language-Image Pre-training with Frozen Image Encoders and Large Language Models, Li etal 2023

Models	#Trainable Params	Open- sourced?	Visual Question Answering VQAv2 (test-dev) VQA acc.	Image C NoCaj CIDEr	aptioning ps (val) SPICE	Image-Te Flick TR@1	t Retrieval r (test) IR@1
BLIP (Li et al., 2022)	583M	1	-	113.2	14.8	96.7	86.7
SimVLM (Wang et al., 2021b)	1.4 B	X	-	112.2	-	-	
BEIT-3 (Wang et al., 2022b)	1.9B	X	-	-	-	94.9	81.5
Flamingo (Alayrac e. al., 2022)	10.2 B	X	56.3				
BLIP-2	188M	✓	65.0	121.6	15.8	97.6	89.7



InstructBLIP

- BLIP v2 architecture
- Instruction tuning for visual reasoning benchmarks
 - 26 datasets (13 train)
 - Meticulously crafted instruction templates
- Good results



InstructBLTP

- BLIP v2 architecture
- Instruction tuning for visual reasoning benchmarks
 - 26 datasets (13 train)
 - Meticulously crafted instruction templates
- Good results

Task	Instruction Temp'are
linage Captioning	<image/> A short image caption: <image/> A short image description: <image/> A photo of <image/> An image that shows <image/> Write a short description for the image. <image/> Write a description for the photo. <image/> Provide a description of what is presented in the photo. <image/> Briefly describe the content of the image. <image/> Can you briefly explain what you see in the Image? <image/> Could you use a few words to describe what you perceive in the photo? <image/> Using language, provide a short depiction of the image. <image/> Using language, provide a short depiction of the image.
VQA	<pre><image/>{Question} <image/>{Question} <image/>{Question} <image/>{Question} A short answer to the question is <image/>Question} A short answer to the question is <image/>Question: {Question} A: <image/>Question: {Question} Short answer: <image/>Given the image, answer the following question with no more than three words, {Question} <image/>Based on the image, respond to this question with a short answer: {Question}. Answer: <image/>Use the provided image to answer the question: {Question} Provide your answer as short as possible: <image/>What is the answer to the following question? '{Question}'' <image/>The question "{Question}' can be answered using the image. A short answer is</pre>
VQG	<image/> Given the image, generate a question whose answer is: {Answer}. Question: <image/> Based on the image, provide a question with the answer (Answer) Question: <image/> Given the visual representation, create a question for which the answer is "{Answer}'. <image/> From the image provided, craft a question that leads to the reply: {Answer}. Question: <image/> Considering the picture, come up with a question where the answer is: {Answer}. <image/> Classed into account, generate an question that has the snewer: {Answer}.

Table 5: Instruction templates used for transforming held-in datasets into instruction tuning data. For datasets with OCR tokens, we simply add "OCR tokens:" after the image query embeddings.

GQA, VizWiz, iYQA, MSVD, MSRVTT <Image> Question: [] Short answer

NoCaps, Flickr30k <Image> A short image description:

TextVQA <Image> OCR tokens: {}. Question: {} Short answer:

IconQA <Image> Question: {} Options: {}. Short answer:

HatefulMennes <- Image> This is an image with: "{}" written on it. Is it hateful? Answer:

VSR <Image> Based on the image, is this statement true or false? "{}" Answer:

Visual Dialog < Image> Dialog history: {}\n Question: {} Short answer:

	ScienceQA IMG	OCR-VQA	OKVQA	A-OKVQA Direct Answer Multi-choid Val Test Val Tes					
Previous SOTA	LLaVA [25]	GIT [<mark>43</mark>]	PaLM-E(562B) [9]	[15]	[<mark>37</mark>]	[<mark>15</mark>]	[<mark>37</mark>]		
	89.0	70.3	66.1	56.3	61.6	73.2	73.6		
BLIP-2 (FlanT5 _{XXL})	89.5	72.7	54.7	57.6	53.7	80.2	76.2		
InstructBLIP (FlanT5 _{XXL})	90.7	73.3	55.5	57.1	54.8	81.0	76.7		
BLIP-2 (Vicuna-7B)	77.3	69.1	59.3	60.0	58.7	72.1	69.0		
InstructBLIP (Vicuna-7B)	79.5	72.8	62.1	64.0	62.1	75.7	73.4		

