Flamingo:

a Visual Language Model for Few-Shot Learning

Andrea Wynn and Xindi Wu 11/21/2022





Overview

Motivation

Flamingo Model Architecture

Training Data & Objective

In-Context Learning & Fine Tuning

Evaluation & Ablation Results

Limitations

Related Work: CM3 & Frozen

Discussion



GPT-3 VIT VisualBERT CLIP





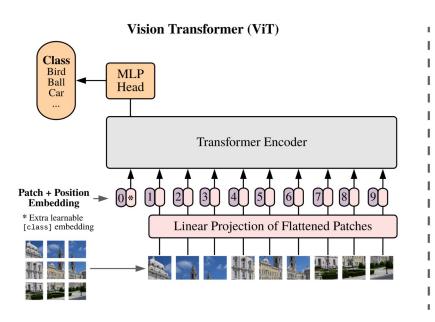
GPT-3

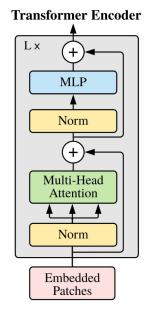
VIT

VisualBERT

CLIP

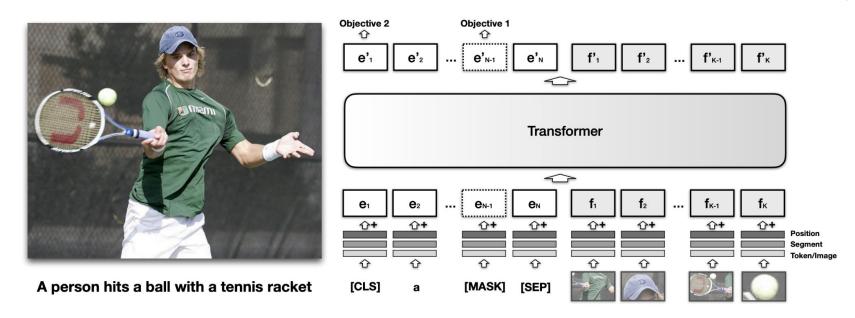
?







GPT-3 VIT VisualBERT CLIP



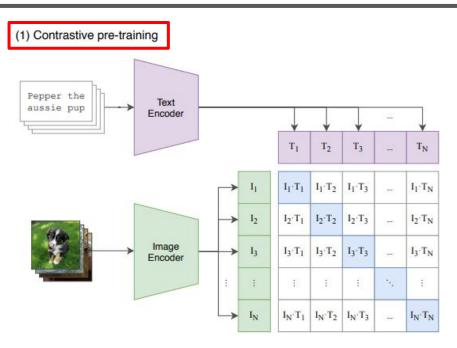


GPT-3

VIT

VisualBERT

CLIP





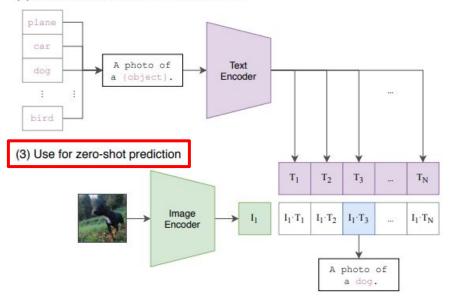
GPT-3

VIT

VisualBERT

CLIP

(2) Create dataset classifier from label text





GPT-3

VIT

VisualBERT

CLIP



The first vision-language model that has in-context learning ability



Motivation | Challenges

GPT-3

VIT

VisualBERT

CLIP

Flamingo

Challenges of multimodal generative modelling

- Unifying strong single-modal models
 - Interleave **cross-attention** layers with language only self-attention layers



Motivation | Challenges

GPT-3

VIT

VisualBERT

CLIP

Flamingo

Challenges of multimodal generative modelling

- Unifying strong single-modal models
 - Interleave **cross-attention** layers with language only self-attention layers
- Supporting images and videos
 - Perceiver-based architecture with a fixed number of visual tokens



Motivation | Challenges

GPT-3

VIT

VisualBERT

CLIP

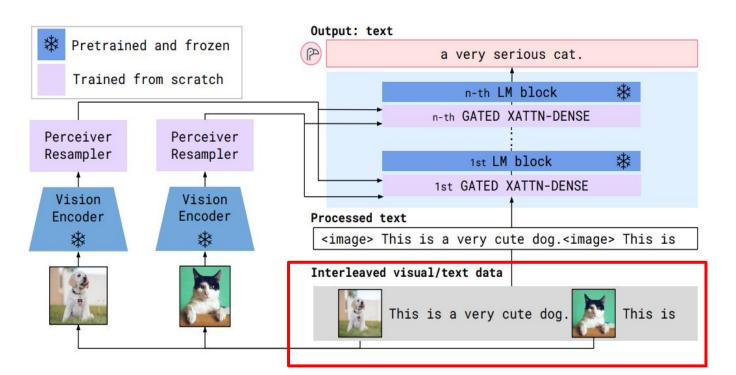
Flamingo

Challenges of multimodal generative modelling

- Unifying strong single-modal models
 - Interleave **cross-attention** layers with language only self-attention layers
- Supporting images and videos
 - Perceiver-based architecture with a fixed number of visual tokens
- Heterogeneous training data
 - Combine web scraping with existing image-text or video-text datasets.



