

---

---

# Few and zero shot learning

— Flamingo & Kosmos-1 —

CS 6804: Multimodal Vision | Deval Srivastava

---

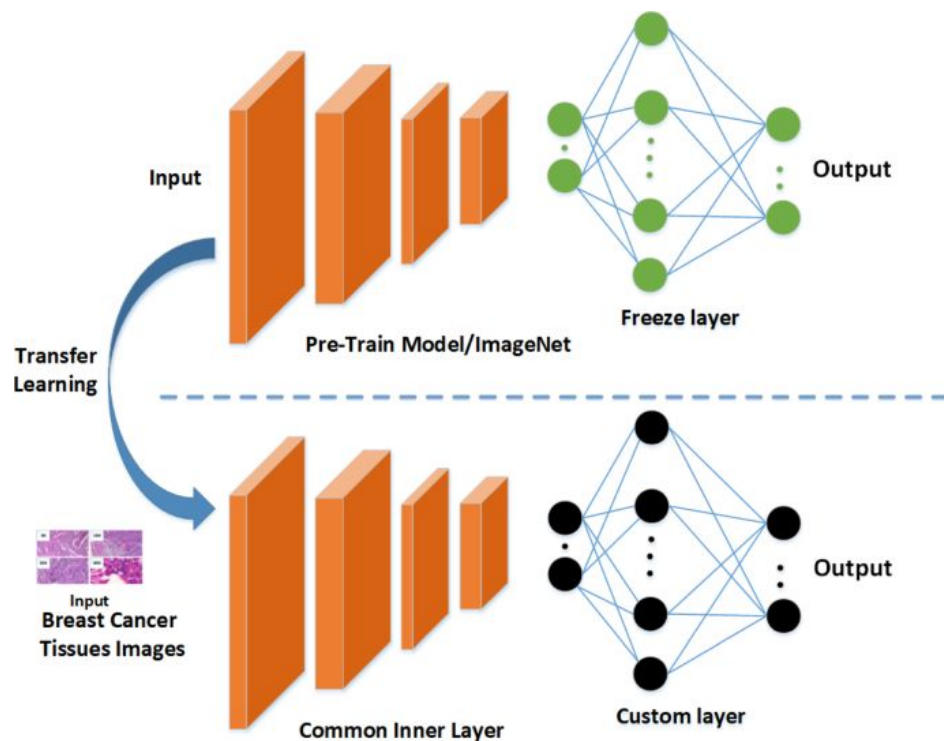
---

# Background

- Few and zero shot learning can be seen as a measure of intelligence.

# Background

- Few and zero shot learning can be seen as a measure of intelligence.
- Different from how currently models learn.



# Background

- Few and zero shot learning can be seen as a measure of intelligence.
- Different from how currently models learn.
- LLMs are able to do this to some extent[1]

## Zero-shot

The model predicts the answer given only a natural language description of the task. No gradient updates are performed.

```
1 Translate English to French: ← task description
2 cheese => ..... ← prompt
```

## One-shot

In addition to the task description, the model sees a single example of the task. No gradient updates are performed.

```
1 Translate English to French: ← task description
2 sea otter => loutre de mer ← example
3 cheese => ..... ← prompt
```

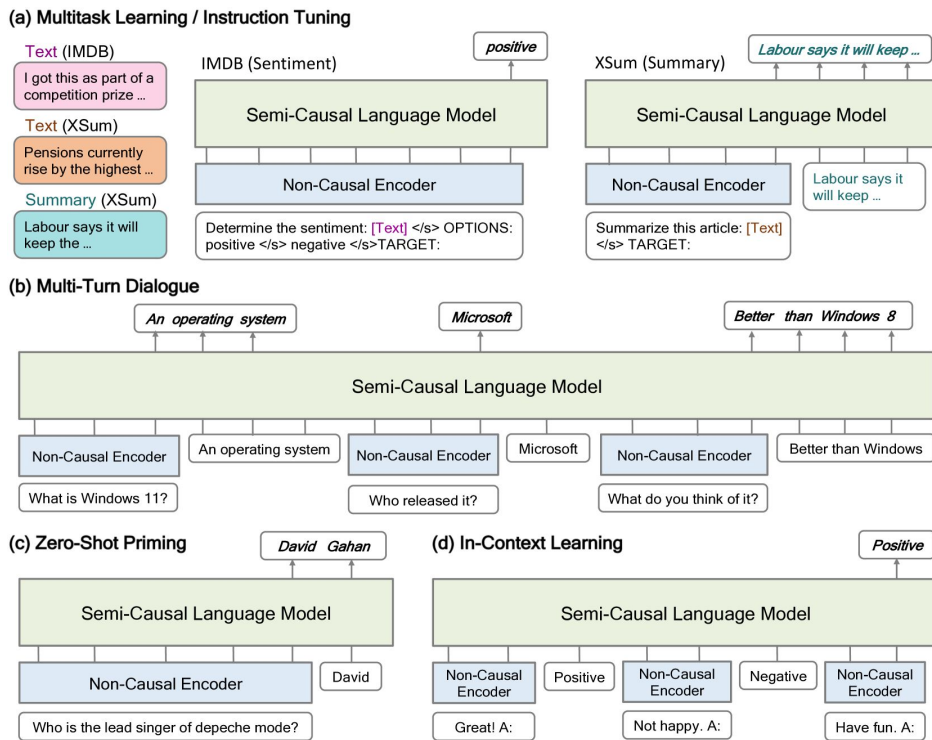
## Few-shot

In addition to the task description, the model sees a few examples of the task. No gradient updates are performed.

```
1 Translate English to French: ← task description
2 sea otter => loutre de mer ← examples
3 peppermint => menthe poivrée ←
4 plush girafe => girafe peluche ←
5 cheese => ..... ← prompt
```

# Background

- Few and zero shot learning can be seen as a measure of intelligence.
- Different from how currently models learn.
- LLMs are able to do this to some extent[1].
- LLMs can do a number of tasks through their versatile text interface.



# Motivation

## GOAL:

- Design a multimodal LLM that can perform effective few shot and zero shot learning from prompts.
- This multimodal LLM would be trained on a variety of sources including interleaved images and text.

Two papers follow this methodology

- FLAMINGO[3]
- Kosmos-1[4]



# Flamingo: a Visual Language Model for Few-Shot Learning

# Flamingo: Key Takeaways

- A new large VLM that can ingest a sequence of text/image or interleaved tokens then output text.
- Sets a new state of the art on a variety of V+L tasks by being prompted with few input / output samples
- Introduces a novel architecture that bridges two frozen pretrained vision and language models.

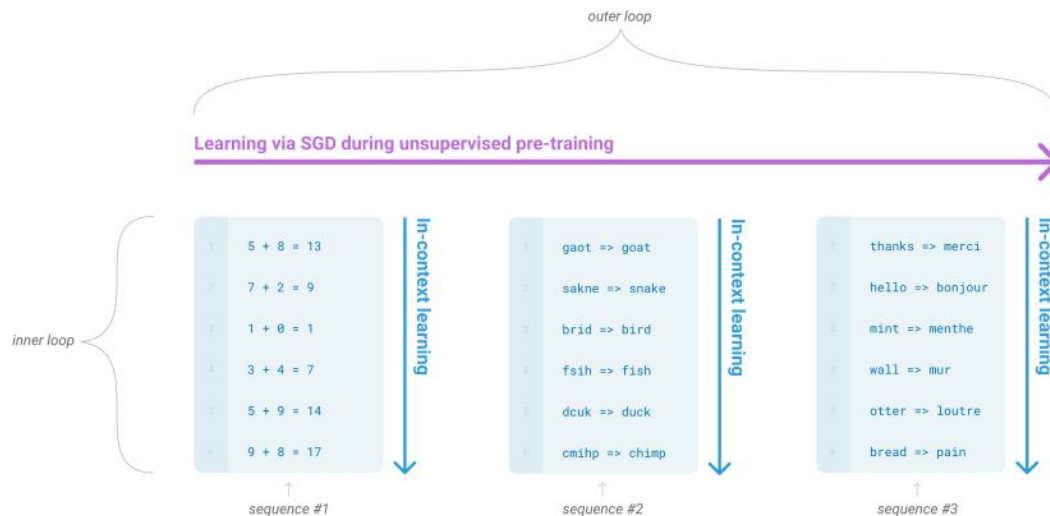


# Related Work

- **Chinchilla[5]**
  - Finetuning a LLM has become an effective strategy to use it for downstream tasks.

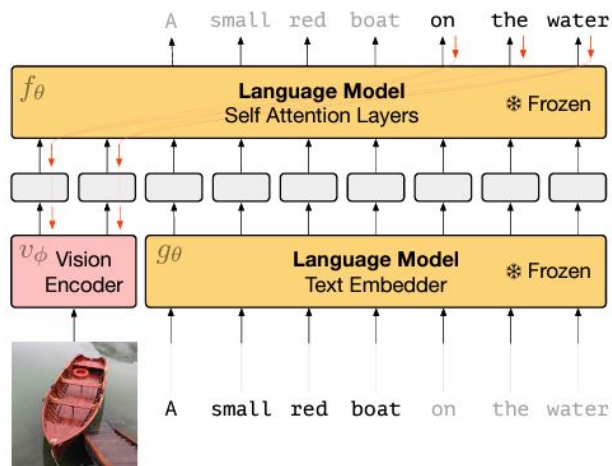
# Related Work

- Chinchilla[5]
- GPT-3[1]
  - Introduces a few shot in-context learning technique.



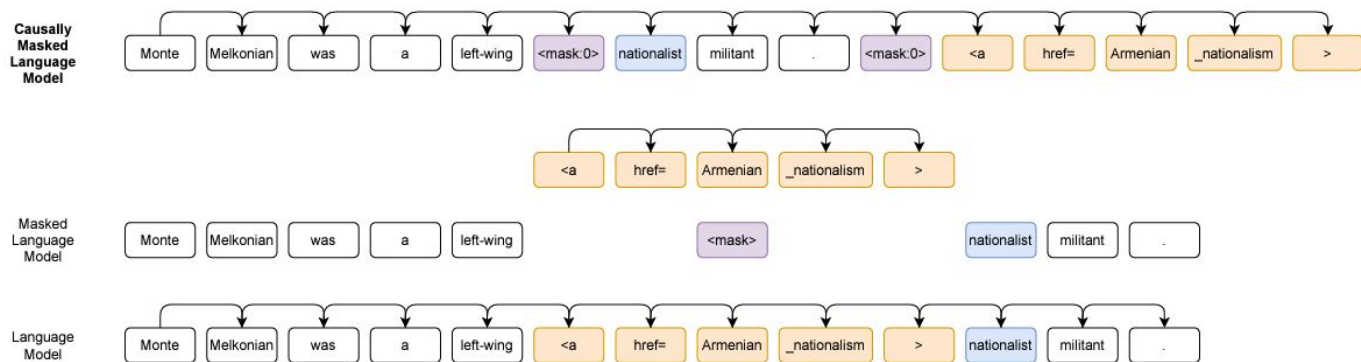
# Related Work

- Chinchilla[5]
- GPT-3[1]
- **Multimodal Few-Shot Learning with Frozen Language Models[6]**
  - Proposes to train frozen LLMs with few learnable layers on interleaved data for V+L tasks.



# Related Work

- Chinchilla[5]
- GPT-3[1]
- Multimodal Few-Shot Learning with Frozen Language Models[6]
- CM3[7]
  - Proposes to train a masked LLM on extracted HTML data for language tasks



# Flamingo:

Interleaved visual/text data



This is a very cute dog.



This is

# Flamingo Overview

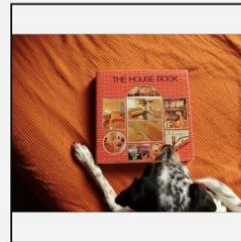
Input Prompt



Question: What is happening here? Answer:

Input Styles

Input Prompt



Question: What is the title of the book? Answer:

Interleaved visual/text data



This is a very cute dog.



This is

# Flamingo Overview

**Processed text**

<image> This is a very cute dog.<image> This is

**Interleaved visual/text data**



This is a very cute dog.