LLAVA v1

- Instruction tuning (visual dialogue) data is scarce
 - Let's use GPT4 (not GPT4-V) to create more data
 - Subset of CC3M (<600k)
 - Learn a single projection W on top of CLIP-VIT into Vicuna (llama2)

Visual Instruction Tuning, Liu etal 2023

Context type 1: Captions

A group of people standing outside of a black vehicle with various luggage. Luggage surrounds a vehicle in an underground parking area

People try to fit all of their luggage in an SUV.

The sport utility vehicle is parked in the public garage, being packed for a trip. Some people with luggage near a van that is transporting it.



Context type 2: Boxes

person: [0.681, 0.242, 0.774, 0.694], backpack: [0.384, 0.696, 0.485, 0.914], suitease: ...<omitted>

Response type 1: conversation

Question: What type of vehicle is featured in the image?

Answer: The image features a black sport utility vehicle (SUV) ...<omitted>

Response type 2: detailed description

The image is an underground parking area with a black sport utility vehicle (SUV) parked. There are three people in the scene, with one person standing closer to the left side of the vehicle, another person in the middle, and the third person on the right side. They are all working together to pack their luggage into the SUV for a trip. ...<

Response type 3: complex reasoning

Question: What challenges do these people face?

Answer: In the image, a group of people is standing outside a black SUV in a parking area, surrounded by various pieces of luggage, including suitcases and backpacks. They are facing the challenge of fitting all their luggage into the black SUV. There are multiple suitcases and backpacks to be packed, which suggests that the group has a significant amount of belongings ...<omitted>



LLAVA v1

- Stage 1 training: Adapter only
 - CC3M subset
 - GPT-4 ask question for caption
- Stage 2 training: Adapter + LLM
 - GPT-4 data (ask GPT-4 to do QA using caption and box inputs)
 - ScienceQA
- Quite easy to train (<1 day on 8 GPUs)

Visual Instruction Tuning, Liu etal 2023

```
messages = [{"role":"system", "content": f"""You are an AI visual assistant, and you are
seeing a single image. What you see are provided with five sentences, describing the same image you
are looking at. Answer all questions as you are seeing the image.
Design a conversation between you and a person asking about this photo. The answers should be in a
tone that a visual AI assistant is seeing the image and answering the question. Ask diverse questions
and give corresponding answers.
Include questions asking about the visual content of the image, including the object types, counting
the objects, object actions, object locations, relative positions between objects, ctc. Only include
questions that have definite answers:
(1) one can see the content in the image that the question asks about and can answer confidently;
(2) one can determine confidently from the image that it is not in the image. Do not ask any question
that cannot be answered confidently.
Also include complex questions that are relevant to the content in the image, for example, asking
about background knowledge of the objects in the image, asking to discuss about events happening in
the image, etc. Again, do not ask about uncertain details. Provide detailed answers when answering
complex questions. For example, give detailed examples or reasoning steps to make the content more
convincing and well-organized. You can include multiple paragraphs if necessary."""}
for sample in fewshot_samples:
    messages.append({"role":"user", "content":sample['context']})
   messages.append({"role":"assistant", "content":sample['response']} )
messages.append({"role":"user", "content": '\n'.join(query)})
```

Table 13: For each query, we illustrate the prompt construction process for ChatGPT/GPT-4 to collect query ['response'] from query ['context'], using few-shot in-context-learning, where examples are from fewahot_samples, each example including input sample['context'] and output sample['response']. Note that messages is the final prompt. In this example, we provide the prompt used to generate the conversation response, please see also see its in-context-learning examples in Table 15 and Table 16 for details. We recommend readers to check out the codebase for the prompts to generated two other types of responses, including detailed decription and complex reasoning.



