

Flamingo: a Visual Language Model for Few-Shot Learning




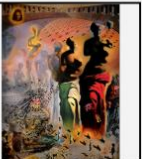

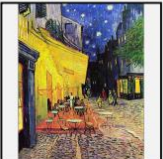



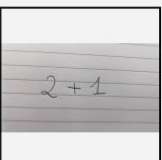
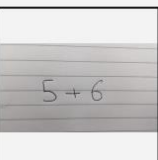
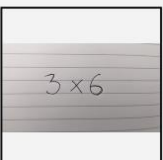

Jean-Baptiste Alayrac^{*,‡}, Jeff Donahue^{*}, Pauline Luc^{*}, Antoine Miech^{*}, Iain Barr[†], Yana Hasson[†],
Karel Lenc[†], Arthur Mensch[†], Katie Millican[†], Malcolm Reynolds[†], Roman Ring[†], Eliza Rutherford[†],
Serkan Cabi, Tengda Han, Zhitao Gong, Sina Samangooei, Marianne Monteiro, Jacob Menick,
Sebastian Borgeaud, Andrew Brock, Aida Nematzadeh, Sahand Sharifzadeh, Mikolaj Binkowski,
Ricardo Barreira, Oriol Vinyals, Andrew Zisserman, Karen Simonyan^{*,‡}

^{*}Equal contributions, ordered alphabetically, [†]Equal contributions, ordered alphabetically, [‡]Equal senior contributions



Flamingo: a Visual Language Model for Few-Shot Learning

Input Prompt

	This is a chinchilla. They are mainly found in Chile.		This is a shiba. They are very popular in Japan.		This is
	What is the title of this painting? Answer: The Hallucinogenic Toreador.		Where is this painting displayed? Answer: Louvres Museum, Paris.		What is the name of the city where this was painted? Answer:
	Output: "Underground"		Output: "Congress"		Output:
	$2+1=3$		$5+6=11$		
					What happens to the man after hitting the ball? Answer:

Completion



a flamingo. They are found in the Caribbean and South America.

Arles.

"Soulomes"

$3 \times 6 = 18$

he falls down.



This is a picture of two teddy bears on the moon.

What are they doing?

They are having a conversation.

What object are they using?

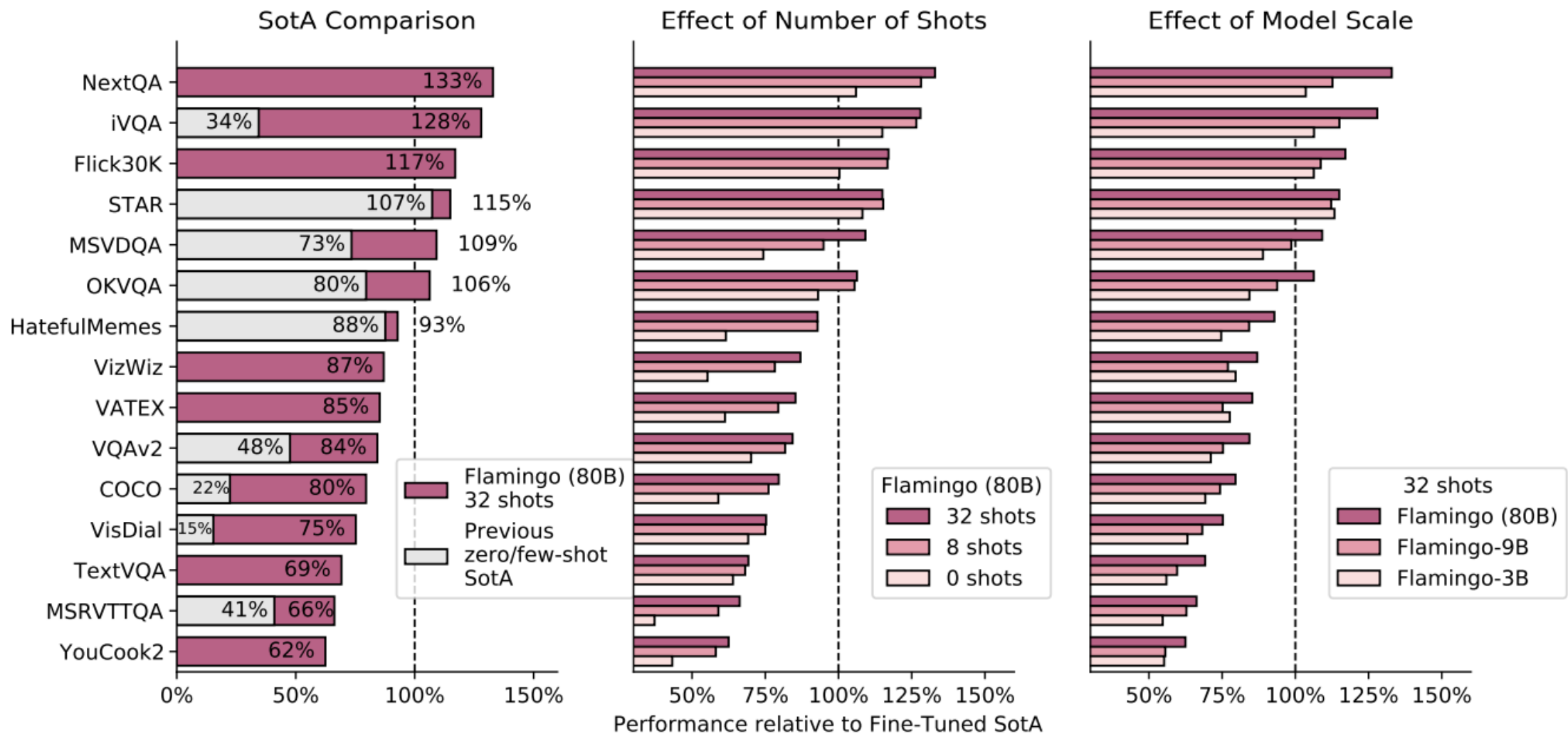
It looks like a computer.

Is this surprising?

Yes, it is surprising.

Why is this picture surprising to you?

I think it is surprising because teddy bears are not usually found on the moon.