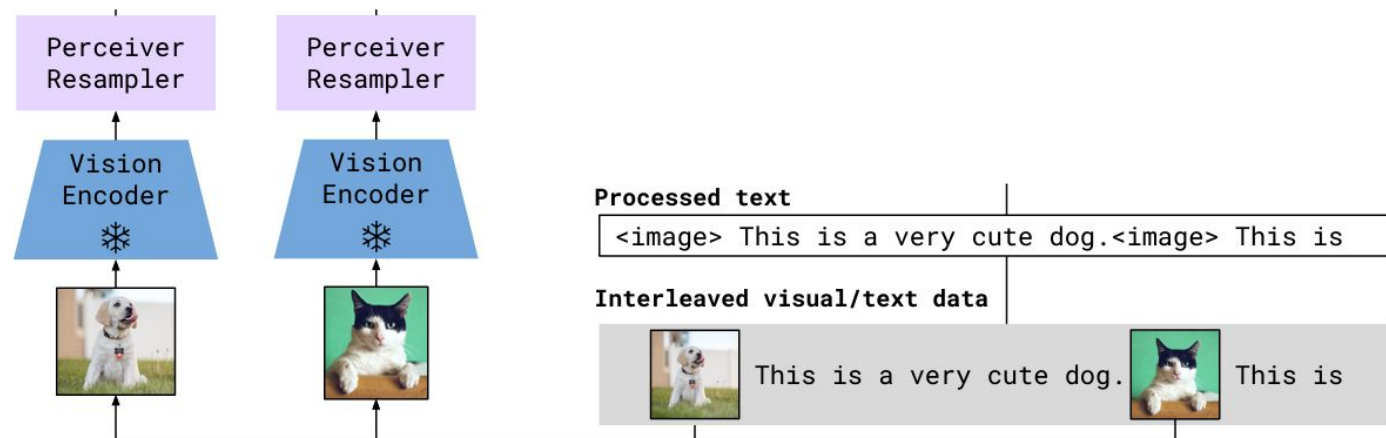
This is

# Flamingo Overview

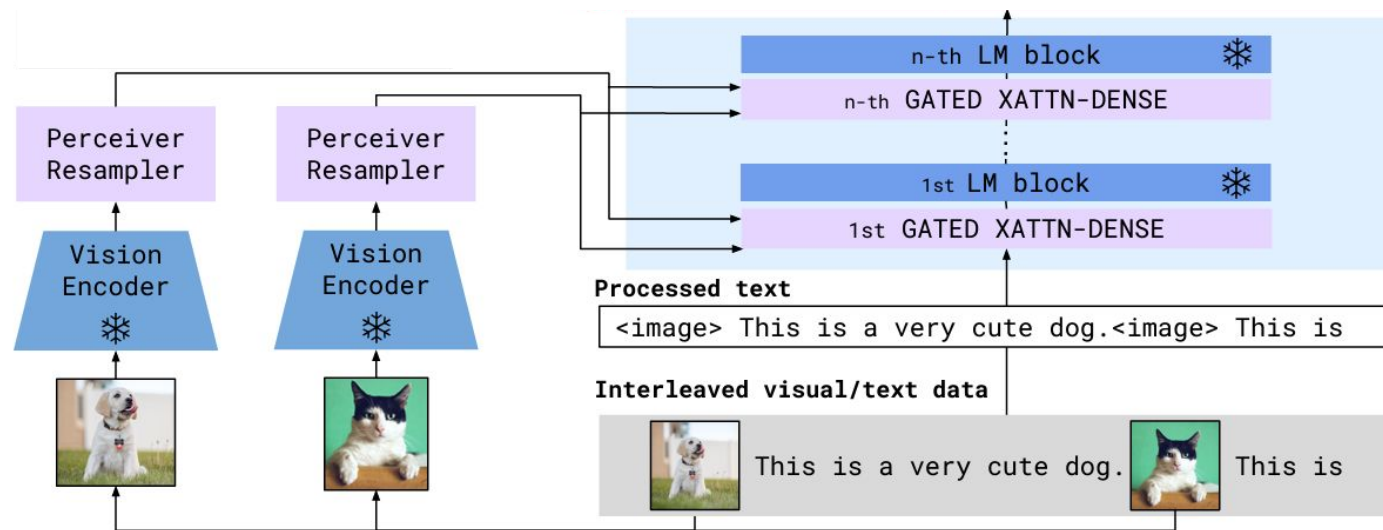




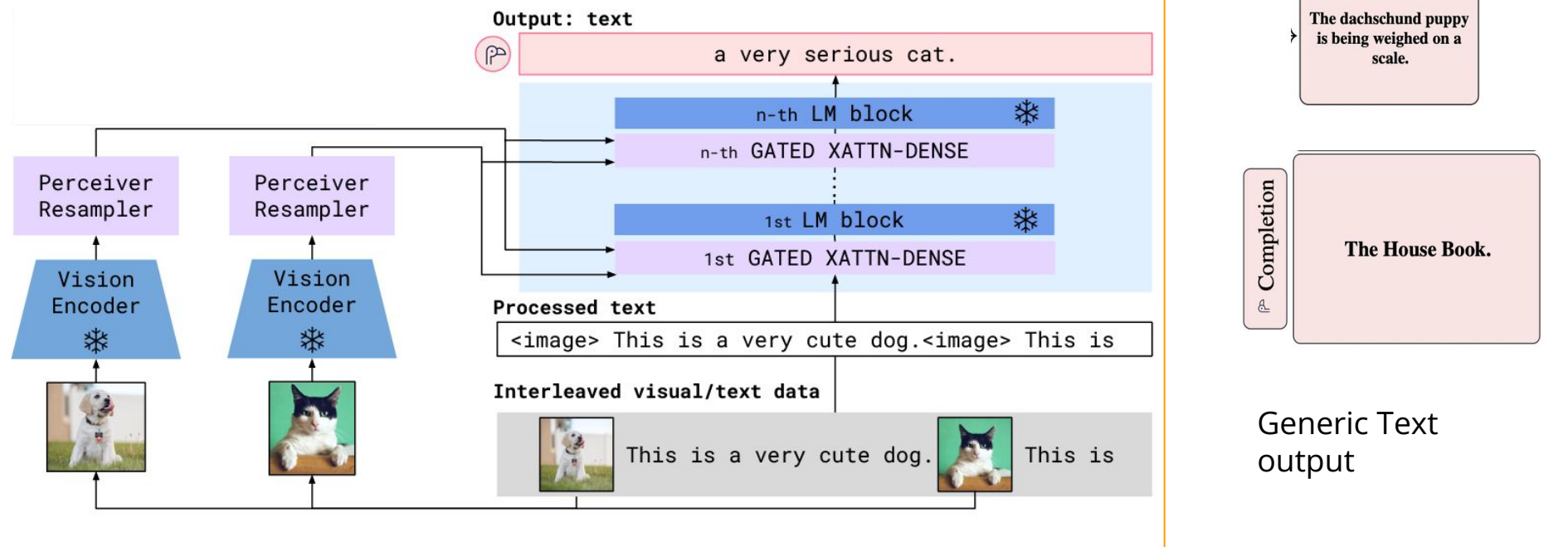
# Flamingo Overview



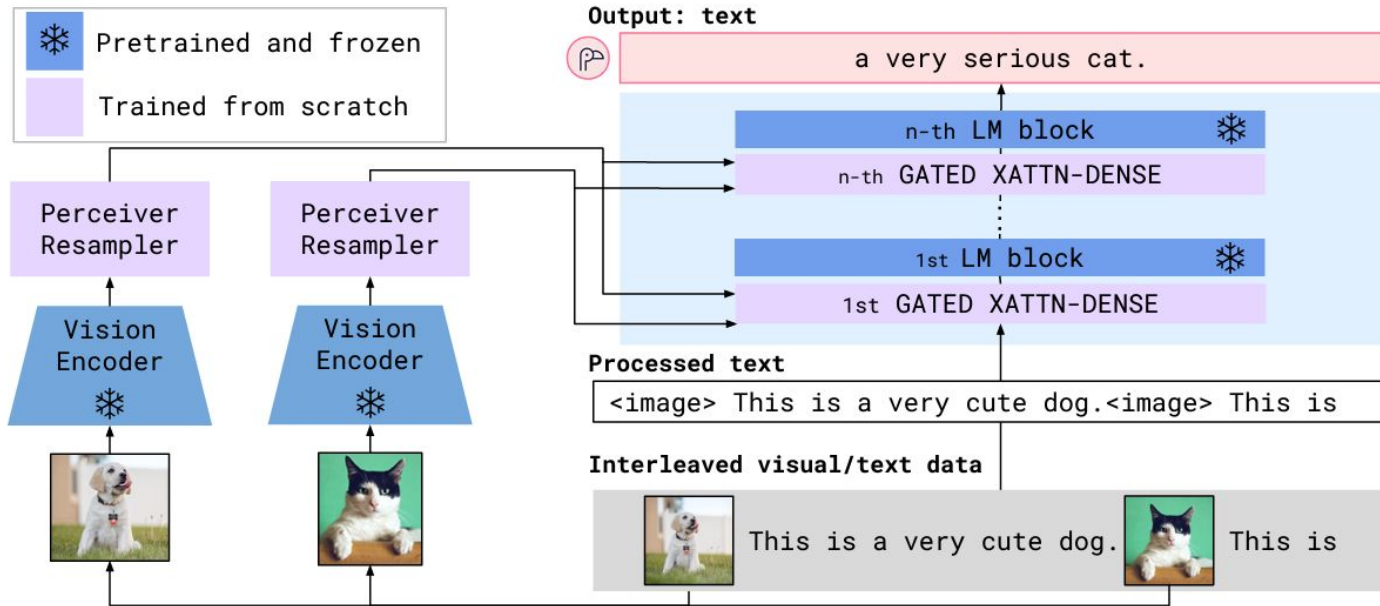
# Flamingo Overview



# Flamingo Overview

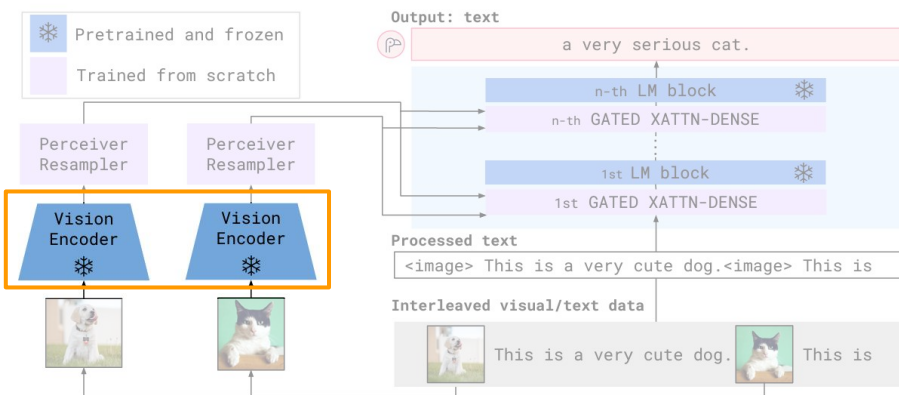


# Flamingo Overview

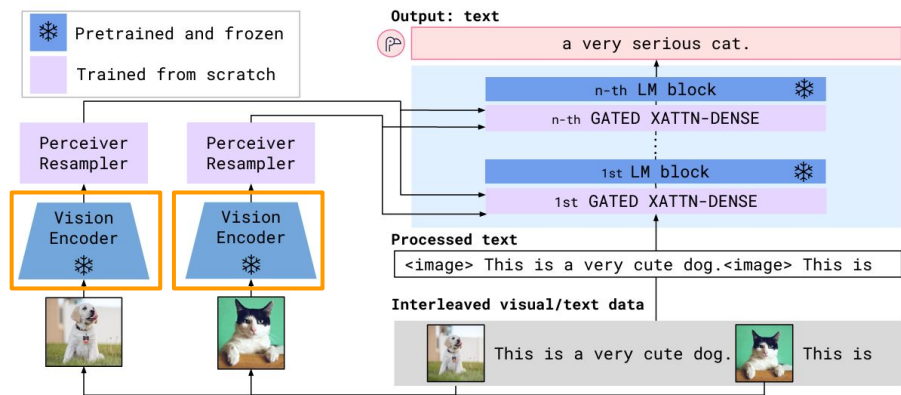


# Vision Encoder

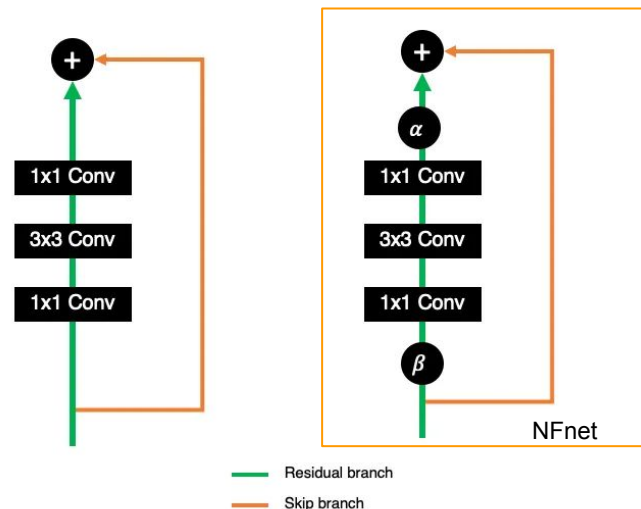
- **NFNet F6[8]** pretrained using the CLIP contrastive loss.