Llava 1.5

- Better (higher resolution) image encoder
- Larger LLM
- More datasets
- Still quite fast to train
 - <1 day on 8-GPUs

Improved Baselines with Visual Instruction Tuning, Liu etal 2023



LLAVA-NEXT

Higher resolution, More data (OCR), Better 32 A100s)

Open-Source	Proprietary								
Data (PT)	Data (IT)	Model	MMMU (val)	Math-Vista	MMB-ENG	MMB-CN	MM-Vet	LLaVA-Wild	SEED-IMG
N/A	N/A	GPT-4V	56.8	49.9	75.8	73.9	67.6	-	71.6
N/A	N/A	Gemini Ultra	59.4	53	-	-	-	-	-
N/A	N/A	Gemini Pro	47.9	45.2	73.6	74.3	64.3	-	70.7
1.4B	50M	Qwen-VL-Plus	45.2	43.3	-	-	55.7	-	65.7
1.5B	5.12M	CogVLM-30B	32.1	-	-	-	56.8	-	-
125M	~1M	Yi-VL-34B	45.9	-	-	-	-	-	-
558K	665K	LLaVA-1.5-13B	36.4	27.6	67.8	63.3	36.3	72.5	68.2
558K	760K	LLaVA-NeXT-34B	51.1	46.5	79.3	79	57.4	89.6	75.9

• Higher resolution, More data (OCR), Better model, full model fine-tuning, still fast (1 day with



- LLaVA at multiple granularities
 - Granularity Controller = 2x2 average pooling
- Good results

Matryoshka Multimodal Models, Cai etal 2024



 X_{S_1}

Describe this image for me.





In the heart of a bustling restaurant, a young girl finds solace at a table.

In the heart of a bustling restaurant, a young girl with vibrant hair is seated at a wooden table, her attention captivated by the camera...

. . .



In the heart of a bustling restaurant, a young girl with long, dark hair is the center of attention. She's dressed in a blue and white striped sweater,. ... The table is adorned with a white paper bag, perhaps holding her meal. A blue Pepsi cup rests on the table .

Matryoshka Multimodal Models



some of them are holding a green flag while other are ...

Approach		# Tokens	MMBench	GQA	POPE	VizWiz	SEEDBench
Qwen-VL [7]		256	38.2	59.3	-	35.2	56.3
Qwen-VL-Chat [7	7]	256	60.6	57.5	-	38.9	58.2
InstructBLIP-7B	58]	32	36.0	49.2	-	34.5	53.4
InstructBLIP-13B [58]		32	-	49.5	78.9	33.4	-
LLaVA-1.5-7B [5]		576	64.8	62.0	85.9	54.4	60.5
		576	65.9	61.9	87.4	54.9	60.6
		144	66.4	61.3	87.0	53.1	59.7
LLaVA-1.5- M^3		36	64.8	60.3	85.5	52.8	58.0
		9	63.1	58.0	83.4	51.9	55.4
		1	59.5	52.6	78.4	49.4	50.1
Approach	# Toker	ns MSVD	MSRVTT	ActivityN	et NextQ	A IntentQ	A EgoSchem
Video-LLaMA [60]	-	51.6	29.6	12.4	-	_	
LLaMA-Adapter [61]	-	54.9	43.8	34.2	-	-	-
Video-ChatGPT [62]	-	64.9	49.3	35.2	-	-	-
Video-LLaVA [63]	2048	70.7	59.2	45.3	-	-	-
InternVideo [64]	-	-	-	-	59.1	-	32.1
LLaVA-NeXT-7B [4]	2880	78.8	63.7	54.3	63.1	60.3	35.8
	2880	78.2	64.5	53.9	63.1	58.8	36.8
	720	79.0	64.5	55.0	62.6	59.6	37.2
LLaVA-NeXT-7B-M ³	180	77.9	63.7	55.0	61.4	59.3	37.6
	15	75 8	62.0	52.2	50.5	587	38.8
	43	13.0	03.0	55.2	59.5	50.7	30.0