Separately trained image + language models, with novel layers in between





Input/Output

Interleaved inputs: text/images/video

Selected single image samples

Question: What is the title of

the book? Answer:

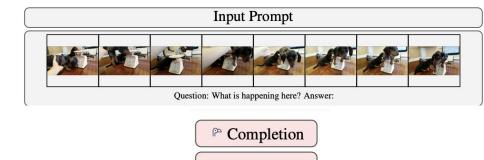
Completion
The House Book.

Selected dialogue samples



Outputs: free-form text

Selected video samples.



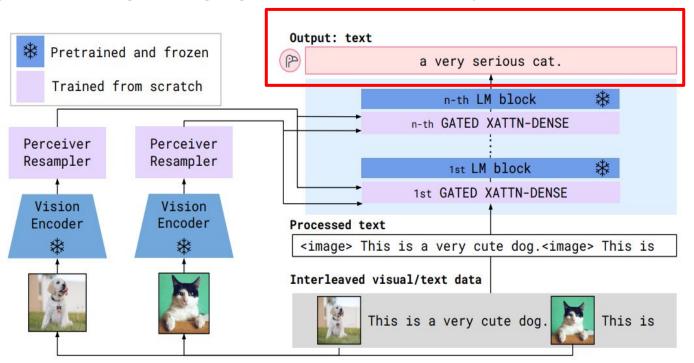
The dachschund puppy

is being weighed on a

scale.

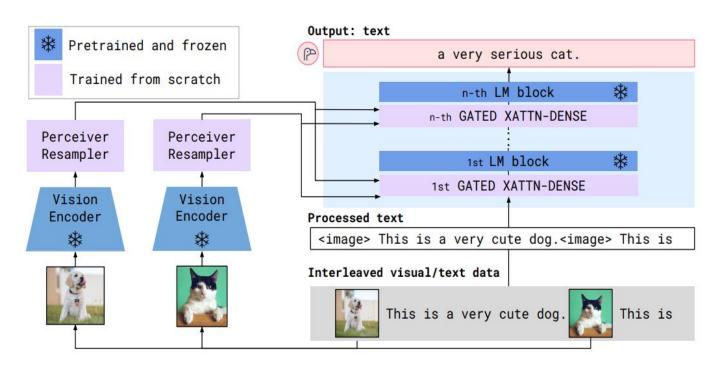


Separately trained image + language models, with novel layers in between





Separately trained image + language models, with novel layers in between



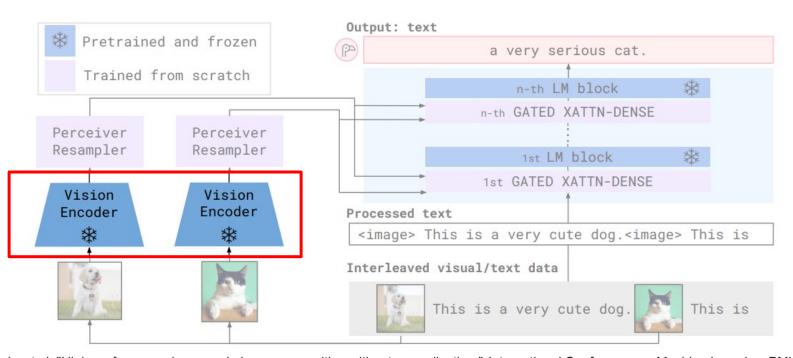


$$p(y|x) = \prod_{\ell=1}^{L} p(y_{\ell}|y_{<\ell}, x_{\leq \ell}),$$



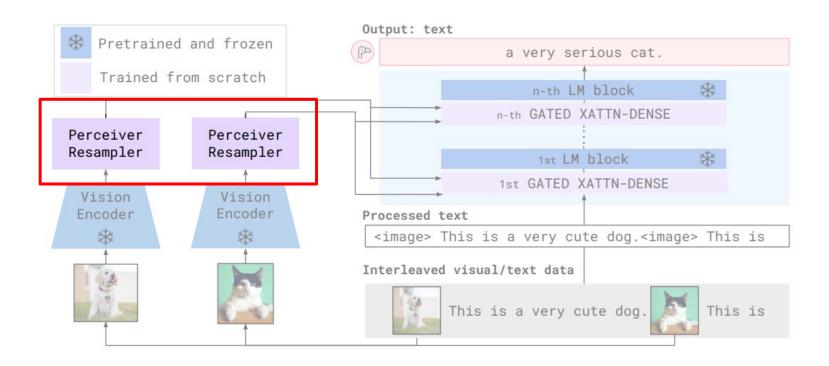
Vision Encoder

Pretrained and frozen Normalizer Free ResNet (NFNet)