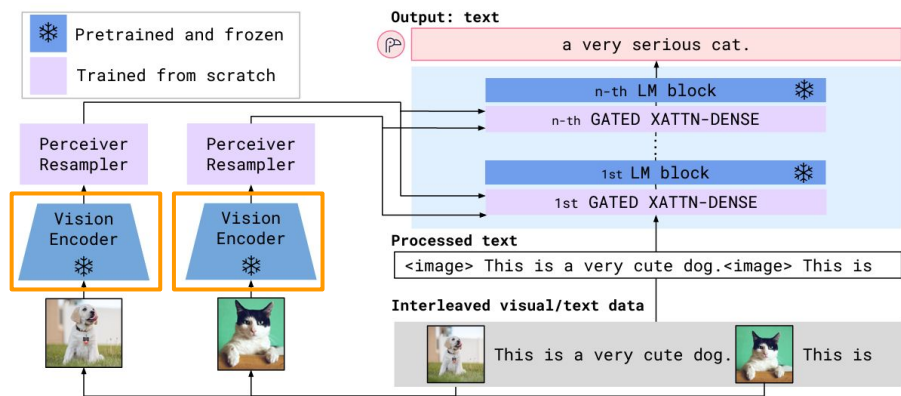
# Vision Encoder



- **NFNet F6[8]** pretrained using the CLIP contrastive loss.

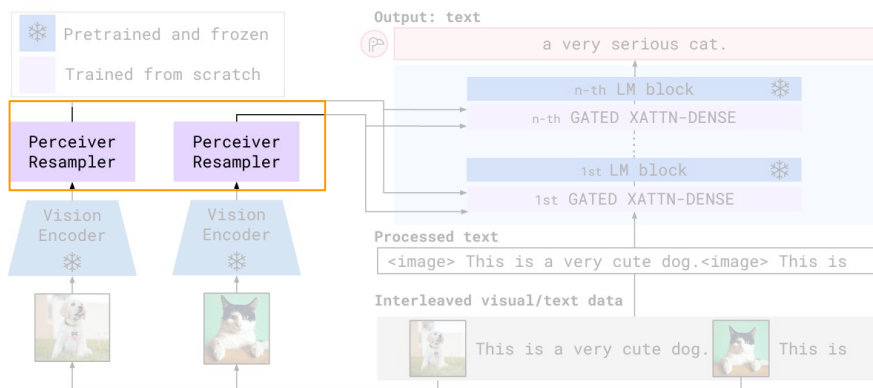


# Vision Encoder



- **NFNet F6[8]** pretrained using the CLIP contrastive loss.
- Trained on ALIGN and LTIP
- **Input:** 288 x 288 image
- **Output:** 2D grid Flattened to 1D
- 1FPS sampling for Videos
- Model is Frozen After Pretraining

# Perceiver Resampler





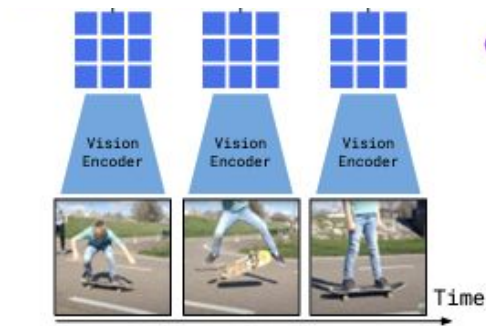
# Perceiver Resampler

- Consumes **variable** number of input frames



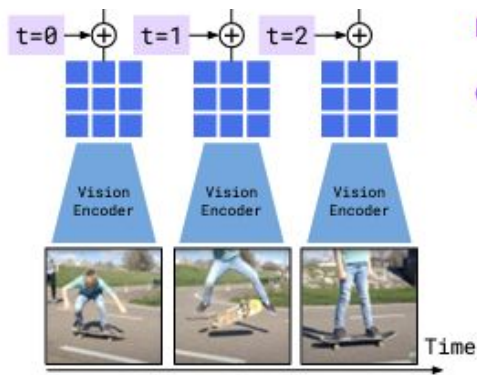
# Perceiver Resampler

- Consumes **variable** number of input frames



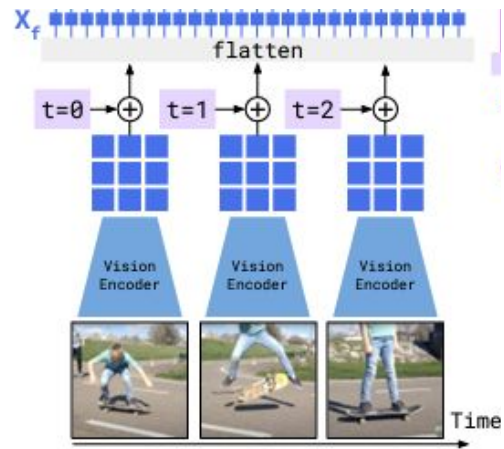
# Perceiver Resampler

- Consumes **variable** number of input frames.
- Appends **temporal** encodings.



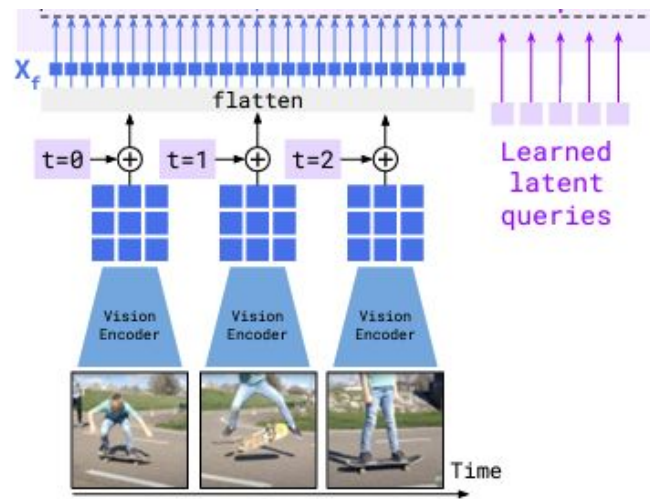
# Perceiver Resampler

- Consumes **variable** number of input frames.
- Appends **temporal** encodings.
- Flattens the image grid.



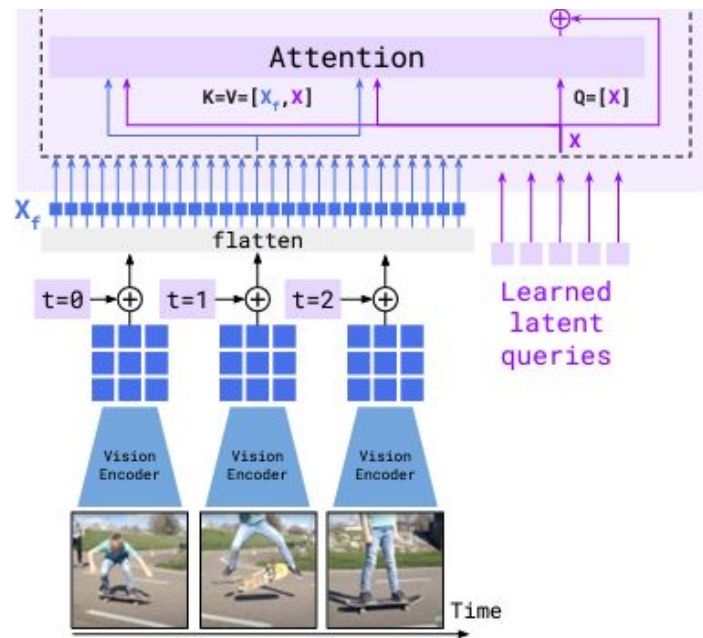
# Perceiver Resampler

- Consumes **variable** number of input frames.
- Appends **temporal** encodings.
- Flattens the image grid.
- Combined with **fixed number** of latent queries



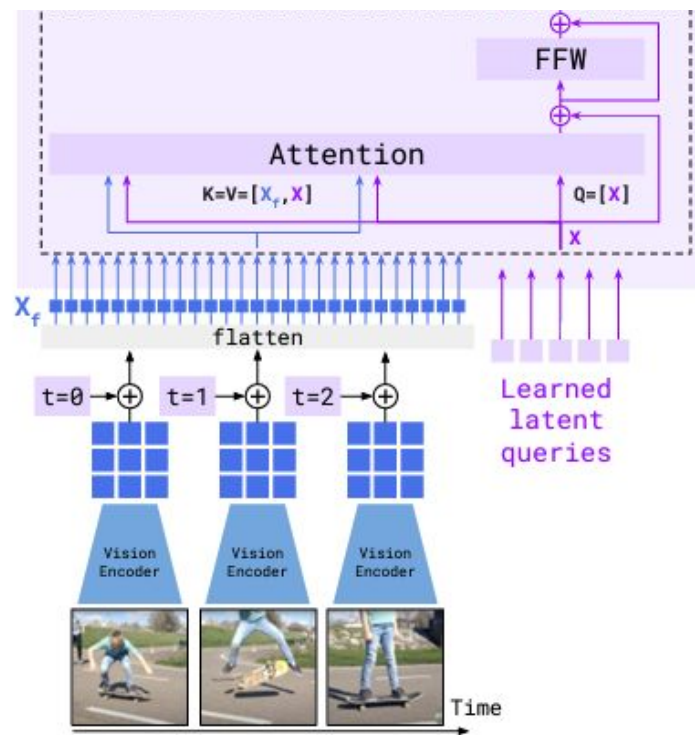
# Perceiver Resampler

- Consumes **variable** number of input Frames.
- Appends **temporal** encodings.
- Flattens the image grid.
- Combined with **fixed number** of latent queries.
- Attention layer with  $Q = \text{latent Queries}$ , and  $K, V = [\text{Image vector}, \text{latent Queries}]$