- Vision encoder (1.9B)
- Adapter (0.8B)
- LLM (7.7B)
- Inputs: Image, Bounding Box, Text
- Outputs: Bounding Box, Text
- Bounding box: Text format with special <box> and <ref> tokens

Qwen-VL: A Versatile Vision-Language Model for Understanding, Localization, Text Reading, and Beyond, Bai etal 2023







- Paper has lots of details on dataset setup and training recipes (base models etc)
- Great results for the time (Fall 2023)

Qwen-VL: A Versatile Vision-Language Model for Understanding, Localization, Text Reading, and Beyond, Bai etal 2023

Stage 1:

Language	Dataset	Original	Cleaned	Remaining%
	LAION-en	2B	280M	14%
	LAION-COCO	600M	300M	50%
	DataComp	1.4B	300M	21%
English	Соуо	700M	200M	28%
English	CC12M	12M	8M	66%
	CC3M	3M	3M	100%
	SBU	1M	0.8M	80%
	COCO Caption	0.6M	0.6M	100%
Chinese	LAION-zh	108M	105M	97%
Chinese	In-house Data	220M	220M	100%
	Total	5B	1.4B	28%

Task	# Samples	Dataset
Captioning	19.7M	LAION-en & zh, DataComp, Coyo, CC12M & 3M, SBU, COCO, In-house Data
VQA	3.6M	GQA, VGQA, VQAv2, DVQA, OCR-VQA, DocVQA, TextVQA, ChartQA, AI2D
Grounding ²	3.5M	GRIT
Ref Grounding	8.7M	GRIT, Visual Genome, RefCOCO, RefCOCO+, RefCOCOg
Grounded Cap.	8.7M	GRIT, Visual Genome, RefCOCO, RefCOCO+, RefCOCOg
OCR	24.8M	SynthDoG-en & zh, Common Crawl pdf & HTML
Pure-text Autoregression	7.8M	In-house Data

The Dataset Format Example of ChatML

Stage 3:

Stage 2:

<im_start>user Picture 1: vg/VG_100K_2/649.jpgWhat is the sign in the picture?<im_end> <im_start>assistant The sign is a road closure with an orange rhombus.<im_end> <im_start>user How is the weather in the picture?<im_end> <im start>assistant The shape of the road closure sign is an orange rhombus.<im_end>



- VLM with Flamingo-style fusion (but on parallel not gated sequence)
- Stage 1 Pre-training: Captioning 1.5B images (LAION-2B, COYO-700M filtered)
- Stage 2 Pre-training: Add in Referring Expression Comprehension (same data, boxes form open-vocab detector)

CogVLM: Visual Expert for Pretrained Language Models, Wang etal 2023



(a) The input of visual language model



(b) The visual expert built on the language model

CogVLM

- Alignment
 - Chat-style data (following LLaVa)
 - Grounded Captioning, Referring Expression (generation and comprehension), grounded VQA