Flamingo: a Visual Language Model for Few-Shot Learning

Dominant computer vision paradigm:

Large-scale pretraining

+

Task-specific fine-tuning

Limitation:

- Thousands of training samples
- Careful per-task hyperparameter tuning
- Significant computational resource

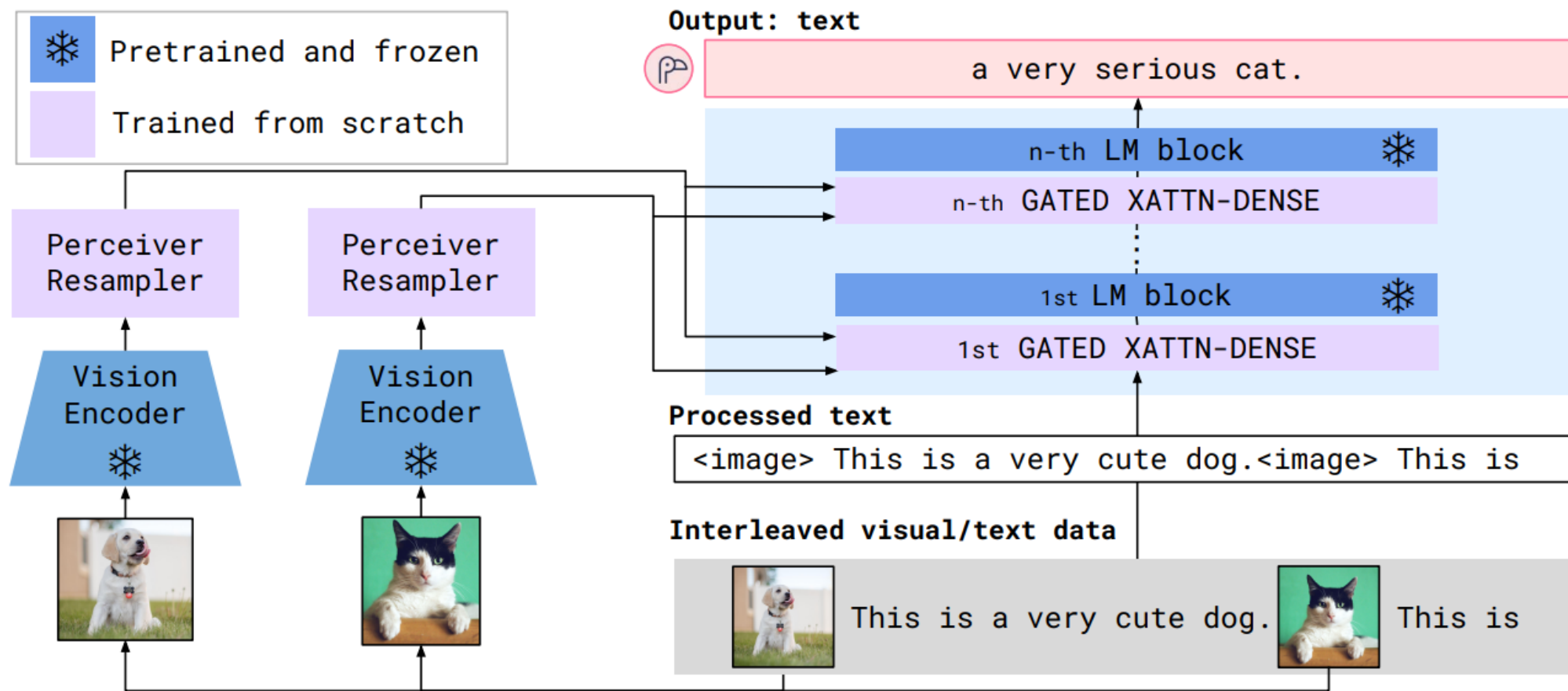
Current multimodal models (CLIP, VL-T5)

- Show promising zero-shot performance
- The trade-off between flexibility and few-shot performance

Challenges and contributions

- Flamingo
 - Using strong unimodal models
 - Fixed pretrained language model and vision model
 - Train fusion layers.
 - Supporting images and videos
 - Perceiver-based architecture with a fixed number of visual tokens
 - Heterogeneous training data
 - Combined web scraping with existing paired (image, text) and (video, text) datasets