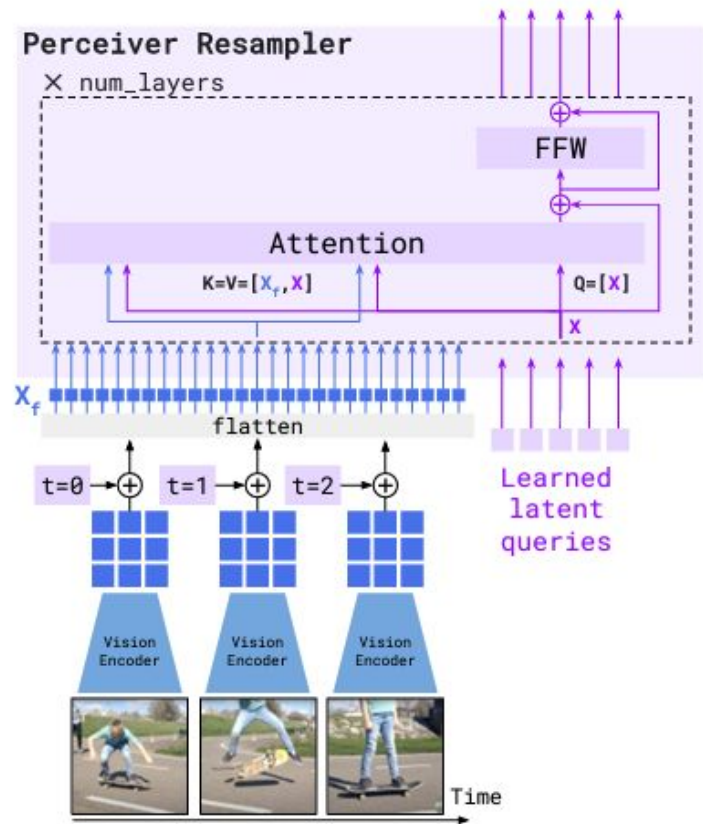
# Perceiver Resampler

- Consumes **variable** number of input Frames.
- Appends **temporal** encodings.
- Flattens the image grid.
- Combined with **fixed number** of latent queries.
- Attention layer with  $Q =$  latent Queries, and  $K, V = [\text{Image vector}, \text{latent Queries}]$



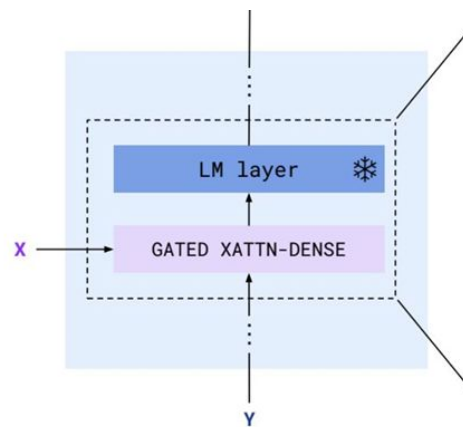
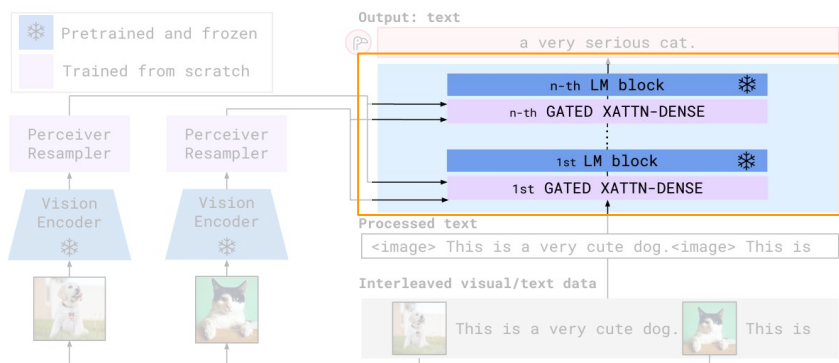
# Perceiver Resampler

- Consumes **variable** number of input Frames.
- Appends **temporal** encodings.
- Flattens the image grid.
- Combined with **fixed number** of latent queries.
- Attention layer with  $Q =$  latent Queries, and  $K, V = [X_r, X]$
- Outputs a Fixed number of visual tokens





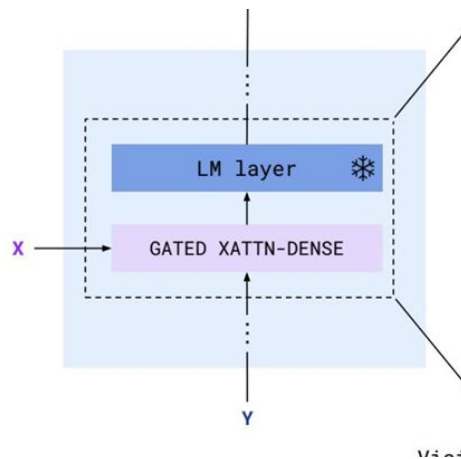
# Conditioning the Language model



Visual

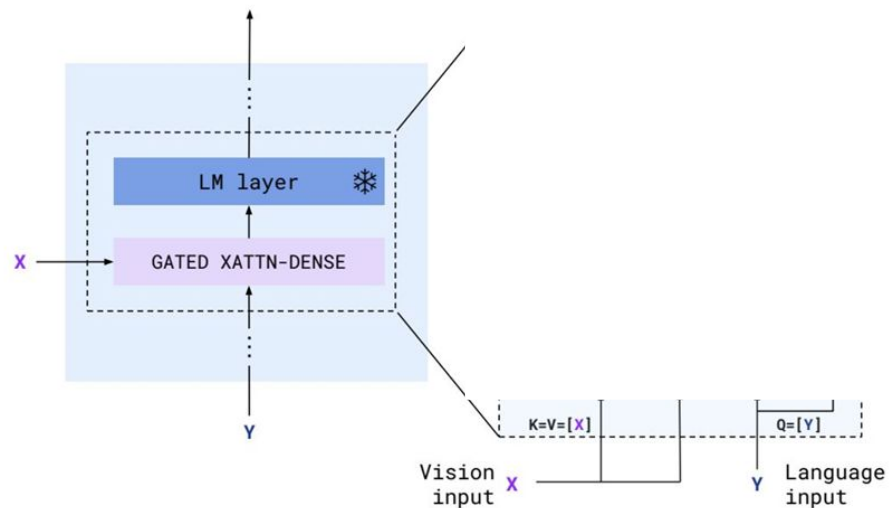
# Conditioning the Language model

- Flamingo uses **Chinchilla** class of LLMs.



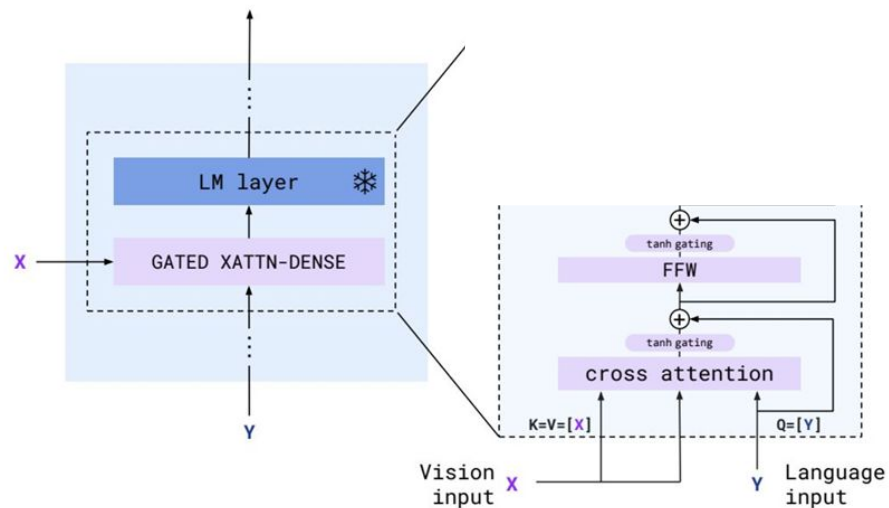
# Conditioning the Language model

- Flamingo uses **Chinchilla** class of LLMs.
- Vision (X) and language (Y) input to a **XATTN** block



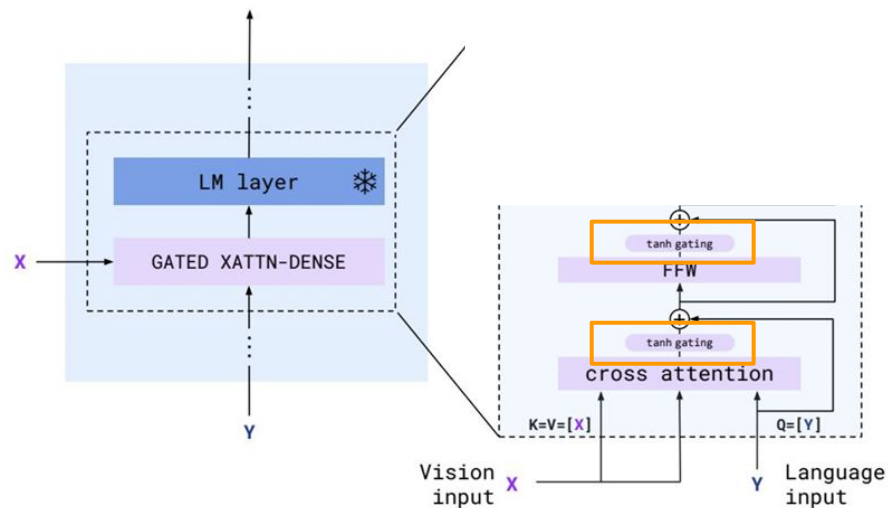
# Conditioning the Language model

- Flamingo uses **Chinchilla** class of LLMs.
- Vision (X) and language (Y) input to a **XATTN** block



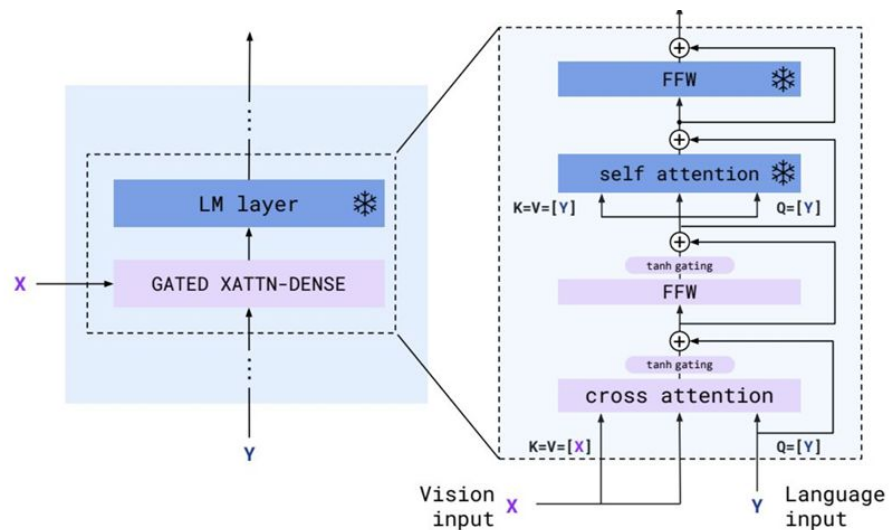
# Conditioning the Language model

- Flamingo uses **Chinchilla** class of LLMs.
- Vision (X) and language (Y) input to a **XATTN** block
- Uses TanH gating with layer learnable **alpha**
- **Alpha** initialized to 0 for stability



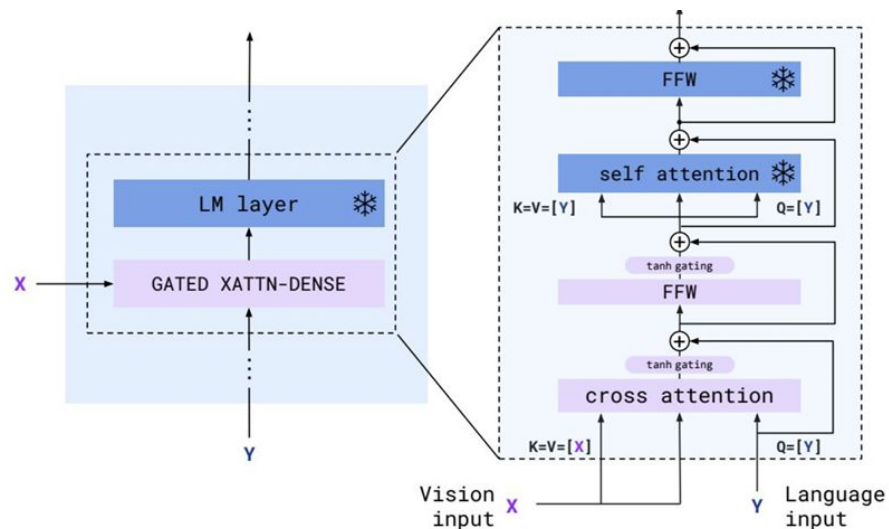
# Conditioning the Language model

- Flamingo uses **Chinchilla** class of LLMs.
- Vision (X) and language (Y) input to a **XATTN** block
- Uses TanH gating with layer learnable **alpha**
- **Alpha** initialized to 0 for stability.