Method	Train	NoCa	aps val	NoCa	ips test	Flickr	COCO	TextCaps
	Data	OOD	overall	OOD	overall	Karp.	Karp.	test
Human	-	95.7	87.1	91.6	85.3	-	-	125.1
VinVL (Zhang et al., 2021)	8.9M	83.8	94.3	78.0	92.5	-	130.8	-
SimVLM (Wang et al., 2021)	1.8B	115.2	112.2	109.5	110.3	-	143.3	-
CoCa (Yu et al., 2022)	4.8B	-	122.4	-	120.6	-	143.6	-
LEMON (Hu et al., 2022)	2B	120.2	117.3	110.1	114.3	-	139.1	-
Flamingo (Alayrac et al., 2022)	2.3B	-	-	-	-	67.2	138.1	-
Prismer (Liu et al., 2023d)	12.7M	113.5	112.9	-	110.8	-	136.5	-
BLIP-2 (Li et al., 2023b)	129M	124.8	121.6	-	-	-	144.5	-
InstructBLIP (Dai et al., 2023)	129M	-	123.1	-	-	82.4	-	-
UniversalCap (Cornia et al., 2021)	35M	123.4	122.1	114.3	119.3	-	143.4	-
GIT (Wang et al., 2022a)	0.8B	127.1	125.5	122.0	123.4	49.6	144.8	138.2
GIT2 (Wang et al., 2022a)	12.9B	<u>130.6</u>	<u>126.9</u>	122.3	124.8	50.7	145.0	<u>145.0</u>
Qwen-VL (Bai et al., 2023)	1.4B	-	121.4	-	-	85.8	-	-
PaLI-17B (Chen et al., 2022b)	1.6B	-	127.0	-	124.4	-	<u>149.1</u>	135.4
PaLI-X-55B (Chen et al., 2023b)	-	-	126.3	-	124.3	-	149.2	147.0
CogVLM (ours)	1.5B	132.6	128.3	128.0	126.4	94.9	148.7	144.9

Mathad	IIM	VQA						LVLM-Benchmark						
Method		VQAv2	OKVQA	TextVQA	OCRVQA	ScienceQA	MM-Ve	t SEED	MMBencl	h LLaVA	POPE	MMMU	MathVista	
MiniGPT-4	Vicuna-7B	-	-	0.6	-	39.6	22.1	47.4	23.0	45.1	-	-	23.1	
IDEFICS-Instruct	LLaMA-65B	37.4	36.9	32.9	-	61.8	39.7	53.2	54.5	56.9	-	-	26.2	
OpenFlamingo	MPT-7B	53.0	38.3	28.3	-	44.8	24.8	42.7	5.7	34.2	-	26.3	18.6	
DreamLLM	Vicuna-7B	56.6	44.3	34.9	-	-	35.9	-	49.9	-	-	-	-	
InstructBLIP	Vicuna-7B	-	-	50.1	-	60.5	26.2	58.8	33.9	59.8	53.8	-	25.3	
Fuyu	Fuyu-8B	74.2*	60.6*	-	-	-	-	-	-	-	-	27.4	-	
Qwen-VL-Chat	Qwen-7B	78.2*	56.6*	61.5*	70.5*	68.8	-	65.4	61.8	67.7	-	32.9	<u>33.8</u>	
LLaVA-1.5	Vicuna-7B	78.5*	-	58.2*	-	66.8	30.5	58.6	64.3	60.7	85.9	-	23.6	
mPLUG-Owl2	LLaMA2-7B	79.4*	57.7*	58.2*	-	68.7	36.2	64.1	64.5	25.0	86.2	32.1	25.3	
Unified-IO2	UIO-2XXL	79.4*	55.5*	-	-	86.2*	-	65.6	71.5	-	<u>87.7</u>	-	-	
LLaVA-1.5	Vicuna-13B	80.0*	-	61.3*	-	71.6	35.4	61.6	67.7	64.6	85.9	33.6	26.1	
SPHINX-2k	LLaMA2 13B	80.7*	62.6*	61.2*	67.8*	70.6	40.2	<u>71.6</u>	65.9	-	87.2	32.9	27.8	
Emu2-Chat	LLaMA-33B	84.9*	<u>64.8*</u>	<u>66.6*</u>	-	-	<u>48.5</u>	62.8	63.6	56.4	-	<u>34.1</u>	-	
CogVLM-Chat	Vicuna-7B	82.3*	64.8*	70.4*	73.8*	91.2*	51.1	72.5	77.6	77.8	87.9	41.1	34.5	