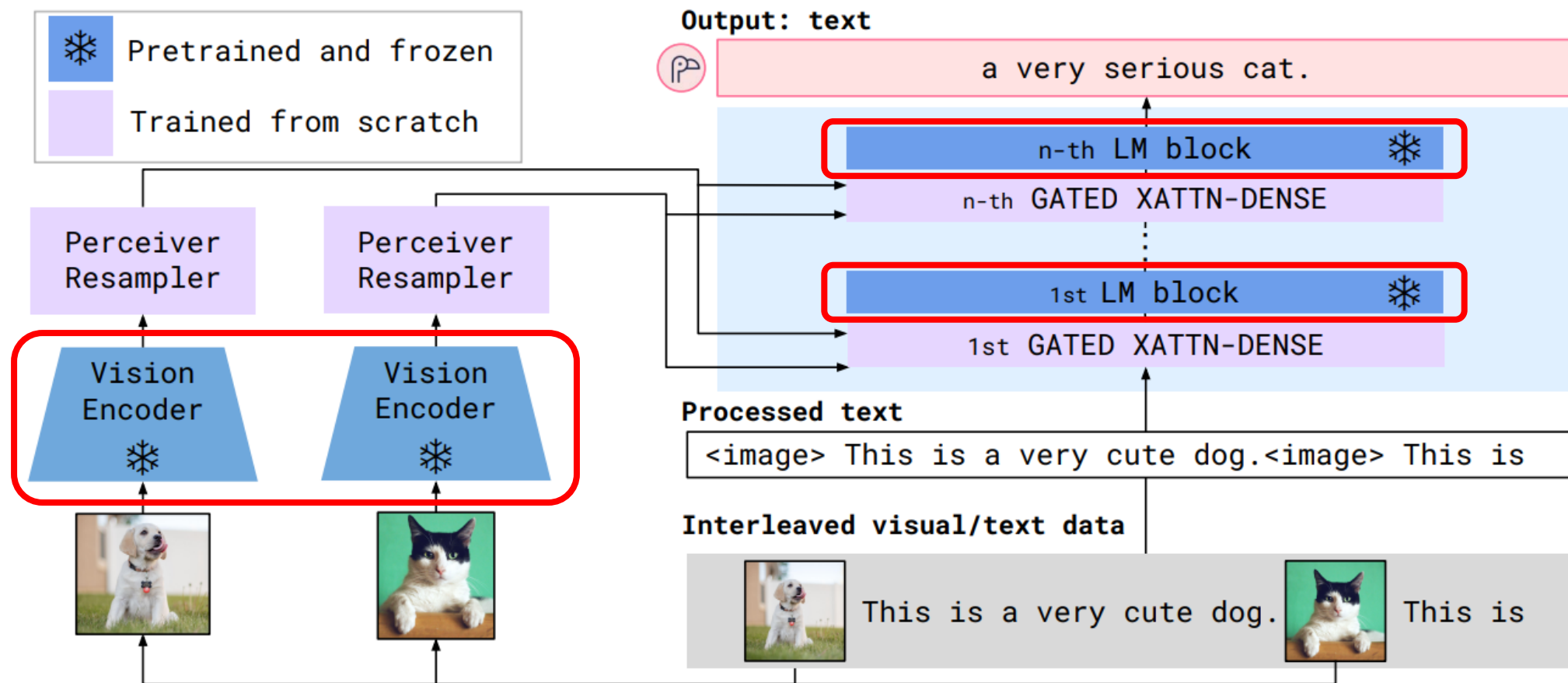
Overview



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

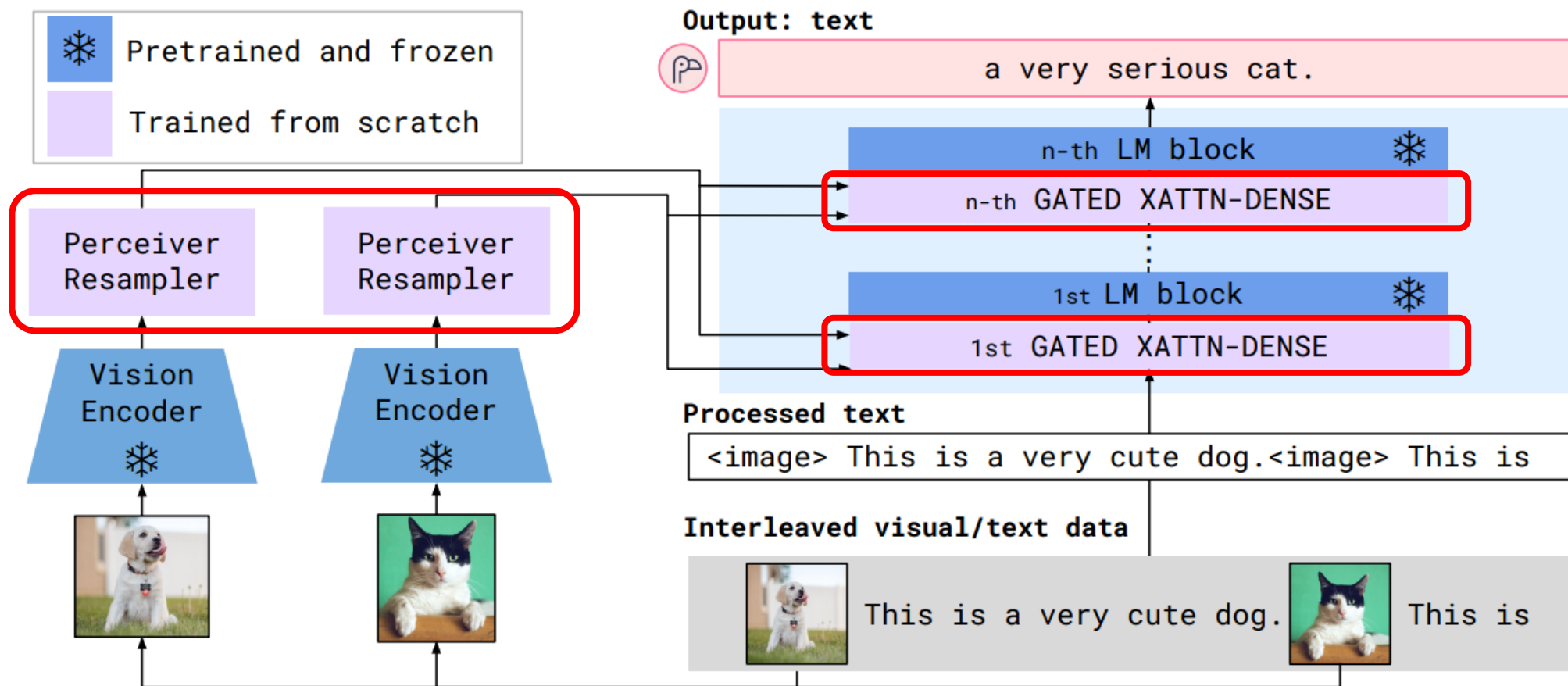
Overview

Leverage pretrained models to save compute



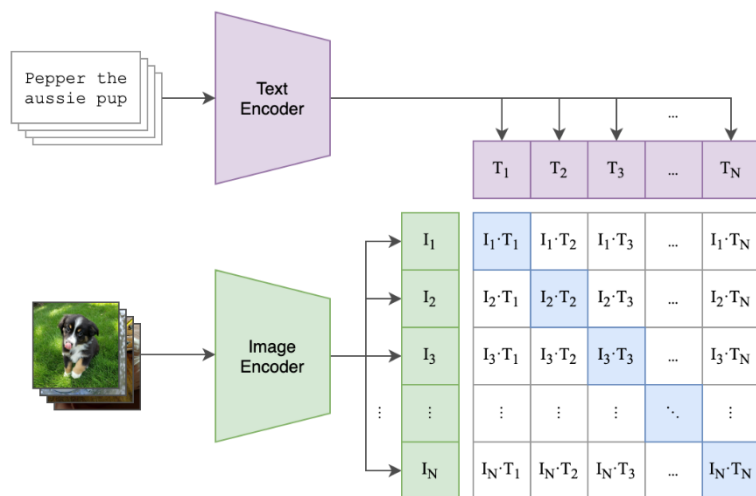
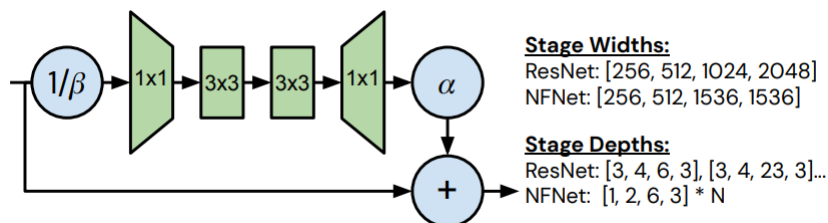
Overview