# Conditioning the Language model

- Flamingo uses **Chinchilla** class of LLMs.
- Vision (X) and language (Y) input to a **XATTN** block
- Uses TanH gating with layer learnable **alpha**
- **Alpha** initialized to 0 for stability.
- Flamingo model variations are introduced through XATTN layers only.



# Multi-visual input support

- Tags are added to input text.





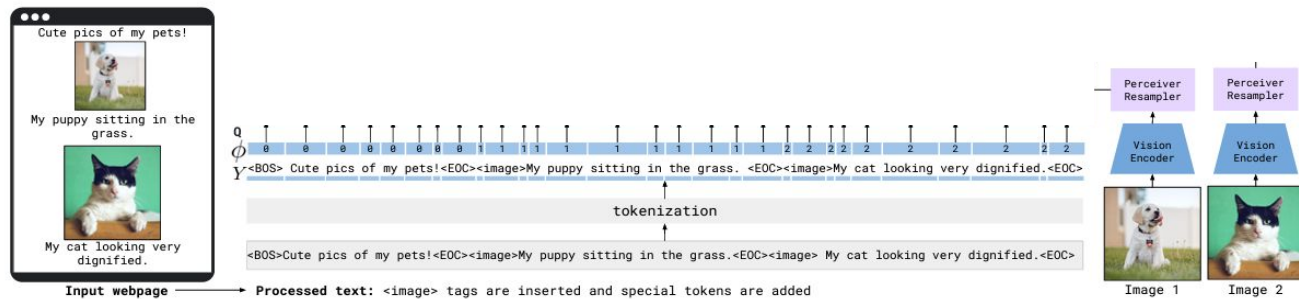
# Multi-visual input support

- Tags are added to input text.
- Images are processed.



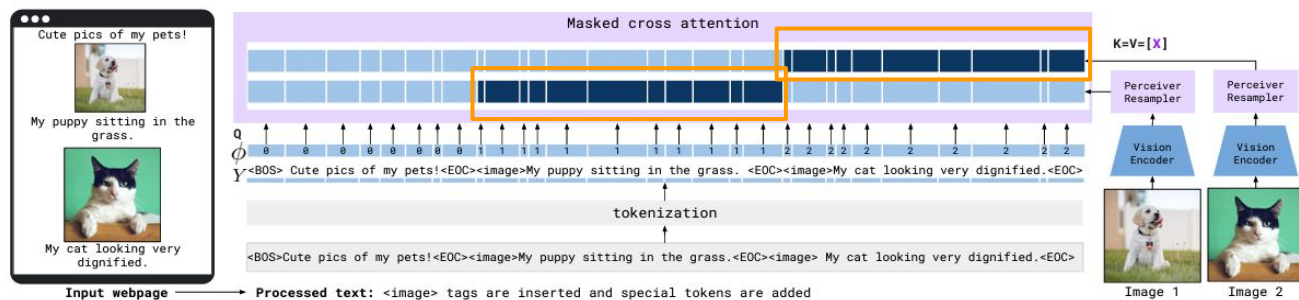
# Multi-visual input support

- Tags are added to input text.
- Images are processed.
- Function  $\phi$  that maps each text token to the last image token.



# Multi-visual input support

- Tags are added to input text.
- Images are processed.
- Function  $\phi$  that maps each text token to the last image token.
- Each token only attends to the last seen image token



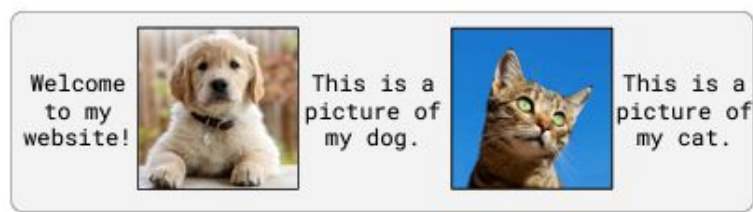
# Training Datasets



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]

- **ALIGN:** 1.8B pairs with 12.4 tokens on average
- **LTIP:** 312M pairs with 20.5 tokens on average

## VTP Dataset:

- 27M short Videos
- 22S duration on average

## M3W Dataset:

- 185M Images
- 182GB of Text

# Training

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

- Flamingo is trained by minimizing log likelihood of text given the previous input (text or image)
- The loss is weighted sum of all the datasets, where  $\mathcal{D}_m$  and  $\lambda_m$  are the  $m$ th dataset and its weight.

# Flamingo tasks

	Dataset	DEV	Gen.	Custom prompt	Task description	Eval set	Metric
Image	ImageNet-1k [94]	✓			Object classification	Val	Top-1 acc.
	MS-COCO [15]	✓	✓		Scene description	Test	CIDEr
	VQAv2 [3]	✓	✓		Scene understanding QA	Test-dev	VQA acc. [3]
	OKVQA [69]	✓	✓		External knowledge QA	Val	VQA acc. [3]
	Flickr30k [139]		✓		Scene description	Test (Karpathy)	CIDEr
	VizWiz [35]		✓		Scene understanding QA	Test-dev	VQA acc. [3]
	TextVQA [100]		✓		Text reading QA	Val	VQA acc. [3]
	VisDial [20]				Visual Dialogue	Val	NDCG
	HatefulMemes [54]			✓	Meme classification	Seen Test	ROC AUC
Video	Kinetics700 2020 [102]	✓			Action classification	Val	Top-1/5 avg
	VATEX [122]	✓	✓		Event description	Test	CIDEr
	MSVDQA [130]	✓	✓		Event understanding QA	Test	Top-1 acc.
	YouCook2 [149]		✓		Event description	Val	CIDEr
	MSRVTTQA [130]		✓		Event understanding QA	Test	Top-1 acc.
	iVQA [135]		✓		Event understanding QA	Test	iVQA acc. [135]
	RareAct [73]			✓	Composite action retrieval	Test	mWAP
	NextQA [129]		✓		Temporal/Causal QA	Test	WUPS
	STAR [128]				Multiple-choice QA	Test	Top-1 acc.

# Zeroshot and few shot results

Method	FT	Shot	OKVQA (I)	VQAv2 (I)	COCO (I)	MSVDQA (V)	VATEX (V)	VizWiz (I)	Flick30K (I)	MSRVTTQA (V)	iVQA (V)	YouCook2 (V)	STAR (V)	VisDial (I)	TextVQA (I)	NextQA (I)	HatefulMemes (I)	RareAct (V)
Zero/Few shot SOTA	<b>X</b>	(X)	[34] 43.3 (16)	[114] 38.2 (4)	[124] 32.2 (0)	[58] 35.2 (0)	-	-	-	[58] 19.2 (0)	[135] 12.2 (0)	-	[143] 39.4 (0)	[79] 11.6 (0)	-	-	[85] 66.1 (0)	[85] 40.7 (0)
<i>Flamingo</i> -3B	<b>X</b>	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
	<b>X</b>	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	-
	<b>X</b>	32	45.9	57.1	99.0	42.6	59.2	45.5	71.2	25.6	37.7	76.7	41.6	47.3	30.6	26.1	56.3	-
<i>Flamingo</i> -9B	<b>X</b>	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
	<b>X</b>	4	49.3	56.3	93.1	36.2	51.7	34.9	72.6	18.2	37.7	70.8	<b>42.8</b>	50.4	33.6	24.7	62.7	-
	<b>X</b>	32	51.0	60.4	106.3	47.2	57.4	44.0	72.8	29.4	40.7	77.3	<u>41.2</u>	50.4	32.6	28.4	63.5	-
<i>Flamingo</i>	<b>X</b>	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	<b>60.8</b>
	<b>X</b>	4	57.4	63.1	103.2	41.7	56.0	39.6	75.1	23.9	44.1	74.5	42.4	<b>55.6</b>	36.5	30.8	68.6	-
	<b>X</b>	32	<b>57.8</b>	<b>67.6</b>	<b>113.8</b>	<b>52.3</b>	<b>65.1</b>	<b>49.8</b>	<b>75.4</b>	<b>31.0</b>	<b>45.3</b>	<b>86.8</b>	42.2	<b>55.6</b>	<b>37.9</b>	<b>33.5</b>	<b>70.0</b>	-
Pretrained FT SOTA	✓	(X)	54.4 [34] (10K)	80.2 [140] (444K)	143.3 [124] (500K)	47.9 [28] (27K)	76.3 [153] (500K)	57.2 [65] (20K)	67.4 [150] (30K)	46.8 [51] (130K)	35.4 [135] (6K)	138.7 [132] (10K)	36.7 [128] (46K)	75.2 [79] (123K)	54.7 [137] (20K)	25.2 [129] (38K)	79.1 [62] (9K)	-

# Zeroshot and few shot results

Method	FT	Shot	OKVQA (I)	VQA <sub>v2</sub> (I)	COCO (I)	MSVDQA (V)	VATEX (V)	VizWiz (I)	Flick30K (I)	MSRVTTQA (V)	iVQA (V)	YouCook2 (V)	STAR (V)	VisDial (I)	TextVQA (I)	NextQA (I)	HatefulMemes (I)	RareAct (V)
Zero/Few shot SOTA	X	(X)	[34] 43.3 (16)	[114] 38.2 (4)	[124] 32.2 (0)	[58] 35.2 (0)	-	-	-	[58] 19.2 (0)	[135] 12.2 (0)	-	[143] 39.4 (0)	[79] 11.6 (0)	-	-	[85] 66.1 (0)	[85] 40.7 (0)
Flamingo-3B	X	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
	X	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	-
	X	32	45.9	57.1	99.0	42.6	59.2	45.5	71.2	25.6	37.7	76.7	41.6	47.3	30.6	26.1	56.3	-
Flamingo-9B	X	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
	X	4	49.3	56.3	93.1	36.2	51.7	34.9	72.6	18.2	37.7	70.8	<b>42.8</b>	50.4	33.6	24.7	62.7	-
	X	32	51.0	60.4	106.3	47.2	57.4	44.0	72.8	29.4	40.7	77.3	41.2	50.4	32.6	28.4	63.5	-
Flamingo	X	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	<b>60.8</b>
	X	4	57.4	63.1	103.2	41.7	56.0	39.6	75.1	23.9	44.1	74.5	42.4	<b>55.6</b>	36.5	30.8	68.6	-
	X	32	<b>57.8</b>	<b>67.6</b>	<b>113.8</b>	<b>52.3</b>	<b>65.1</b>	<b>49.8</b>	<b>75.4</b>	<b>31.0</b>	<b>45.3</b>	<b>86.8</b>	42.2	<b>55.6</b>	<b>37.9</b>	<b>33.5</b>	<b>70.0</b>	-
Pretrained FT SOTA	✓	(X)	54.4 [34] (10K)	80.2 [140] (444K)	143.3 [124] (500K)	47.9 [28] (27K)	76.3 [153] (500K)	57.2 [65] (20K)	67.4 [150] (30K)	46.8 [51] (130K)	35.4 [135] (6K)	138.7 [132] (10K)	36.7 [128] (46K)	75.2 [79] (123K)	54.7 [137] (20K)	25.2 [129] (38K)	79.1 [62] (9K)	-