Туре	Model	I	RefCOC	0	R	efCOCC)+	RefC	OCOg	Visual7W
-, -, -, -, -, -, -, -, -, -, -, -, -, -		val	test-A	test-B	val	test-A	test-B	val	test	test
	OFA-L* (Wang et al., 2022b)	79.96	83.67	76.39	68.29	76.00	61.75	67.57	67.58	-
Generalist	VisionLLM-H (Wang et al., 2023b)	-	86.70	-	-	-	-	-	-	-
	Shikra-7B (Chen et al., 2023a)	87.01	90.61	80.24	81.60	87.36	72.12	82.27	82.19	-
	Shikra-13B (Chen et al., 2023a)	87.83	91.11	81.81	82.89	87.79	74.41	82.64	83.16	85.33
	Qwen-VL (Bai et al., 2023)	89.36	92.26	85.34	83.12	88.25	77.21	85.58	85.48	-
	Ferret-13B (You et al., 2023)	89.48	92.41	84.36	82.81	88.14	75.17	85.83	86.34	-
	CogVLM-Grounding	92.76	94.75	88.99	88.68	92.91	83.39	89.75	90.79	91.05
Constantine (G-DINO-L (Liu et al., 2023e)	90.56	93.19	88.24	82.75	88.95	75.92	86.13	87.02	-
	UNINEXT-H (Lin et al., 2023a)	92.64	94.33	91.46	85.24	89.63	79.79	88.73	89.37	-
specialisi	ONE-PEACE (Wang et al., 2023a)	92.58	94.18	89.26	88.77	92.21	83.23	89.22	89.27	-

OtterHD



- Built on Perssimon-8B and Fuyu-8B (decoder only transformers)
 - Images fed in tokenized
 - Image new-line token
- Dynamic image resolution
- Instruction-tuned

							How ma Model It is image visib	any cam . Resp hard t . Howev le in t	els are onse: o exami er, I d he imag	e insid ine the find th ge.	de thi e came nere a	s image els in s are thre	e? Such in Se came	ntrica els	te	mon
						Π	F	U YU mer De	coder							
Linear projection ——	↑ ↑	†	^ -	n) N	↑ 1	<u> </u>	 _\n	^ 		Î		This	is)	a	 Ецуа)	
Image patch		1			- All											

Models	I/D Pairs	Train Rec	Fyal Ras	MagBench		M	\mathbf{ME}^1	POPE		MMR	 M_Vista
WIUUCIS				Multi.	FF.	Cog.	Percep.		141141- V		111- 11510
Idefics-9B _{instruct} [24]	1M	224	224	20.8	13.4	187.9	1165.0	74.6	23.7	45.5	19.8
Otter-9B [25]	150K	224	224	25.7	15.8	306.4	1292.3	72.5	24.7	48.3	19.7
InstructBLIP-7B [13]	1.2M	224	224	5.6	15.2	-	-	-	26.2	36.0	-
InstructBLIP-13B [13]	1.2M	224	224	3.8	16.3	291.8	1212.8	78.9	25.6	33.9	25.3
LLaVA-7B _{1.5} [30]	$3.6M^2$	336	336	26.8	24.7	-	1510.7	85.9	30.5	<u>59.5</u>	-
Qwen-VL-7B _{chat} [4]	1.4B	448	448	14.5	15.9	360.7	1487.5	-	-	61.8	-
Fuyu-8B [5]	-	_	Original	29.3	15.2	237.5	728.6	74.1	21.4	10.7	20.6
		512	512	33.5	31.4	289.8	1359.3	86.1	25.1	58.5	22.3
OtterHD-8B	370K	1024	1024	37.8	37.2	288.5	1313.7	81.5	19.8	53.6	17.3
		Dynamic	Original	42.7	39.9	331.4	1223.4	<u>86.0</u>	<u>26.3</u>	58.3	<u>23.5</u>

Dataset	LLaVA-DD/CR	VQAv2	GQA	OKVQA	OCRVQA	A-OKVQA	COCO-GOI
Avg. W	577	581	495	617	352	587	586
Avg. H	481	482	409	448	490	482	476
Pairs	53240	20000	30000	18018	16354	34112	20000
Dataset	COCO-Caption	TextQA	RefCOCO	COCO-ITM	ImageNet	LLaVA-RLHF	Combined
Avg. W	578	950	591	577	469	340	542
Avg. H	484	811	486	484	387	572	467
Pairs	20000	19293	20000	20000	50000	50000	371017