Bridge pretrained models harmoniously



Vision Encoder

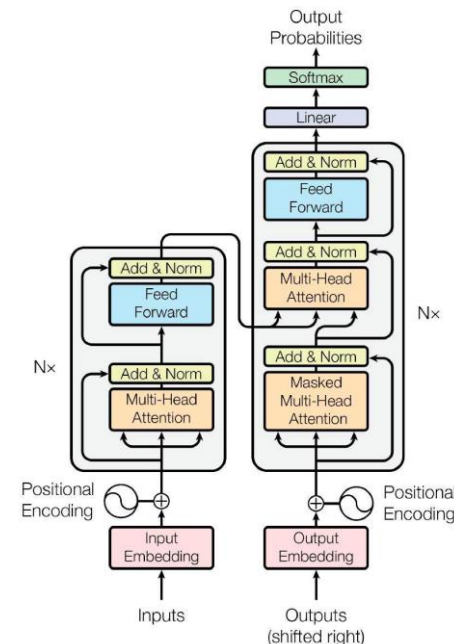
NFNet



Contrastive training

- Brock, A., De, S., Smith, S. L., & Simonyan, K. (2021, July). High-performance large-scale image recognition without normalization. In *International Conference on Machine Learning* (pp. 1059-1071). PMLR.
- Radford, A., Kim, J. W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., ... & Sutskever, I. (2021, July). Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning* (pp. 8748-8763). PMLR.

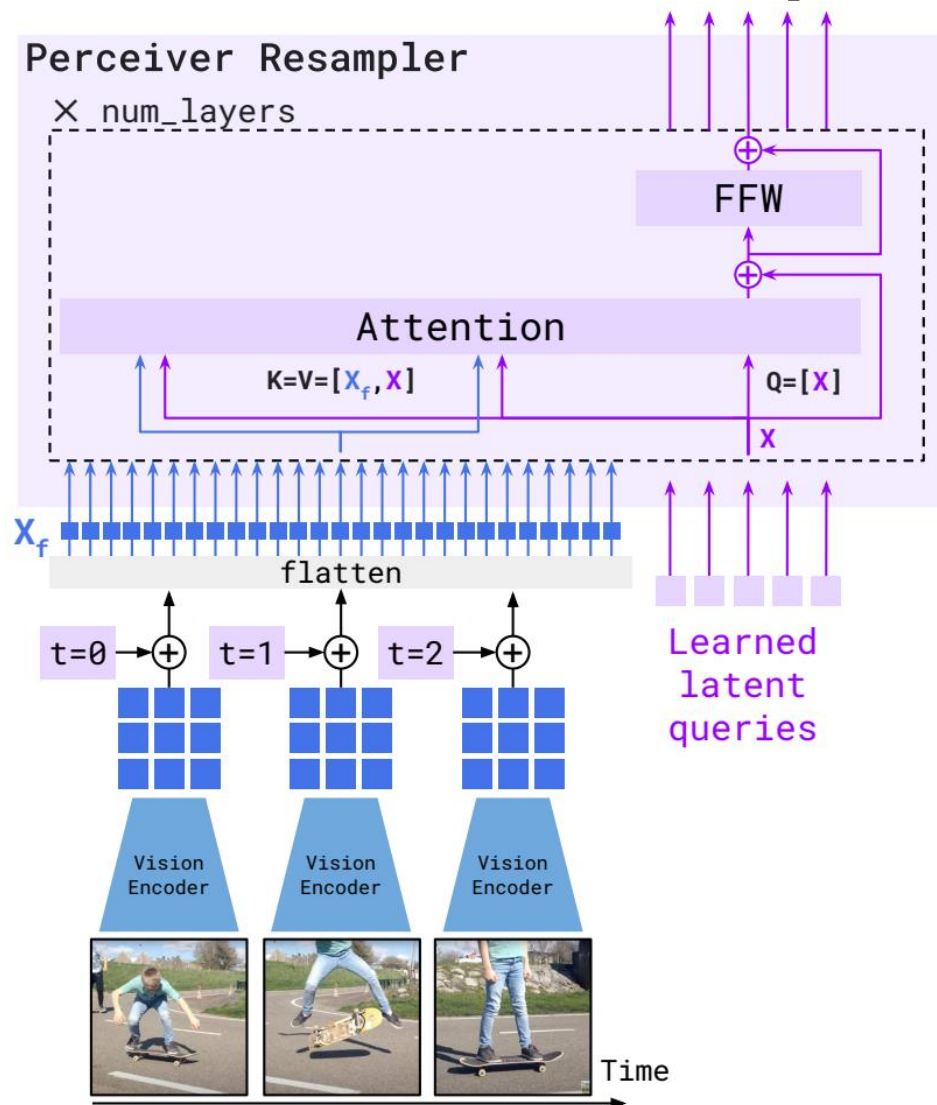
Pretrained language model



Model	Size (# Parameters)	Training Tokens
LaMDA (Thoppilan et al., 2022)	137 Billion	168 Billion
GPT-3 (Brown et al., 2020)	175 Billion	300 Billion
Jurassic (Lieber et al., 2021)	178 Billion	300 Billion
Gopher (Rae et al., 2021)	280 Billion	300 Billion
MT-NLG 530B (Smith et al., 2022)	530 Billion	270 Billion
Chinchilla	70 Billion	1.4 Trillion

- Hoffmann, J., Borgeaud, S., Mensch, A., Buchatskaya, E., Cai, T., Rutherford, E., ... & Sifre, L. (2022). Training Compute-Optimal Large Language Models. arXiv preprint arXiv:2203.15556.