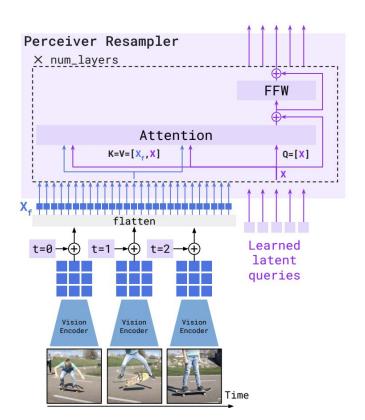


Brock, Andy, et al. "High-performance large-scale image recognition without normalization." International Conference on Machine Learning. PMLR, 2021. 17



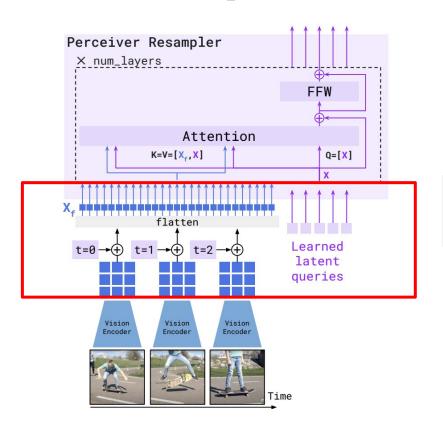






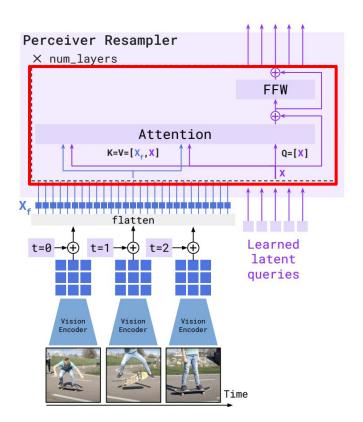
```
def perceiver_resampler(
    x_f, # The [T, S, d] visual features (T=time, S=space)
    time_embeddings, # The [T, 1, d] time pos embeddings.
    x, # R learned latents of shape [R, d]
    num_layers, # Number of layers
  # Add the time position embeddings and flatten.
  x_f = x_f + time_embeddings
  x_f = flatten(x_f) \# [T, S, d] \rightarrow [T * S, d]
  # Apply the Perceiver Resampler layers.
    # Attention.
    x = x + attention_i(q=x, kv=concat([x_f, x]))
   x = x + ffw_i(x)
```





```
x_f, # The [T, S, d] visual features (T=time, S=space)
  time_embeddings, # The [T, 1, d] time pos embeddings.
  x, # R learned latents of shape [R, d]
  num_layers, # Number of layers
# Add the time position embeddings and flatten.
x_f = x_f + time_embeddings
x_f = flatten(x_f) \# [T, S, d] \rightarrow [T * S, d]
# Apply the Perceiver Resampler layers.
  # Attention.
  x = x + attention_i(q=x, kv=concat([x_f, x]))
 # Feed forward.
 x = x + ffw_i(x)
return x
```