Flamingo is better than current SOTA few shot/zero shot

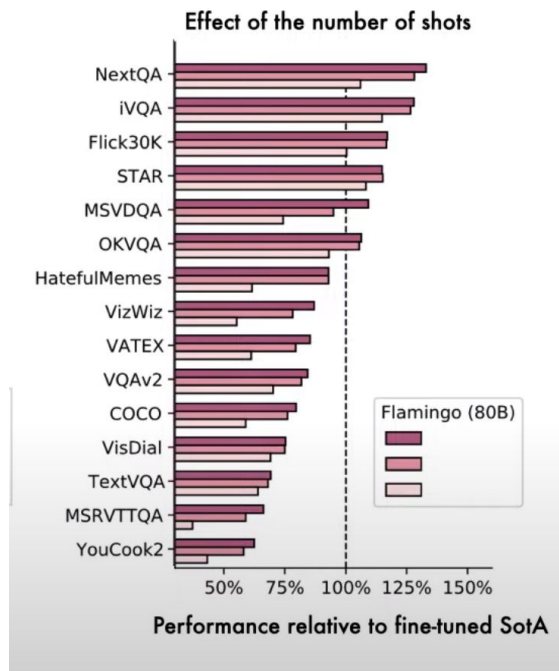


# Zeroshot and few shot results

Method	FT	Shot	OKVQA (I)	VQAv2 (I)	COCO (I)	MSVDQA (V)	VATEX (V)	VizWiz (I)	Flick30K (I)	MSRVTTQA (V)	iVQA (V)	YouCook2 (V)	STAR (V)	VisDial (I)	TextVQA (I)	NextQA (I)	HatefulMemes (I)	RareAct (V)
Zero/Few shot SOTA	X	(X)	[34] 43.3 (16)	[114] 38.2 (4)	[124] 32.2 (0)	[58] 35.2 (0)	-	-	-	[58] 19.2 (0)	[135] 12.2 (0)	-	[143] 39.4 (0)	[79] 11.6 (0)	-	-	[85] 66.1 (0)	[85] 40.7 (0)
Flamingo-3B	X	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
	X	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	-
	X	32	45.9	57.1	99.0	42.6	59.2	45.5	71.2	25.6	37.7	76.7	41.6	47.3	30.6	26.1	56.3	-
Flamingo-9B	X	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
	X	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	-
	X	32	51.0	60.4	106.3	47.2	57.4	44.0	72.8	29.4	40.7	77.3	41.2	50.4	32.6	28.4	63.5	-
Flamingo	X	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
	X	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	-
	X	32	57.8	67.6	113.8	52.3	65.1	49.8	75.4	31.0	45.3	86.8	42.2	55.6	37.9	33.5	70.0	-
Pretrained FT SOTA	✓	(X)	54.4 [34] (10K)	80.2 [140] (444K)	143.3 [124] (500K)	47.9 [28] (27K)	76.3 [153] (500K)	57.2 [65] (20K)	67.4 [150] (30K)	46.8 [51] (130K)	35.4 [135] (6K)	138.7 [132] (10K)	36.7 [128] (46K)	75.2 [79] (123K)	54.7 [137] (20K)	25.2 [129] (38K)	79.1 [162] (9K)	-

It achieves SOTA on 6 tasks

# Zeroshot and few shot results



Performance increases generally when the number of shots are increased.

# Finetuning results

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 [124]	80.0	80.3	<b>143.3</b>	-	-	-	-	-	-	-	-	-	-
OFA [119]	79.9	80.0	<u>149.6</u>	-	-	-	-	-	-	-	-	-	-
Florence [140]	80.2	80.4	-	-	-	-	-	-	-	-	-	-	-
🦩 <i>Flamingo</i> Fine-tuned	<b>82.0</b>	<b>82.1</b>	138.1	<b>84.2</b>	<b>65.7</b>	<b>65.4</b>	<b>47.4</b>	61.8	59.7	118.6	<b>57.1</b>	54.1	<b>86.6</b>
Restricted SotA <sup>†</sup>	80.2	80.4	<b>143.3</b>	76.3	-	-	46.8	<u>75.2</u>	<b>74.5</b>	<b>138.7</b>	54.7	<u>73.7</u>	79.1
	[140]	[140]	[124]	[153]	-	-	[51]	[79]	[79]	[132]	[137]	[84]	[62]
Unrestricted SotA	81.3	81.3	<u>149.6</u>	81.4	57.2	60.6	-	-	<u>75.4</u>	-	-	-	84.6
	[133]	[133]	[119]	[153]	[65]	[65]	-	-	[123]	-	-	-	[152]

# Model Scaling

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

# Ablation studies

Ablated setting	<i>Flamingo</i> -3B original value	Changed value	Param. count ↓	Step time ↓	COCO CIDEr ↑	OKVQA top1 ↑	VQAv2 top1 ↑	MSVDQA top1 ↑	VATEX CIDEr ↑	Overall score ↑	
<b><i>Flamingo</i>-3B model</b>			3.2B	1.74s	86.5	42.1	55.8	36.3	53.4	<b>70.7</b>	
(i)	Training data	All data	w/o Video-Text pairs	3.2B	1.42s	84.2	43.0	53.9	34.5	46.0	67.3
			w/o Image-Text pairs	3.2B	0.95s	66.3	39.2	51.6	32.0	41.6	60.9
			Image-Text pairs → LAION	3.2B	1.74s	79.5	41.4	53.5	33.9	47.6	66.4
			w/o M3W	3.2B	1.02s	54.1	36.5	52.7	31.4	23.5	53.4
(ii)	Optimisation	Accumulation	Round Robin	3.2B	1.68s	76.1	39.8	52.1	33.2	40.8	62.9
(iii)	Tanh gating	✓	✗	3.2B	1.74s	78.4	40.5	52.9	35.9	47.5	66.5
(iv)	Cross-attention architecture	GATED XATTN-DENSE	VANILLA XATTN	2.4B	1.16s	80.6	41.5	53.4	32.9	50.7	66.9
			GRAFTING	3.3B	1.74s	79.2	36.1	50.8	32.2	47.8	63.1
(v)	Cross-attention frequency	Every	Single in middle	2.0B	0.87s	71.5	38.1	50.2	29.1	42.3	59.8
			Every 4th	2.3B	1.02s	82.3	42.7	55.1	34.6	50.8	68.8
			Every 2nd	2.6B	1.24s	83.7	41.0	55.8	34.5	49.7	68.2
(vi)	Resampler	Perceiver	MLP	3.2B	1.85s	78.6	42.2	54.7	35.2	44.7	66.6
			Transformer	3.2B	1.81s	83.2	41.7	55.6	31.5	48.3	66.7
(vii)	Vision encoder	NFNet-F6	CLIP ViT-L/14	3.1B	1.58s	76.5	41.6	53.4	33.2	44.5	64.9
			NFNet-F0	2.9B	1.45s	73.8	40.5	52.8	31.1	42.9	62.7
(viii)	Freezing LM	✓	✗ (random init)	3.2B	2.42s	74.8	31.5	45.6	26.9	50.1	57.8
			✗ (pretrained)	3.2B	2.42s	81.2	33.7	47.4	31.0	53.9	62.7

# Ablation studies

Ablated setting	<i>Flamingo</i> -3B original value	Changed value	Param. count ↓	Step time ↓	COCO CIDEr ↑	OKVQA top1 ↑	VQAv2 top1 ↑	MSVDQA top1 ↑	VATEX CIDEr ↑	Overall score ↑	
<b><i>Flamingo</i>-3B model</b>			3.2B	1.74s	86.5	42.1	55.8	36.3	53.4	<b>70.7</b>	
(i)	Training data	All data	w/o Video-Text pairs	3.2B	1.42s	84.2	43.0	53.9	34.5	46.0	67.3
			w/o Image-Text pairs	3.2B	0.95s	66.3	39.2	51.6	32.0	41.6	60.9
			Image-Text pairs → LAION	3.2B	1.74s	79.5	41.4	53.5	33.9	47.6	66.4
			w/o M3W	3.2B	1.02s	54.1	36.5	52.7	31.4	23.5	53.4
(ii)	Optimisation	Accumulation	Round Robin	3.2B	1.68s	76.1	39.8	52.1	33.2	40.8	62.9
(iii)	Tanh gating	✓	✗	3.2B	1.74s	78.4	40.5	52.9	35.9	47.5	66.5
(iv)	Cross-attention architecture	GATED XATTN-DENSE	VANILLA XATTN	2.4B	1.16s	80.6	41.5	53.4	32.9	50.7	66.9
			GRAFTING	3.3B	1.74s	79.2	36.1	50.8	32.2	47.8	63.1
(v)	Cross-attention frequency	Every	Single in middle	2.0B	0.87s	71.5	38.1	50.2	29.1	42.3	59.8
			Every 4th	2.3B	1.02s	82.3	42.7	55.1	34.6	50.8	68.8
			Every 2nd	2.6B	1.24s	83.7	41.0	55.8	34.5	49.7	68.2
(vi)	Resampler	Perceiver	MLP	3.2B	1.85s	78.6	42.2	54.7	35.2	44.7	66.6
			Transformer	3.2B	1.81s	83.2	41.7	55.6	31.5	48.3	66.7
(vii)	Vision encoder	NFNet-F6	CLIP ViT-L/14	3.1B	1.58s	76.5	41.6	53.4	33.2	44.5	64.9
			NFNet-F0	2.9B	1.45s	73.8	40.5	52.8	31.1	42.9	62.7
(viii)	Freezing LM	✓	✗ (random init)	3.2B	2.42s	74.8	31.5	45.6	26.9	50.1	57.8
			✗ (pretrained)	3.2B	2.42s	81.2	33.7	47.4	31.0	53.9	62.7

# Ablation studies

Ablated setting	<i>Flamingo</i> -3B original value	Changed value	Param. count ↓	Step time ↓	COCO CIDEr↑	OKVQA top1↑	VQAv2 top1↑	MSVDQA top1↑	VATEX CIDEr↑	Overall score↑	
<b><i>Flamingo</i>-3B model</b>			3.2B	1.74s	86.5	42.1	55.8	36.3	53.4	<b>70.7</b>	
(i)	Training data	All data	w/o Video-Text pairs	3.2B	1.42s	84.2	43.0	53.9	34.5	46.0	67.3
			w/o Image-Text pairs	3.2B	0.95s	66.3	39.2	51.6	32.0	41.6	60.9
			Image-Text pairs → LAION	3.2B	1.74s	79.5	41.4	53.5	33.9	47.6	66.4
			w/o M3W	3.2B	1.02s	54.1	36.5	52.7	31.4	23.5	53.4
(ii)	Optimisation	Accumulation	Round Robin	3.2B	1.68s	76.1	39.8	52.1	33.2	40.8	62.9
(iii)	Tanh gating	✓	✗	3.2B	1.74s	78.4	40.5	52.9	35.9	47.5	66.5
(iv)	Cross-attention architecture	GATED XATTN-DENSE	VANILLA XATTN	2.4B	1.16s	80.6	41.5	53.4	32.9	50.7	66.9
			GRAFTING	3.3B	1.74s	79.2	36.1	50.8	32.2	47.8	63.1
(v)	Cross-attention frequency	Every	Single in middle	2.0B	0.87s	71.5	38.1	50.2	29.1	42.3	59.8
			Every 4th	2.3B	1.02s	82.3	42.7	55.1	34.6	50.8	68.8
			Every 2nd	2.6B	1.24s	83.7	41.0	55.8	34.5	49.7	68.2
(vi)	Resampler	Perceiver	MLP	3.2B	1.85s	78.6	42.2	54.7	35.2	44.7	66.6
			Transformer	3.2B	1.81s	83.2	41.7	55.6	31.5	48.3	66.7
(vii)	Vision encoder	NFNet-F6	CLIP ViT-L/14	3.1B	1.58s	76.5	41.6	53.4	33.2	44.5	64.9
			NFNet-F0	2.9B	1.45s	73.8	40.5	52.8	31.1	42.9	62.7
(viii)	Freezing LM	✓	✗ (random init)	3.2B	2.42s	74.8	31.5	45.6	26.9	50.1	57.8
			✗ (pretrained)	3.2B	2.42s	81.2	33.7	47.4	31.0	53.9	62.7