Perceiver Resampler

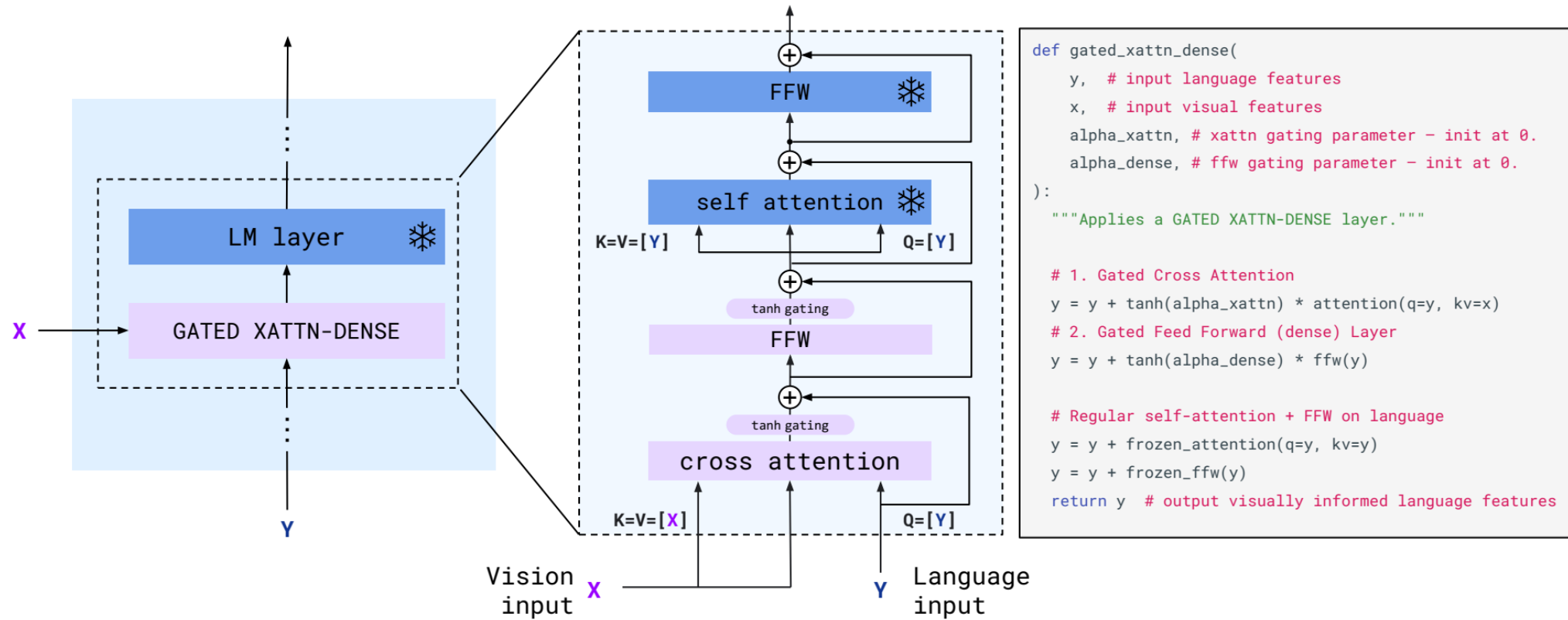


```
def perceiver_resampler(
    x_f, # The [T, S, d] visual features (T=time, S=space)
    time_embeddings, # The [T, 1, d] time pos embeddings.
    x, # R learned latents of shape [R, d]
    num_layers, # Number of layers
):
    """The Perceiver Resampler model."""

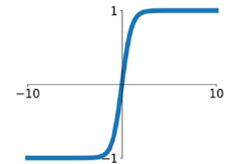
    # Add the time position embeddings and flatten.
    x_f = x_f + time_embeddings
    x_f = flatten(x_f) # [T, S, d] -> [T * S, d]
    # Apply the Perceiver Resampler layers.
    for i in range(num_layers):
        # Attention.
        x = x + attention_i(q=x, kv=concat([x_f, x]))
        # Feed forward.
        x = x + ffw_i(x)
    return x
```

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

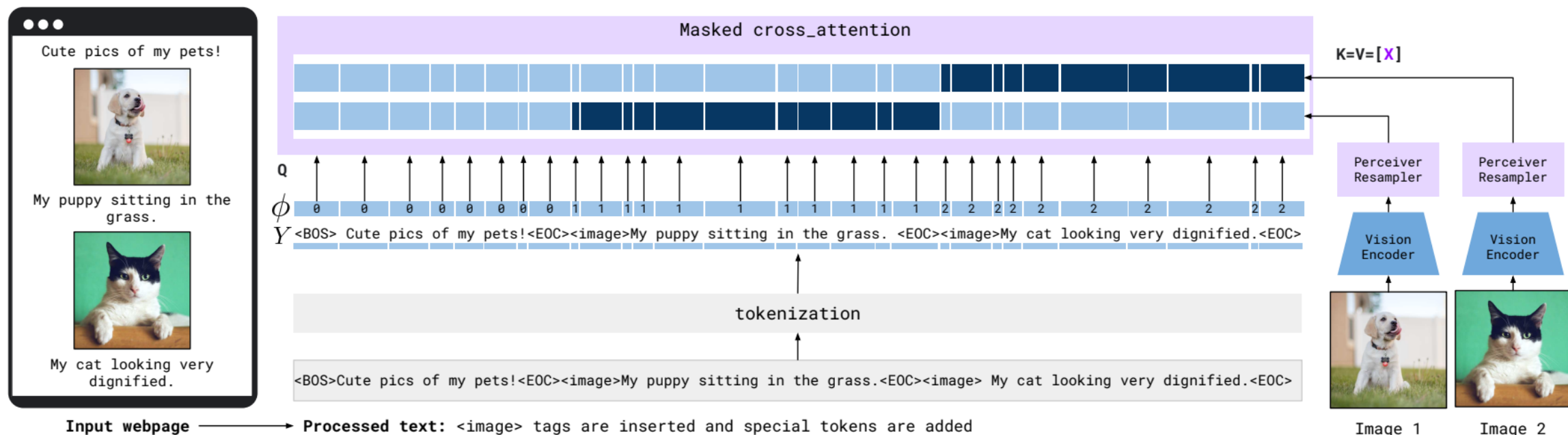
Gated xattn-dense block



tanh
tanh(x)



Per-image/video attention masking



Training data



This is an
image of a
flamingo.

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

ALIGN: 1.8B
LTIP: 312M



A kid
doing a
kickflip.

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

VTP: 27M



Welcome
to my
website!

This is a
picture of
my dog.



This is a
picture of
my cat.