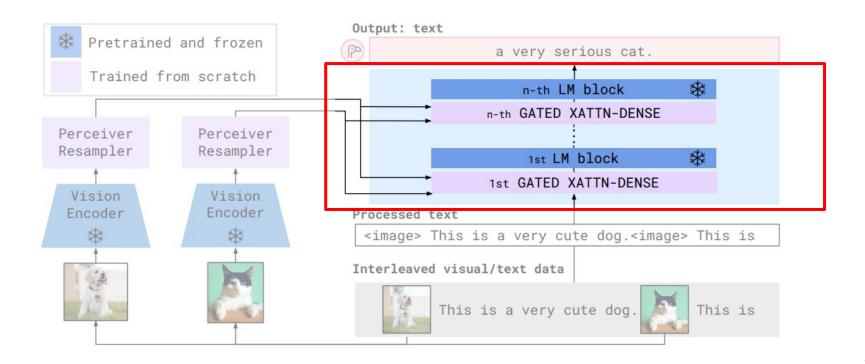




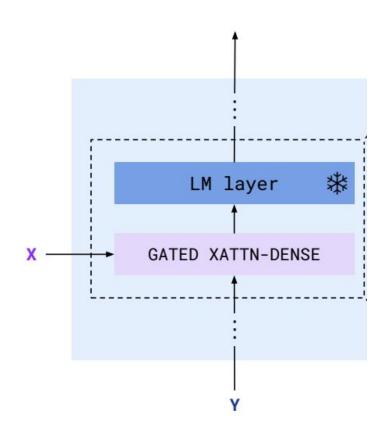
```
x_f, # The [T, S, d] visual features (T=time, S=space)
  time_embeddings, # The [T, 1, d] time pos embeddings.
 x, # R learned latents of shape [R, d]
  num_layers, # Number of layers
# Add the time position embeddings and flatten.
x_f = x_f + time_embeddings
x_f = flatten(x_f) \# [T, S, d] \rightarrow [T * S, d]
# Apply the Perceiver Resampler layers.
for i in range(num_layers):
  # Attention.
  x = x + attention_i(q=x, kv=concat([x_f, x]))
  # Feed forward.
 x = x + ffw_i(x)
return x
```