OtterHD



- More data
- Faster base-model
- New High-res VLM benchmark (MagnifierBench)

OtterHD: A High-Resolution Multi-modality Model, Li etal 2023

Data Mixture We compiled a total of 370K instruction/response pairs sourced from the following public datasets: LLaVA-Instruct [30], VQAv2 [2], GQA [23], OKVQA [36], OCRVQA [38], A-OKVQA [45], COCO-GOI [33], COCO-Caption [10], TextQA [48], RefCOCO [58], COCO-ITM [28], ImageNet [17], and LLaVA-RLHF [51]. The data mixture and specific prompt strategies are motivated by LLaVA-1.5 [30] and Idefics-Instruct [24] to achieve better text formatting control. All the datasets were organized into instruction/response pairs, aggregated into a single dataloader and uniformly sampled during the training phase to ensure representational integrity.







VILA

- Lot's of tricks and ablations
 - How to interleave multiple image-text pairs
 - Blending text-only and image-text data





Method	LLM	Res.	PT	IT	VQA ^{v2}	GQA	VisWiz	SQA ^I	VQA ^T	POPE	MME	MMB	MMB ^{CN}	SEED	LLaVA ^W	MM-Vet
BLIP-2 [35]	Vicuna-13B	224	129M	-	41.0	41	19.6	61	42.5	85.3	1293.8	_	_	46.4	38.1	22.4
InstructBLIP [18]	Vicuna-7B	224	129M	1.2M	_	49.2	34.5	60.5	50.1	-	_	36	23.7	53.4	60.9	26.2
InstructBLIP [18]	Vicuna-13B	224	129M	1.2M	_	49.5	33.4	63.1	50.7	78.9	1212.8	_	_	_	58.2	25.6
Shikra [12]	Vicuna-13B	224	600K	5.5M	77.4*	_	_	_	_	-	_	58.8	-	_	_	_
IDEFICS-9B [30]	LLaMA-7B	224	353M	1 M	50.9	38.4	35.5	_	25.9	-	_	48.2	25.2	_	_	_
IDEFICS-80B [30]	LLaMA-65B	224	353M	1 M	60.0	45.2	36.0	_	30.9	-	_	54.5	38.1	_	_	_
Qwen-VL [9]	Qwen-7B	448	1.4B	50M	78.8*	59.3*	35.2	67.1	63.8	-	_	38.2	7.4	56.3	_	_
Qwen-VL-Chat [9]	Qwen-7B	448	1.4B	50M	78.2*	57.5*	38.9	68.2	61.5	-	1487.5	60.6	56.7	58.2	_	_
LLaVA-1.5 [38]	Vicuna-1.5-7B	336	0.6M	0.7M	78.5*	62.0*	50.0	66.8	58.2	85.9	1510.7	64.3	58.3	58.6	63.4	30.5
LLaVA-1.5 [38]	Vicuna-1.5-13B	336	0.6M	0.7M	<u>80.0</u> *	63.3*	53.6	<u>71.6</u>	61.3	85.9	1531.3	67.7	<u>63.6</u>	<u>61.6</u>	<u>70.7</u>	<u>35.4</u>
VILA-7B (ours)	Llama-2-7B	336	50M	1M	79.9*	<u>62.3</u> *	<u>57.8</u>	68.2	<u>64.4</u>	<u>85.5</u>	<u>1533.0</u>	<u>68.9</u>	61.7	61.1	69.7	34.9
VILA-13B (ours)	Llama-2-13B	336	50M	1 M	80.8*	63.3*	60.6	73.7	66.6	84.2	1570.1	70.3	64.3	62.8	73.0	38.8
+ShareGPT4V	Llama-2-13B	336	50M	1 M	80.6*	63.2*	62.4	73.1	65.3	84.8	1556.5	70.8	65.4	61.4	78.4	45.7