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

M3W: 185M images
and 182GB of text

Training objective

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

Experiments

- Model

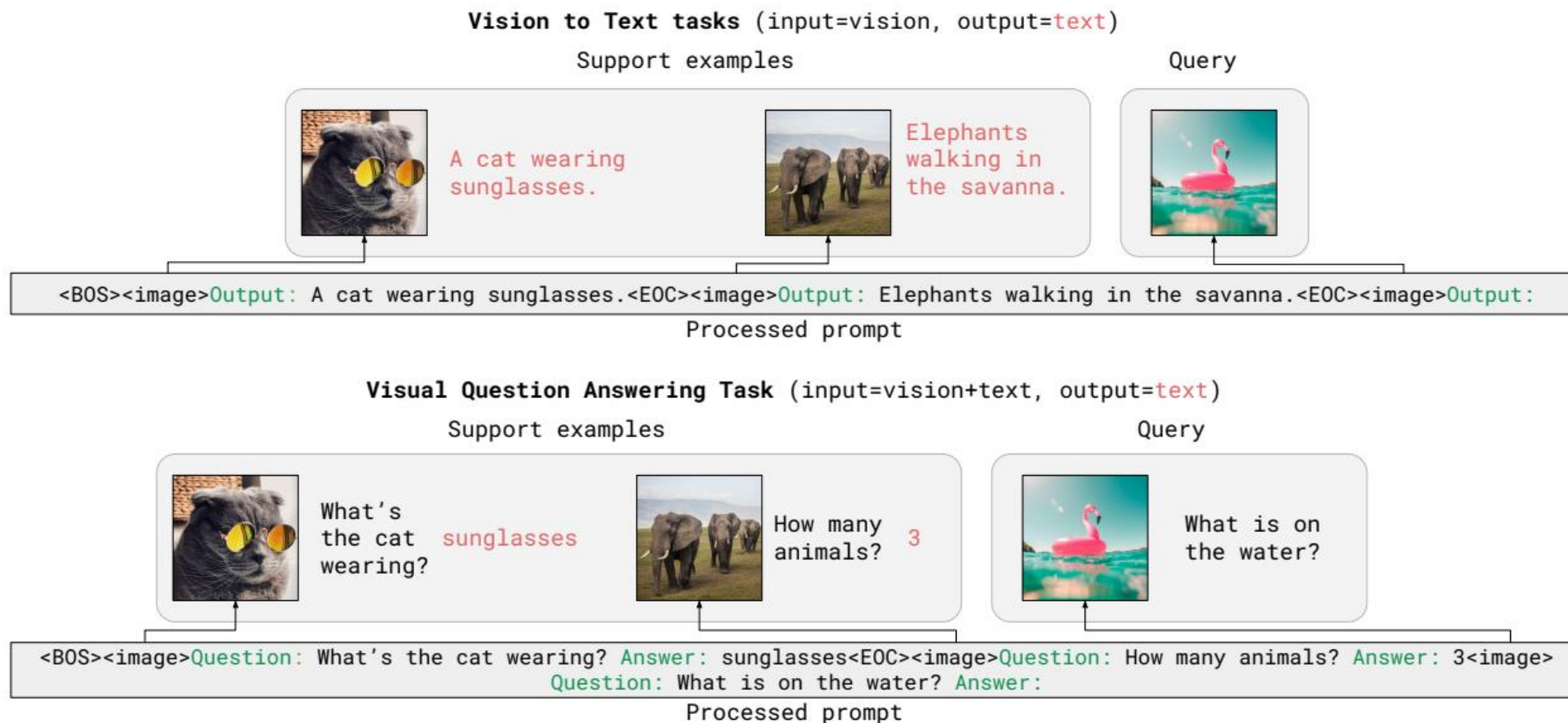
	Requires model sharding	Frozen		Trainable			Total count
		Language	Vision	GATED	XATTN-DENSE	Resampler	
<i>Flamingo</i> -3B	✗	1.4B	435M	1.2B (every)		194M	3.2B
<i>Flamingo</i> -9B	✗	7.1B	435M	1.6B (every 4th)		194M	9.3B
<i>Flamingo</i>	✓	70B	435M	10B (every 7th)		194M	80B

	Resampler				xattn dense				Frozen LM			
	L	D	H	Act.	L	D	H	Act.	L	D	H	Act.
<i>Flamingo</i> -3B	6	1536	16	Sq. ReLU	24	2048	16	Sq. ReLU	24	2048	16	GeLU
<i>Flamingo</i> -9B	6	1536	16	Sq. ReLU	10	4096	32	Sq. ReLU	40	4096	32	GeLU
<i>Flamingo</i>	6	1536	16	Sq. ReLU	12	8192	64	Sq. ReLU	80	8192	64	GeLU

1536 TPUv4 for 15 days

Experiments

- Few-shot interleaved prompt generation.