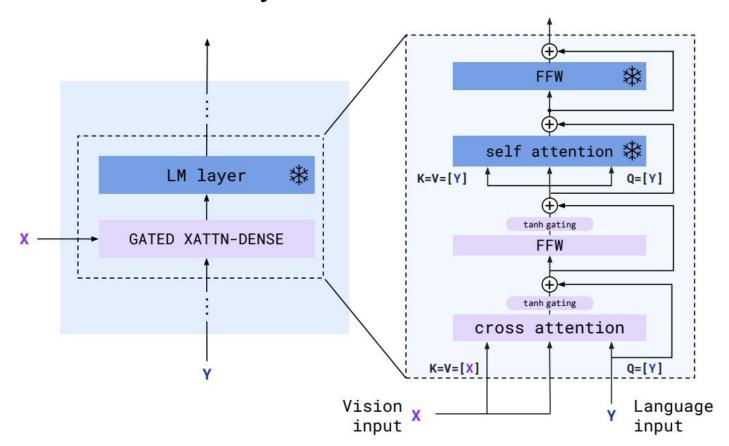
Conditioning the Language Model



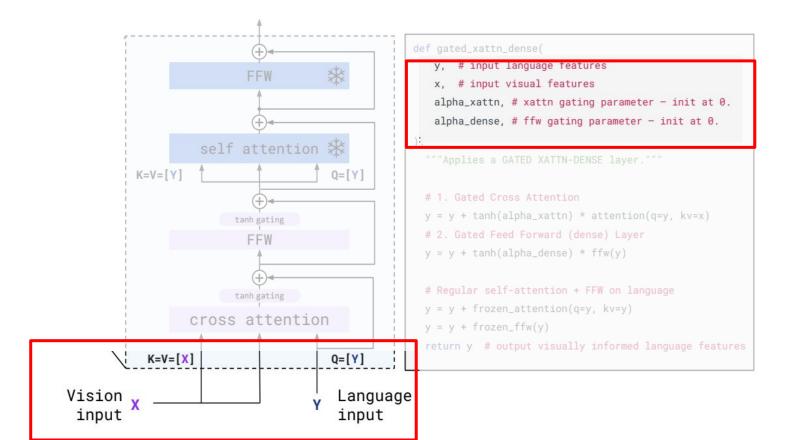




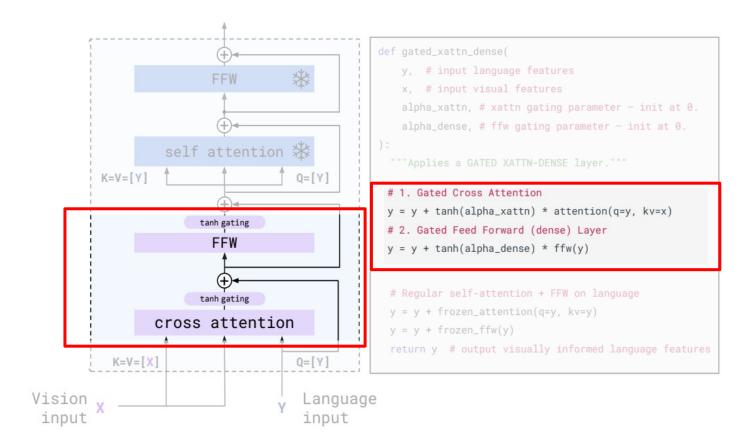




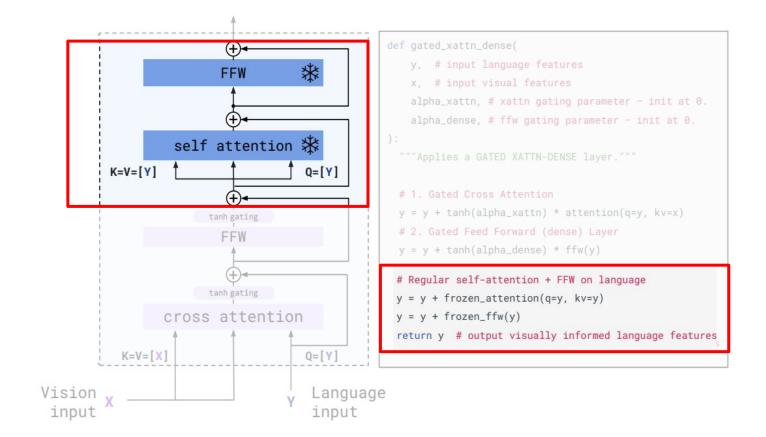




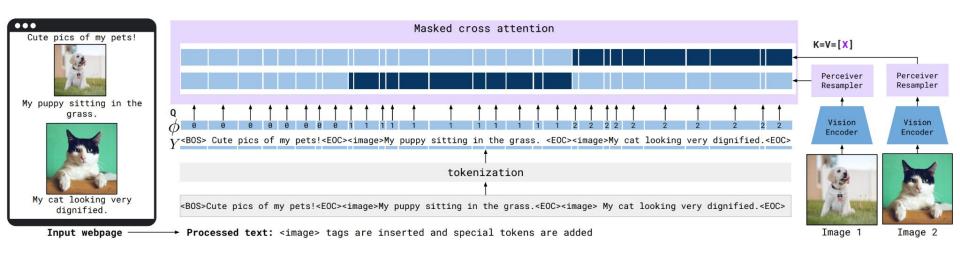












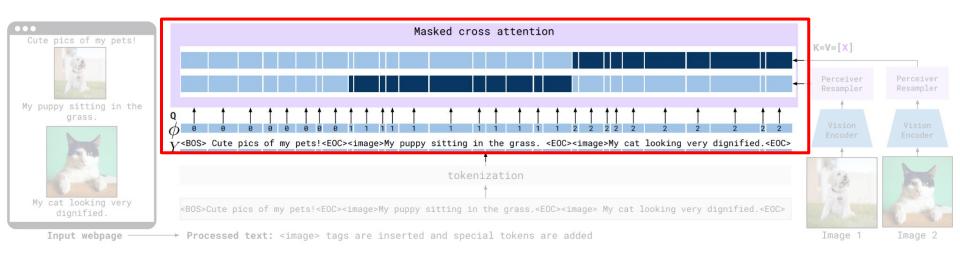














Pre-Lecture Question

Describe how Flamingo handles input sequences of arbitrarily interleaved textual and visual data, and combines pre-trained text-only and vision-only models.

Answer:

For example, the input contains an image of a dog together with a text description and an image of a cat with an incomplete text description. The text is parsed from the input with images replaced with placeholders the images are also extracted from the input passed through a frozen vision encoder and then mapped through the perceiver resampler to produce a fixed number of visual tokens per input.



Training Data



Mixture of Datasets



This is an image of a flamingo.



A kid doing a kickflip.





This is a picture of my dog.



This is a picture of my cat.

Image-Text Pairs dataset [N=1, T=1, H, W, C]

Video-Text Pairs dataset [N=1, T>1, H, W, C]

Multi-Modal Massive Web (M3W) dataset [N>1, T=1, H, W, C]

- N: Number of visual inputs for a single example
- T: Number of video frames
- H, W, C: height, width, color channels



Interleaved Image/Text: MultiModal MassiveWeb (M3W)

- Interleaved text and image training data
- Compiled from webpage HTML
- Randomly sample 256 token subsequence and extract first 5 images

Example:



Multi-Modal Massive Web (M3W) dataset [N>1, T=1, H, W, C]



Image-Text Pairs: ALIGN



"motorcycle front wheel"



"thumbnail for version as of 21 57 29 june 2010"



"file frankfurt airport skyline 2017 05 jpg"



"file london barge race 2 jpg"



"moustache seamless wallpaper design"



"st oswalds way and shops"

Source: https://arxiv.org/pdf/2102.05918v2.pdf

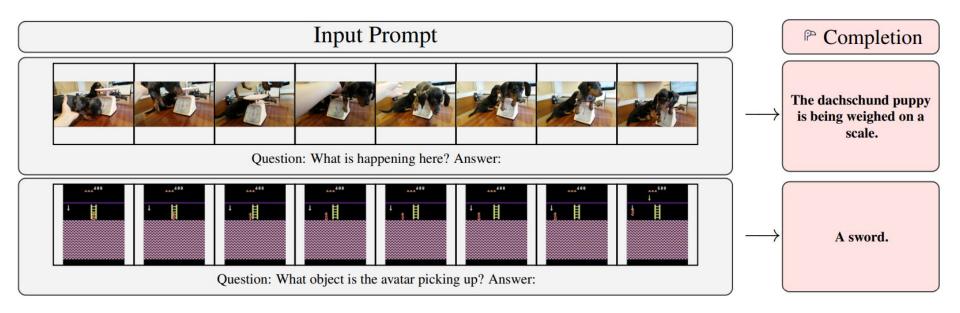


Image-Text Pairs: Long Text & Image Pairs (LTIP)





Video & Text Pairs (VTP)





Data Augmentation & Preprocessing

- Visual inputs resized to 320x320
- M3W Data Augmentation: Randomizing image placement

(a) This is my dog! <dog image>

(b) <dog image> That was my dog!