$V \in (| | | P$

- Better captioning pipeline for VLMs and dataset (VeCap)
 - Ethics check (does the LLM reply: "I am sorry that I cannot ...")
 - LLM Knowledge Fusion: Rephrase the following two sentences into one...
- Better CLIP model

VeCLIP: Improving CLIP Training via Visual-enriched Captions, Lai etal 2023







LLM-VeC Caption:

A red brick house with a while roof and a front porch, surrounded by a garden at 112 Lafayette Dr

AltText:

SURF IN, SURF OUT ~ LARGE condo, Steps to Beach + Parking.

LLM-VeC Caption:

A small, modern home with a patio, outdoo furniture, and a stone wall.

AltText:

More then one million Australian households are already in mortgage stress.

LLM-VeC Caption:

A real estate sign advertises a house in Australian for sale in a residential neighborhood, with a tree in the background.



AltText:

9 Misconceptions About Alcohol.

LLM-VeC Caption: A glass of beer sits on a table with a lit candle. creating a cozy atmosphere.

AltText: Ring Car Cam.

LLM-VeC Caption: A hand is holding a smartphone with Ring Car

Cam, watching a car driving on a road.

AltText:

3 bedroom apartment to rent Croydon.

LLM-VeC Caption: A bedroom with a bed with white mattress, a window, and a wooden dresser, providing a comfortable and well-organized space.





Vision language models Trend 1

Architectures stabilized

- Image -> Sequence (vision tokens)
- Test -> Sequence (text tokens)
- Transformer



Vision language models Trend 2

• Data matters

 Innovations in training and data yield large improvements

LLM-VeC Caption:

A red brick house with a white roof and a front porch, surrounded by a garden at 112 Lafayette Dr

SURF IN, SURF OUT ~ LARGE condo, Steps to Beach + Parking.

LLM-VeC Caption:

A small, modern home with a patio, outdoor furniture, and a stone wall.

AltText:

9 Misconceptions About Alcohol.

LLM-VeC Caption: A glass of beer sits on a table with a lit candle, creating a cozy atmosphere.

AltText:

Ring Car Cam.

LLM-VeC Caption:

A hand is holding a smartphone with Ring Car Cam, watching a car driving on a road.

AltText:

More then one million Australian households are already in mortgage stress.

LLM-VeC Caption:

A real estate sign advertises a house in Australian for sale in a residential neighborhood, with a tree in the background.

3 bedroom apartment to rent Croydon.

AltText:

A bedroom with a bed with white mattress, a window, and a wooden dresser, providing a comfortable and well-organized space.

References

- Unified-IO: A Unified Model for Vision, Language, and Multi-Modal Tasks, Lu etal. 2022
- Flamingo: a Visual Language Model for Few-Shot Learning, Alayrac etal. 2022
- BLIP: Bootstrapping Language-Image Pre-training for Unified Vision-Language Understanding and Generation, Li etal. 2022
- BLIP-2: Bootstrapping Language-Image Pre-training with Frozen Image Encoders and Large Language Models, Li etal. 2023
- InstructBLIP: Towards General-purpose Vision-Language Models with Instruction Tuning, Dai etal. 2023
- Visual Instruction Tuning, Liu etal. 2023
- Improved Baselines with Visual Instruction Tuning, Liu etal. 2023
- Matryoshka Multimodal Models, Cai etal. 2024
- Qwen-VL: A Versatile Vision-Language Model for Understanding, Localization, Text Reading, and Beyond, Bai etal 2023
- CogVLM: Visual Expert for Pretrained Language Models, Wang etal. 2023
- OtterHD: A High-Resolution Multi-modality Model, Li etal. 2023
- VILA: On Pre-training for Visual Language Models, Lin etal. 2023
- VeCLIP: Improving CLIP Training via Visual-enriched Captions, Lai etal. 2023