Experiments

- Few-shot: comparison to SotA

Method	FT	Shot	OKVQA	VQAv2	COCO	MSVDQA	VATEX	VizWiz	Flick30K	MSRVTTQA	iVQA	YouCook2	STAR	VisDial	TextVQA	NextQA	HatefulMemes	RareAct
Zero/Few shot SOTA	✗	(X)	[39] 43.3 (16)	[124] 38.2 (4)	[134] 32.2 (0)	[64] 35.2 (0)	-	-	-	[64] 19.2 (0)	[145] 12.2 (0)	-	[153] 39.4 (0)	[87] 11.6 (0)	-	-	[94] 66.1 (0)	[94] 40.7 (0)
Flamingo-3B	✗	0	41.2	49.2	73.0	27.5	40.1	28.9	60.6	11.0	32.7	55.8	39.6	46.1	30.1	21.3	53.7	58.4
	✗	4	43.3	53.2	85.0	33.0	50.0	34.0	72.0	14.9	35.7	64.6	41.3	47.3	32.7	22.4	53.6	-
	✗	8	44.6	55.4	90.6	37.0	54.5	38.4	71.7	19.6	36.8	68.0	40.6	47.6	32.4	23.9	54.7	-
	✗	16	45.6	56.7	95.4	40.2	57.1	43.3	73.4	23.4	37.4	73.2	40.1	47.5	31.8	25.2	55.3	-
	✗	32	45.9	57.1	99.0	42.6	59.2	45.5	71.2	25.6	37.7	76.7	41.6	OOC	30.6	26.1	56.3	-
Flamingo-9B	✗	0	44.7	51.8	79.4	30.2	39.5	28.8	61.5	13.7	35.2	55.0	41.8	48.0	31.8	23.0	57.0	57.9
	✗	4	49.3	56.3	93.1	36.2	51.7	34.9	72.6	18.2	37.7	70.8	42.8	50.4	33.6	24.7	62.7	-
	✗	8	50.0	58.0	99.0	40.8	55.2	39.4	73.4	23.9	40.0	75.0	43.4	51.2	33.6	25.8	63.9	-
	✗	16	50.8	59.4	102.2	44.5	58.5	43.0	72.7	27.6	41.5	77.2	42.4	51.3	33.5	27.6	64.5	-
	✗	32	51.0	60.4	106.3	47.2	57.4	44.0	72.8	29.4	40.7	77.3	41.2	OOC	32.6	28.4	63.5	-
Flamingo	✗	0	50.6	56.3	84.3	35.6	46.7	31.6	67.2	17.4	40.7	60.1	39.7	52.0	35.0	26.7	46.4	60.8
	✗	4	57.4	63.1	103.2	41.7	56.0	39.6	75.1	23.9	44.1	74.5	42.4	55.6	36.5	30.8	68.6	-
	✗	8	57.5	65.6	108.8	45.5	60.6	44.8	78.2	27.6	44.8	80.7	42.3	56.4	37.3	32.3	70.0	-
	✗	16	57.8	66.8	110.5	48.4	62.8	48.4	78.9	30.0	45.2	84.2	41.1	56.8	37.6	32.9	70.0	-
	✗	32	57.8	67.6	113.8	52.3	65.1	49.8	75.4	31.0	45.3	86.8	42.2	OOC	37.9	33.5	70.0	-
Pretrained FT SOTA	✓	(X)	54.4 [39] (10K)	80.2 [150] (444K)	143.3 [134] (500K)	47.9 [32] (27K)	76.3 [165] (500K)	57.2 [70] (20K)	67.4 [162] (30K)	46.8 [57] (130K)	35.4 [145] (6K)	138.7 [142] (10K)	36.7 [138] (46K)	75.2 [87] (123K)	54.7 [147] (20K)	25.2 [139] (38K)	75.4 [60] (9K)	-

Experiments

- Fine-tuning Flamingo: Flamingo + base visual encoder

Method	VQAV2		COCO	VATEX	VizWiz		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
 <i>Flamingo</i> - 32 shots	67.6	-	113.8	65.1	49.8	-	31.0	56.8	-	86.8	36.0	-	70.0
SimVLM [134]	80.0	80.3	143.3	-	-	-	-	-	-	-	-	-	-
OFA [129]	79.9	80.0	<u>149.6</u>	-	-	-	-	-	-	-	-	-	-
Florence [150]	80.2	80.4	-	-	-	-	-	-	-	-	-	-	-
 <i>Flamingo</i> Fine-tuned	<u>82.0</u>	<u>82.1</u>	138.1	<u>84.2</u>	<u>65.7</u>	<u>65.4</u>	<u>47.4</u>	61.8	59.7	118.6	<u>57.1</u>	54.1	<u>86.6</u>
Restricted SotA [†]	80.2	80.4	143.3	76.3	-	-	46.8	<u>75.2</u>	74.5	<u>138.7</u>	54.7	<u>73.7</u>	75.4
	[150]	[150]	[134]	[165]	-	-	[57]	[87]	[87]	[142]	[147]	[92]	[60]
Unrestricted SotA	81.3	81.3	<u>149.6</u>	81.4	57.2	60.6	-	-	<u>75.4</u>	-	-	-	84.6
	[143]	[143]	[129]	[165]	[70]	[70]	-	-	[133]	-	-	-	[164]

Experiments

- Ablation studies




	Ablated setting	Flamingo 3B value	Changed value	Param. count ↓	Step time ↓	COCO CIDEr↑	OKVQA top1↑	VQAv2 top1↑	ImageNet top1↑	MSVDQA top1↑	VATEX CIDEr↑	Kinetics top1-top5↑	Overall score↑
	Flamingo 3B model (short training)			3.2B	1.74s	86.5	42.1	55.8	59.9	36.3	53.4	49.4	68.4
(i)	Training data	All data	M3W	3.2B	0.68s	58.0	37.2	48.6	35.7	29.5	33.6	34.0	50.7
			w/o VTP	3.2B	1.42s	84.2	43.0	53.9	59.6	34.5	46.0	45.8	65.4
			w/o LTIP/ALIGN	3.2B	0.95s	66.3	39.2	51.6	41.4	32.0	41.6	38.2	56.5
			w/o M3W	3.2B	1.02s	54.1	36.5	52.7	24.9	31.4	23.5	28.3	46.9
(ii)	Optimisation	Grad. accumulation	Round Robin	3.2B	1.68s	76.1	39.8	52.1	50.7	33.2	40.8	39.7	59.7
(iii)	Tanh gating	✓	✗	3.2B	1.74s	78.4	40.5	52.9	54.0	35.9	47.5	46.4	64.0
(iv)	Cross-attention architecture	GATED XATTN-DENSE	VANILLA XATTN GRAFTING	2.4B	1.16s	80.6	41.5	53.4	59.0	32.9	50.7	46.8	65.2
				3.3B	1.74s	79.2	36.1	50.8	47.5	32.2	47.8	27.9	57.4
(v)	Cross-attention frequency	Every	Single in middle	2.0B	0.87s	71.5	38.1	50.2	44.0	29.1	42.3	28.3	54.6
			Every 4th	2.3B	1.02s	82.3	42.7	55.1	57.1	34.6	50.8	45.5	65.9
			Every 2nd	2.6B	1.24s	83.7	41.0	55.8	59.6	34.5	49.7	47.4	66.2
(vi)	Resampler	Perceiver	MLP	3.2B	1.85s	78.6	42.2	54.7	53.6	35.2	44.7	42.1	63.3
			Transformer	3.2B	1.81s	83.2	41.7	55.6	59.0	31.5	48.3	47.4	65.1
(vii)	Resampler size	Medium	Small	3.1B	1.58s	81.1	40.4	54.1	60.2	36.0	50.2	48.9	66.4
			Large	3.4B	1.87s	84.4	42.2	54.4	60.4	35.1	51.4	49.4	67.3