This is my cat! <cat image>

<cat image> That was my cat!



Training Objective

$$\sum_{m=1}^{M} \lambda_m \cdot \mathbb{E}_{(x,y) \sim \mathcal{D}_m} \left[-\sum_{\ell=1}^{L} \log p(y_{\ell}|y_{<\ell}, x_{\leq \ell}) \right]$$

- Weighted sum of dataset specific expected negative log likelihood of text, given some visual inputs
- AdamW optimizer
- No weight decay for Perceiver Resampler
- Weight decay of 0.1 for other parameters



Pre-Lecture Question

Describe what datasets are used for mixed training. How important is each type of dataset empirically?

Answer:

Datasets - M3W (interleaved images and text), ALIGN (large, lower quality image + text pairs), LTIP (image + text pairs), VTP (video + text pairs)

Importance (lambda weights) - 1.0 (M3W), 0.2 (ALIGN), 0.2 (LTIP), 0.03 (VTP)

Number of datasets (M) - 4



Flamingo Evaluation



Benchmark Tasks

	Dataset	DEV	Gen.	Custom prompt	Task description
	ImageNet-1k [94]	/			Object classification
	MS-COCO [15]	1	1		Scene description
	VQAv2 [3]	1	1		Scene understanding QA
O	OKVQA [69]	1	/		External knowledge QA
Image	Flickr30k [139]		1		Scene description
Im	VizWiz [35]		1		Scene understanding QA
	TextVQA [100]		1		Text reading QA
	VisDial [20]				Visual Dialogue
	HatefulMemes [54]			✓	Meme classification
	Kinetics700 2020 [102]	/			Action classification
	VATEX [122]	1	1		Event description
	MSVDQA [130]	1	/		Event understanding QA
	YouCook2 [149]		1		Event description
0	MSRVTTQA [130]		1		Event understanding QA
de	iVQA [135]		1		Event understanding QA
\leq	RareAct [73]			1	Composite action retrieval
	NextQA [129]		/		Temporal/Causal QA
	STAR [128]				Multiple-choice QA



Benchmark Tasks: ImageNet-1k



Source: https://link.springer.com/content/pdf/10.1007/s11263-015-0816-y.pdf



Benchmark Tasks: Visual Question Answering (VQA)



What color are her eyes? What is the mustache made of?



Is this person expecting company? What is just under the tree?



How many slices of pizza are there? Is this a vegetarian pizza?



Does it appear to be rainy?

Does this person have 20/20 vision?

Source: https://link.springer.com/content/pdf/10.1007/s11263-016-0966-6.pdf



Benchmark Tasks: Kinetics700 2020

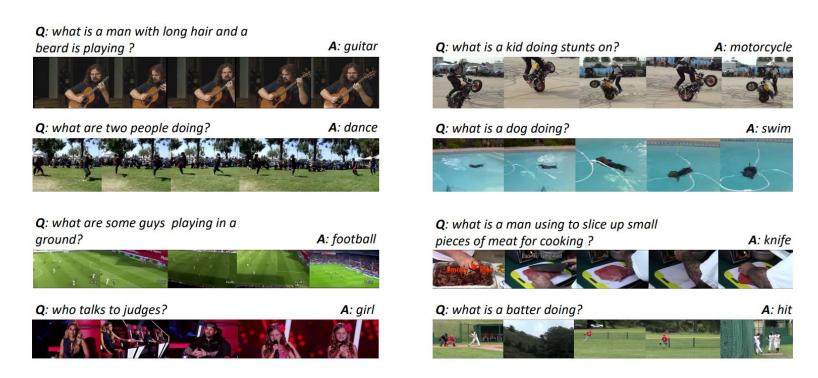
- Taken from YouTube videos
- Format: label, youtube_id, start time, end time

label	youtube_id	time_start	time_end
clay pottery making	0dWlqevl	19	29
javelin throw	07WQ2iBlw	1	11
climbing a rope	ONTAs-fAO	29	39
sipping cup	0l35AkU34	68	78
flipping pancake	33Lscn6sk	4	14
tickling	3OAstUWtU	45	55