# Ablation studies

Ablated setting	<i>Flamingo</i> -3B original value	Changed value	Param. count ↓	Step time ↓	COCO CIDEr↑	OKVQA top1↑	VQAv2 top1↑	MSVDQA top1↑	VATEX CIDEr↑	Overall score↑	
<b><i>Flamingo</i>-3B model</b>			3.2B	1.74s	86.5	42.1	55.8	36.3	53.4	<b>70.7</b>	
(i)	Training data	All data	w/o Video-Text pairs	3.2B	1.42s	84.2	43.0	53.9	34.5	46.0	67.3
			w/o Image-Text pairs	3.2B	0.95s	66.3	39.2	51.6	32.0	41.6	60.9
			Image-Text pairs → LAION	3.2B	1.74s	79.5	41.4	53.5	33.9	47.6	66.4
			w/o M3W	3.2B	1.02s	54.1	36.5	52.7	31.4	23.5	53.4
(ii)	Optimisation	Accumulation	Round Robin	3.2B	1.68s	76.1	39.8	52.1	33.2	40.8	62.9
(iii)	Tanh gating	✓	✗	3.2B	1.74s	78.4	40.5	52.9	35.9	47.5	66.5
(iv)	Cross-attention architecture	GATED XATTN-DENSE	VANILLA XATTN	2.4B	1.16s	80.6	41.5	53.4	32.9	50.7	66.9
			GRAFTING	3.3B	1.74s	79.2	36.1	50.8	32.2	47.8	63.1
(v)	Cross-attention frequency	Every	Single in middle	2.0B	0.87s	71.5	38.1	50.2	29.1	42.3	59.8
			Every 4th	2.3B	1.02s	82.3	42.7	55.1	34.6	50.8	68.8
			Every 2nd	2.6B	1.24s	83.7	41.0	55.8	34.5	49.7	68.2
(vi)	Resampler	Perceiver	MLP	3.2B	1.85s	78.6	42.2	54.7	35.2	44.7	66.6
			Transformer	3.2B	1.81s	83.2	41.7	55.6	31.5	48.3	66.7
(vii)	Vision encoder	NFNet-F6	CLIP ViT-L/14	3.1B	1.58s	76.5	41.6	53.4	33.2	44.5	64.9
			NFNet-F0	2.9B	1.45s	73.8	40.5	52.8	31.1	42.9	62.7
(viii)	Freezing LM	✓	✗ (random init)	3.2B	2.42s	74.8	31.5	45.6	26.9	50.1	57.8
			✗ (pretrained)	3.2B	2.42s	81.2	33.7	47.4	31.0	53.9	62.7



# Classification 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	-	0	<b>85.7 [82]</b>	<b>69.6 [85]</b>
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	$\leq 0.02$	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

# Strengths

- The addition of extra layers while keeping the rest of the model frozen preserves knowledge of both models and is novel.
- The method is able to get very impressive results just using few input samples as demonstrations.
- The paper and appendix include a huge number of studies, justifying most of their model decisions, data decisions, parameter choices etc.

# Weaknesses

- The spotlight flamingo model that gets the best results is exceptionally big at 80B parameters and making it quite cumbersome to use.
- The authors havent released their model and data, this is not a technical weakness but sets a bad precedent within the research community.
- Flamingo performs worse than its vision encoder on image classification.

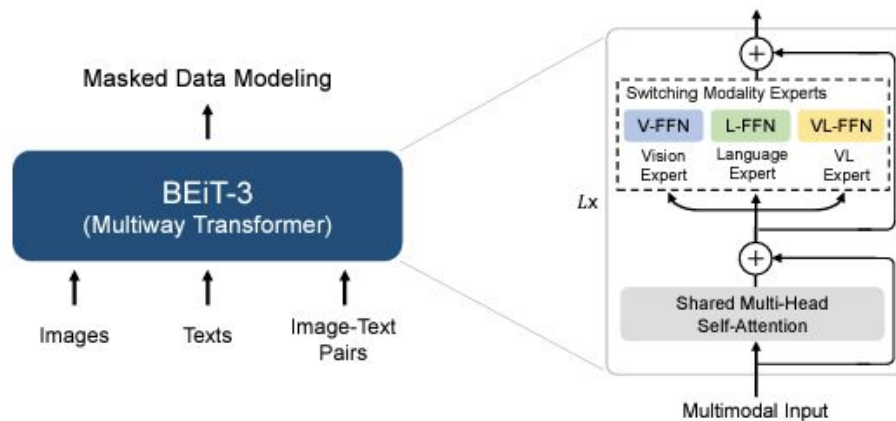
# Language Is Not All You Need: Aligning Perception with Language Models

# Kosmos-1: Key Takeaways

- A large multimodal LLM that can perceive general modalities, perform zero shot and few shot learning.
- Trained on a web scale multimodal corpora containing interleaved text and images.
- Kosmos-1 demonstrates impressive capabilities across, vision, language and perception language tasks.
- They evaluate on unique tasks like multimodal chain of thought reasoning, OCR free NLP and a novel nonverbal reasoning test.

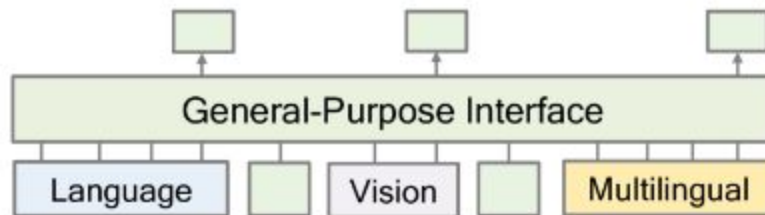
# Related Works

- **MetaLM[2]: LLMs are general purpose interfaces**
  - Any input and output format that can be converted to text token can be a LLM usecase.
- **Extending LLMs to multimodal tasks**
  - Flamingo[3]: Large MLLM for few shot learning
  - BeIT[9]: masked language modelling on images, text, and pairs in a unified manner



# Kosmos - 1: Overview

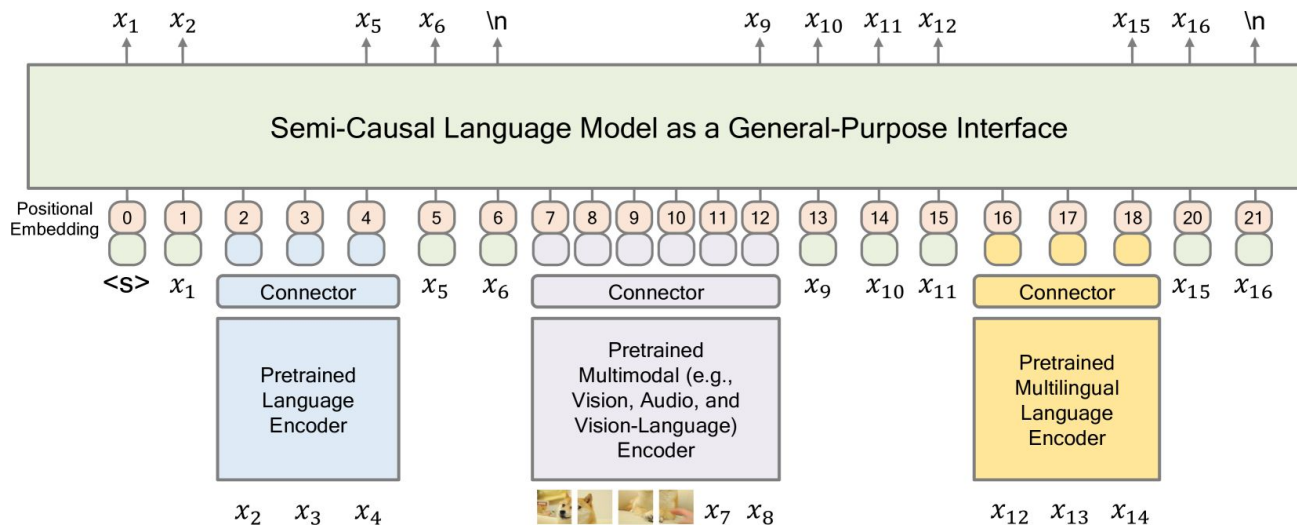
- Kosmos - 1 follows the same philosophy as the MetaLM and treats language models as a universal task layer.
- It builds on MetaLM, trains on more multimodal data, uses interleaved inputs.



MetaLM

# Meta LM: Language models are a general purpose Interface

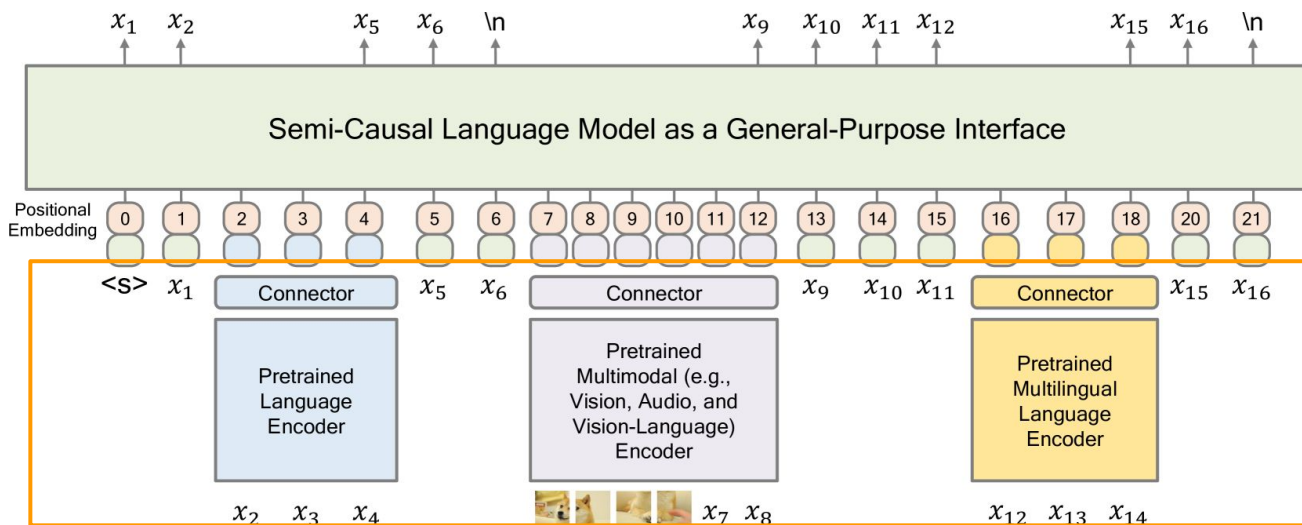
- MetaLM proposes to use LMs as a general interface for all kinds of input like video, images, multilingual etc.



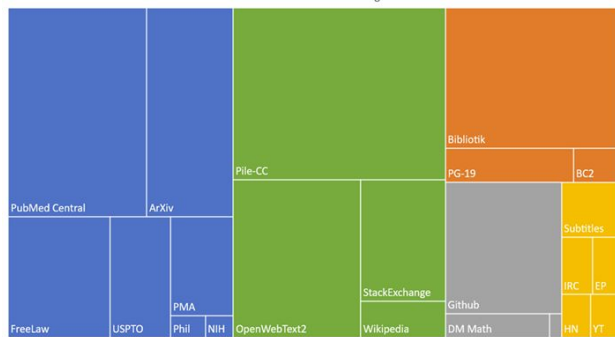


# Meta LM: Language models are a general purpose Interface

- MetaLM proposes to use LMs as a general interface for all kinds of input like video, images, multilingual etc.