Experiments

- Ablation studies

	Ablated setting	Flamingo 3B value	Changed value	Param. count ↓	Step time ↓	COCO CIDEr↑	OKVQA top1↑	VQAv2 top1↑	ImageNet top1↑	MSVDQA top1↑	VATEX CIDEr↑	Kinetics top1-top5↑	Overall score↑
	Flamingo 3B model (short training)			3.2B	1.74s	86.5	42.1	55.8	59.9	36.3	53.4	49.4	68.4
(viii)	Multi-Img att.	Only last	All previous	3.2B	1.74s	70.0	40.9	52.0	52.3	32.1	46.8	42.0	60.8
(ix)	p_{next}	0.5	0.0	3.2B	1.74s	85.0	41.6	55.2	60.3	36.7	50.6	49.9	67.8
			1.0	3.2B	1.74s	81.3	43.3	55.6	57.8	36.8	52.7	47.8	67.6
(x)	Vision encoder	NFNet-F6	CLIP ViT-L/14	3.1B	1.58s	76.5	41.6	53.4	49.5	33.2	44.5	42.3	61.4
			NFNet-F0	2.9B	1.45s	73.8	40.5	52.8	49.8	31.1	42.9	36.6	58.9
(xi)	LM pretraining	MassiveText	C4	3.2B	1.74s	81.3	34.4	47.1	60.6	30.9	53.9	46.9	62.5
(xii)	Freezing Vision	✓	✗ (random init)	3.2B	4.70s*	74.5	41.6	52.7	45.2	31.4	35.8	32.6	56.6
			✗ (pretrained)	3.2B	4.70s*	83.5	40.6	55.1	55.6	34.6	50.7	41.2	64.5
(xiii)	Freezing LM	✓	✗ (random init)	3.2B	2.42s	74.8	31.5	45.6	59.5	26.9	50.1	43.4	58.2
			✗ (pretrained)	3.2B	2.42s	81.2	33.7	47.4	60.7	31.0	53.9	49.9	62.9
(xiv)	Co-train LM on MassiveText	✗	✓ (random init)	3.2B	5.34s*	69.3	29.9	46.1	59.9	28.1	45.5	46.9	57.4
			✓ (pretrained)	3.2B	5.34s*	83.0	42.5	53.3	60.9	35.1	51.1	50.1	67.2

Experiments

Input Prompt	 <p>Question: What do you think the capacities of these are? Answer:</p>	 <p>Question: What is odd about this image? Explain why it is unusual. Answer:</p>	 <p>Question: What country is this? Why do you think so? Answer:</p>
Completion	<p>The floppy disk is 1.44MB and the CD is 700MB.</p>	<p>The image is odd because the elephant is in the back of the truck. It is unusual because elephants are not usually transported in the back of a truck.</p>	<p>It is Canada. I think so because the flag is the Canadian flag.</p>

Experiments



This is a dog. It's a white fluffy dog.

What season is this?

It looks like a summer day.

Why do you think so?

The grass is green and the dog is not wearing a sweater.

What breed is the dog?


It's a Samoyed.

Is the dog running?

No, it's sitting.

Can you describe the pose of its back legs?

The back legs are bent and the dog is sitting on its haunches.



Which city is this?

This is a street sign in New York City.

Which street corner is it?

It is the intersection of Madison Avenue and East 42nd Street.

Which subway station is nearby?

The closest subway station is Grand Central Station.

How about the nearest park?

The nearest park is Bryant Park.

There is another sign. What does it say?

It says "One way".

I would like you to read the color and then say the color it is written in. For example:

PURPLE

Color is "Purple" and it is written in red.

GREEN

Color is "Green" and it is written in blue.

YELLOW

Color is "Yellow" and it is written in green.

BLACK

Color is "Black" and it is written in yellow.

Well done! Do you know the name of the test these images come from?

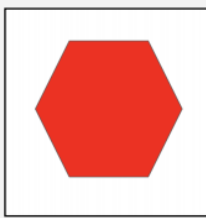
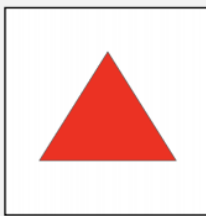
I think it is called the Stroop test.

Can you explain how humans perform in this test?

Humans are slower when the color of the word and the color of the word are different.

How about you?

I am not affected by this difference.



What is the difference between these two images?

The first one is a triangle and the second one is a hexagon.


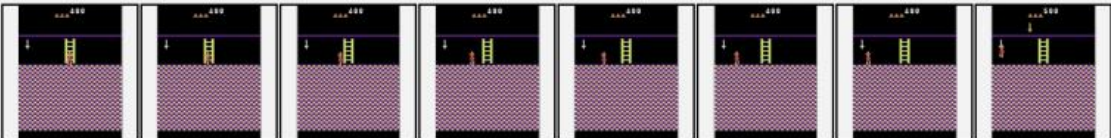




Which one has larger internal angles?

The hexagon has larger internal angles.

What color are they?

They are both red.

Experiments

Input Prompt								Completion
	Question: What is happening here? Answer:							The dachshund puppy is being weighed on a scale.
	Question: What object is the avatar picking up? Answer:							A sword.
	Question: What is the object being lifted? Answer:							The object is a small plastic bowl.
	Question: What are the three objects in this video? Answer:							A scissors, a pen, and a rubber band.
	Question: What is written here? Answer:							Flamingo.
	What happens to the man after hitting the ball? Answer:							he falls down.




Experiments

- Limitation on Classification performance.

Model	Method	Prompt size	shots/class	ImageNet top 1	Kinetics700 avg top1/5
SotA	Fine-tuned	-	full	91.0 [137]	89.0 [144]
SotA	Contrastive	-	0	85.7 [90]	69.6 [94]
NFNetF6	Our contrastive	-	0	77.9	62.9
<i>Flamingo-3B</i>	RICES	8	1	70.9	55.9
		16	1	71.0	56.9
		16	5	72.7	58.3
<i>Flamingo-9B</i>	RICES	8	1	71.2	58.0
		16	1	71.7	59.4
		16	5	75.2	60.9
	Random	16	≤ 0.02	66.4	51.2
<i>Flamingo-80B</i>	RICES	8	1	71.9	60.4
		16	1	71.7	62.7
		16	5	76.0	63.5
	RICES+ensembling	16	5	77.3	64.2

Experiments

- Limitation

Input Prompt	 <p>Question: What is on the phone screen? Answer:</p>	 <p>Question: What can you see out the window? Answer:</p>	 <p>Question: Whom is the person texting? Answer:</p>
Output	<p>A text message from a friend.</p>	<p>A parking lot.</p>	<p>The driver.</p>