Source: https://arxiv.org/pdf/2210.10864.pdf



Benchmark Tasks: MSVDQA



Source: https://dl.acm.org/doi/pdf/10.1145/3123266.3123427



Classification Task Results

Model	Method	Prompt size	shots/class	ImageNet top 1	Kinetics700 avg top1/5
SotA	Fine-tuned	=	full	90.9 [127]	89.0 [134]
SotA	Contrastive	Η.	0	85.7 [82]	69.6 [85]
NFNetF6	Our contrastive	_	0	77.9	62.9
		8	1	70.9	55.9
Flamingo-3B	RICES	16	1	71.0	56.9
		16	5	72.7	58.3
		8	1	71.2	58.0
Flamingo-9B	RICES	16	1	71.7	59.4
· ·		16	5	75.2	60.9
	Random	16	≤ 0.02	66.4	51.2
		8	1	71.9	60.4
Flamingo-80B	RICES	16	1	71.7	62.7
		16	5	76.0	63.5
	RICES+ensembling	16	5	77.3	64.2



Fine Tuning Results

Method	VAOV	7 (42)	0000	VATEX	VizWiz	71 W 71 V	MSRVTTQA		VisDial	YouCook2		TextVQA	HatefulMemes
	test-dev	test-std	test	test	test-dev	test-std	test	valid	test-std	valid	valid	test-std	test seen
* Flamingo - 32 shots	67.6	-	113.8	65.1	49.8	-	31.0	56.8	-	86.8	36.0	-	70.0
SimVLM [124]	80.0	80.3	143.3	-		-	-	-	1=6	=2	1=4	15	-
OFA [119]	79.9	80.0	149.6	-	-	-	-	<u></u>	-	-	-	-	-
Florence [140]	80.2	80.4	-	-	-	3 -	_	-	-	-	-	-	-
Flamingo Fine-tuned	82.0	82.1	138.1	84.2	65.7	65.4	<u>47.4</u>	61.8	59.7	118.6	57.1	54.1	86.6
Restricted SotA [†]	80.2	80.4	143.3	76.3	-	(-	46.8	75.2	74.5	138.7	54.7	73.7	79.1
Restricted SolA	[140]	[140]	[124]	[153]	-	-	[51]	[79]	[79]	[132]	[137]	[84]	[62]
Unrestricted SotA	81.3	81.3	149.6	81.4	57.2	60.6	-	-	<u>75.4</u>	-	-	-	84.6
Uniestricted SolA	[133]	[133]	[119]	[153]	[65]	[65]	-	-	[123]	-	-	-	[152]

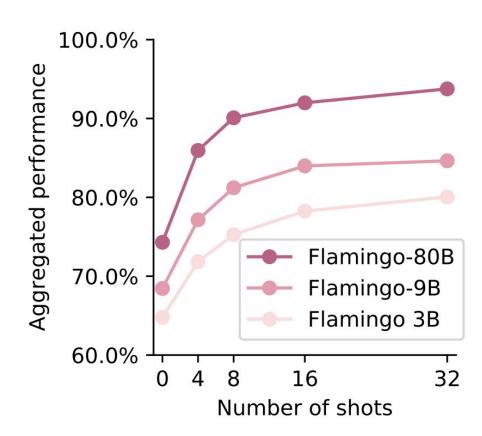


Model Scaling

	Requires	Froze	en	Trainable		Total
	model sharding	Language	Vision	GATED XATTN-DENSE	Resampler	count
Flamingo-3B	Х	1.4B	435M	1.2B (every)	194M	3.2B
Flamingo-9B	X	7.1B	435M	1.6B (every 4th)	194M	9.3B
Flamingo	✓	70B	435M	10B (every 7th)	194M	80B



Number of Shots





Ablation Studies



Ablation Studies

	Ablated setting	Flamingo-3B original value	Changed value	Overall score
		Flamingo-31	3 model	70.7
(i)	Training data	All data	w/o Video-Text pairs w/o Image-Text pairs Image-Text pairs → LAION w/o M3W	67.3 60.9 66.4 53.4
(ii)	Optimisation	Accumulation	Round Robin	62.9
(iii)	Tanh gating	1	X	66.5
(iv)	Cross-attention architecture	GATED XATTN-DENSE	VANILLA XATTN GRAFTING	66.9 63.1
(v)	Cross-attention frequency	Every	Single in middle Every 4th Every 2nd	59.8 68.8 68.2
(vi)	Resampler	Perceiver	MLP Transformer	66.6 66.7
(vii)	Vision encoder	NFNet-F6	CLIP ViT-L/14 NFNet-F0	64.9 62.7
(viii)	Freezing LM	✓	X (random init) X (pretrained)	57.8 62.7



Pre-Training Dataset Ablation

Dataset	Combination	ImageNet			COC			
	strategy	accuracy		nage-to-t			kt-to-ima	
		top-1	R@1	R@5	R@10	R@1	R@5	R@10
LTIP	None	40.8	38.6	66.4	76.4	31.1	57.4	68.4
ALIGN	None	35.2	32.2	58.9	70.6	23.7	47.7	59.4
LTIP + ALIGN	Accumulation	45.6	42.3	68.3	78.4	31.5	58.3	69.0
LTIP + ALIGN	Data merged	38.6	36.9	65.8	76.5	15.2	40.8	55.7
LTIP + ALIGN	Round-robin	41.2	40.1	66.7	77.6	29.2	55.1	66.6



Frozen Language Model

	Ablated setting	Flamingo 3B value	Changed value	Overall score ↑
		Flamingo 3B mode	el (short training)	70.7
(i)	Resampler size	Medium	Small Large	67.9 69.0
(ii)	Multi-Img att.	Only last	All previous	63.5
(iii)	p_{next}	0.5	0.0 1.0	69.6 70.4
(iv)	LM pretraining	MassiveText	C4	62.8
(v)	Freezing Vision	✓	X (random init) X (pretrained)	61.4 68.1
(vi)	Co-train LM on MassiveText	×	✓ (random init) ✓ (pretrained)	55.9 68.6
(vii)	and Vision encoder	and NFNetF6	M3W+LAION400M+VTP and CLIP	54.7 64.9

0-initialized tanh gating

	Ablated setting	Flamingo-3B original value	Changed value	Overall score
		Flamingo-3E	3 model	70.7
(i)	Training data	All data	w/o Video-Text pairs w/o Image-Text pairs Image-Text pairs → LAIO w/o M3W	67.3 60.9 66.4 53.4
(ii)	Optimisation	Accumulation	Round Robin	62.9
(iii)	Tanh gating	✓	Х	66.5
	Cross attention		VANILLA XATTN	66.9
(iv)	Cross-attention architecture	GATED XATTN-DENSE	GRAFTING	63.1
(iv) (v)				



Failures: Hallucinations

Input Prompt



Question: What is on the phone screen? Answer:



Question: What can you see out the window? Answer:



Question: Whom is the person texting? Answer:

P Output

A text message from a friend.

A parking lot.

The driver.

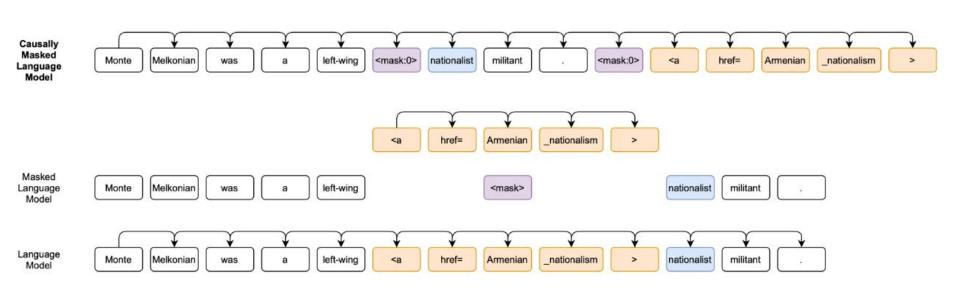


Survey of Visual LMs



CM₃

- Causally Masked Multimodal Modeling
- Images tokenized by VQVAE-GAN (source: https://arxiv.org/abs/2012.09841)

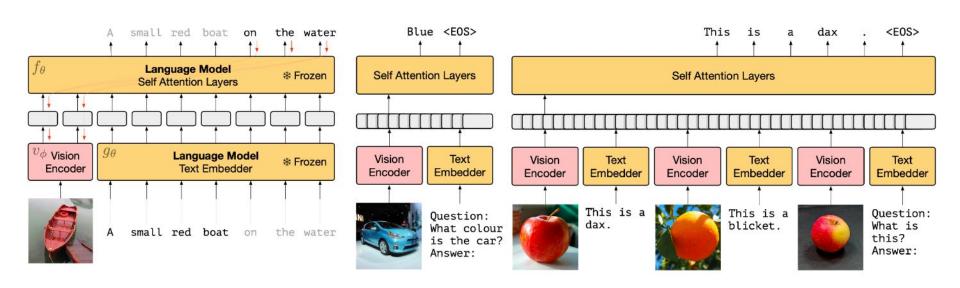


Paper: https://arxiv.org/abs/2201.07520



Learning Image Embeddings on Frozen LM Prefix

Multimodal few shot learning for interleaved vision and text



Paper: https://arxiv.org/abs/2106.13884



Discussion

If you are going to build a visual LM for few-shot learning, what are the other ways of fusing visual and textual data? What pre-training data would you consider?