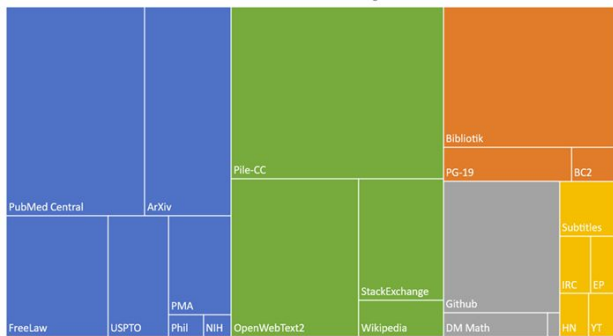


# Meta LM: Pretraining and Evaluation



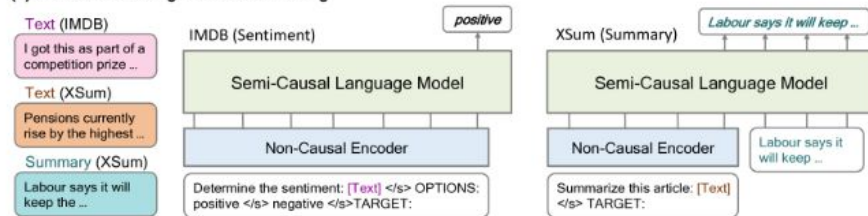
Train on PILE for language only Tasks

# Meta LM: Pretraining and Evaluation

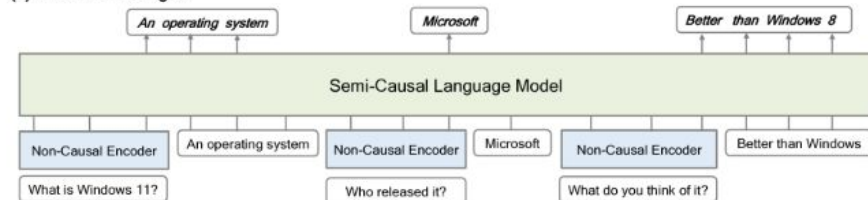


Train on PILE for language only Tasks

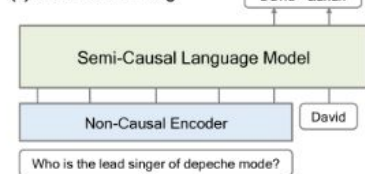
(a) Multitask Learning / Instruction Tuning



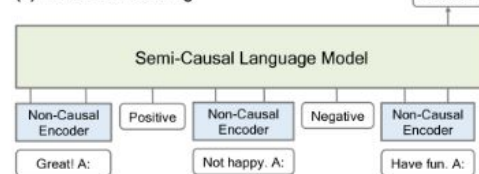
(b) Multi-Turn Dialogue



(c) Zero-Shot Priming

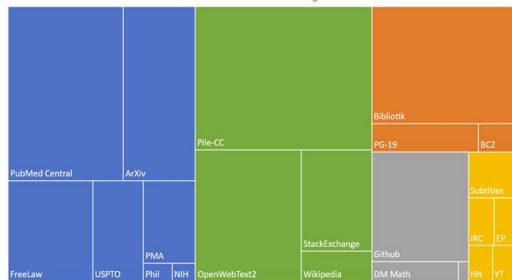


(d) In-Context Learning



Evaluate on a Number of Tasks

# Meta LM: Pretraining and Evaluation



Visual Genome

COCO caption

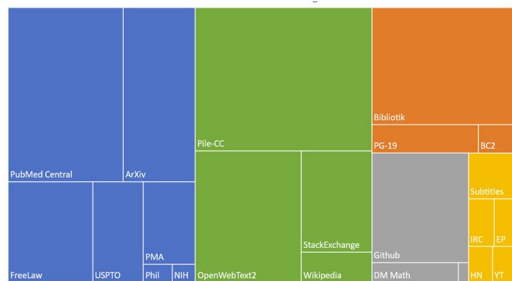
Conceptual  
Captions

SBU Caption

4M images + 10M Pairs

Train on PILE with Image-caption datasets

# Meta LM: Pretraining and Evaluation



Visual Genome

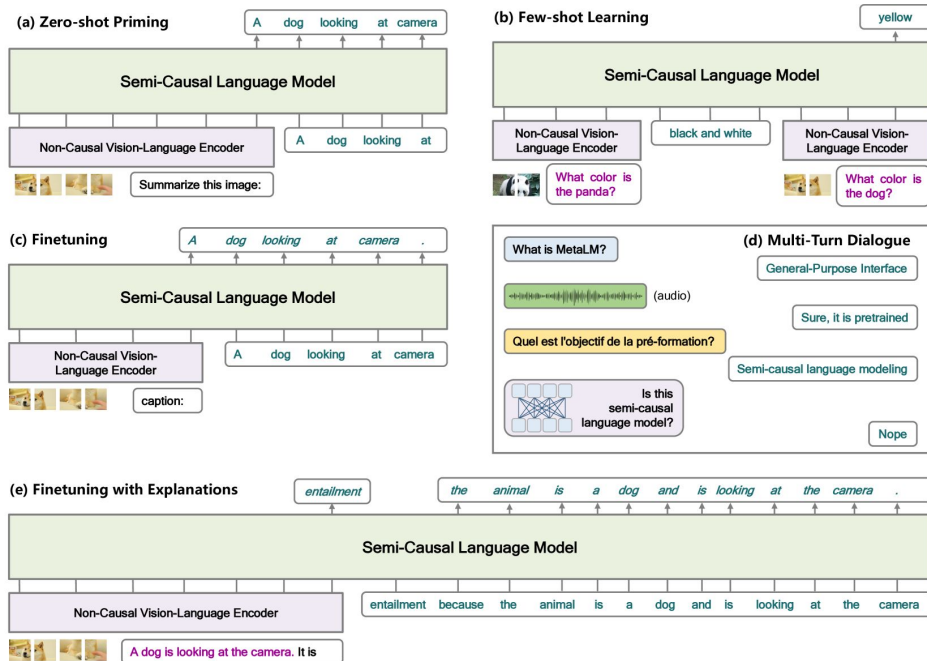
COCO caption

Conceptual Captions

SBU Caption

4M images + 10M Pairs

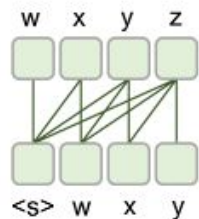
Train on PILE with Image-caption datasets



Evaluate on a number of tasks

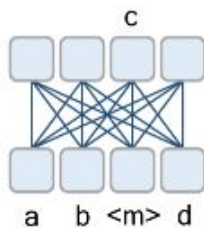
# Meta LM: Model Architecture

- Introduces a new semi-causal architecture that jointly learns with a combination of pretrained encoders each focusing on a modality.



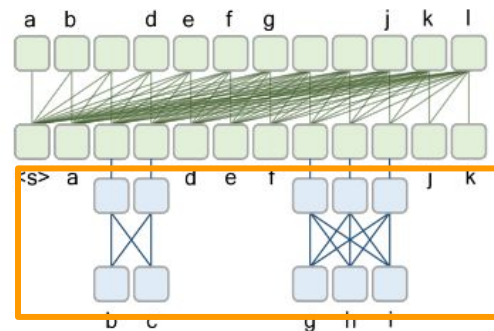
(a) Causal LM  
(Unidirectional)

GPT 3 style decoder



(c) Non-Causal LM  
(Bidirectional)

BERT style encoder

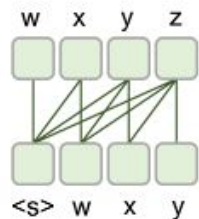


(d) Semi-Causal LM

MetaLM

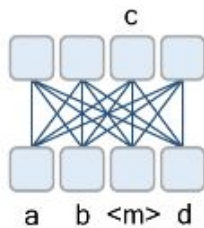
# Meta LM: Model Architecture

- Introduces a new semi-causal architecture that jointly learns with a combination of pretrained encoders each focusing on a modality.



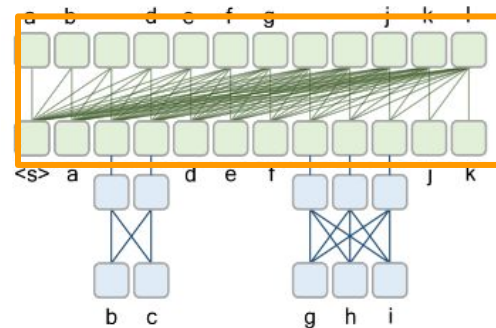
(a) Causal LM  
(Unidirectional)

GPT 3 style decoder



(c) Non-Causal LM  
(Bidirectional)

BERT style encoder



(d) Semi-Causal LM

MetaLM

# Meta LM: Loss Function

- Trains on a semi causal modelling objective, where the model predicts the text token given the representations of previous tokens from bidirectional encoders.

$$\max \sum_{i=0}^k \sum_{t=e_t}^{s_{(i+1)}} \log P(x_t \mid \mathbf{x}_{<t}, \{\mathbf{h}(\mathbf{x}_{s_j}^e)\}_{j < i})$$

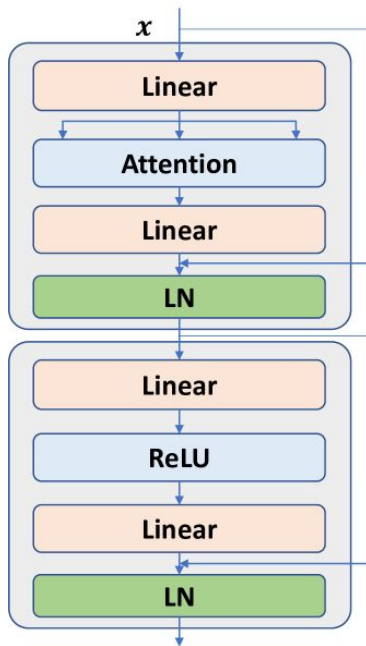
Here  $H(X)$  is the encoder function for each span which can be text or image.



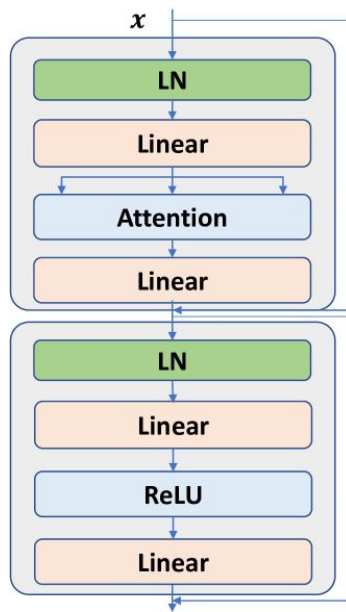
# Going from MetaLM to Kosmos-1

- Kosmos-1 trains a model with 24 layers with a hidden dimension size of 2048 and 32 attention heads, totalling to 1.3B parameters similar to MetaLM.
- They change the default transformer module to the **magenta[10]** module.

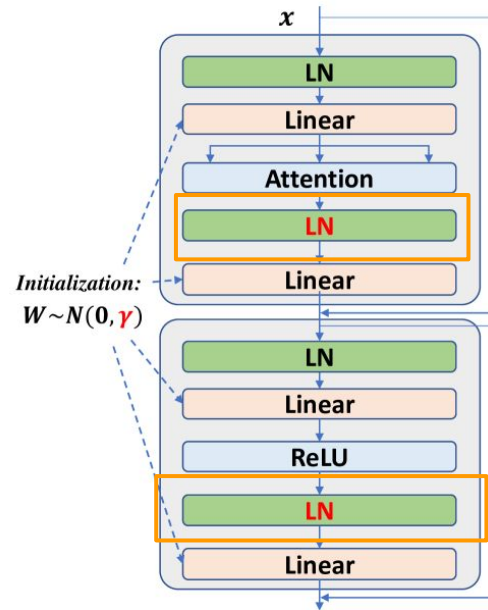
# Background: Magento Module vs Other methods



BERT, NMT,  
Transformer



ViT, GPT



Magento

# Going from MetaLM to Kosmos-1

- Kosmos-1 trains a model with 24 layers with a hidden dimension size of 2048 and 32 attention heads, totalling to 1.3B parameters similar to MetaLM.
- They change the default transformer to the **magento**[10] module
- They use Extrapolatable position embedding (**xPos**)[11] which generalizes better at long term dependencies.

<b>Length</b>	<b>256</b>	<b>512</b>	<b>1024</b>	<b>2048</b>	<b>4096</b>
	Interpolation			Extrapolation	
Transformer	46.34	36.39	29.94	132.63	1283.79
Alibi	37.66	29.92	24.99	23.14	24.26
Roformer	38.09	30.38	25.52	73.6	294.45
LEX Transformer (Ours)	<b>34.3</b>	<b>27.55</b>	<b>23.31</b>	<b>21.6</b>	<b>20.73</b>

# Going from MetaLM to Kosmos-1

- Kosmos-1 trains a model with 24 layers with a hidden dimension size of 2048 and 32 attention heads, totalling to 1.3B parameters similar to MetaLM.
- They change the default transformer to the **magenta**[10] module.
- They use Extrapolatable position embedding (**xPos**)[11] which generalizes better at long term dependencies.
- Trains with the semi causal next token prediction task, minimizing log likelihood.

# Kosmos-1: Input Format

## Text:

`<s>` Kosmos-1 can perceive multimodal input, learn in context, and generate output. `</s>`

## Image-Caption:

`<s>` `<image>` Image Embedding `</image>` WALL-E giving potted plant to EVE. `</s>`

## Multimodal:

`<s>` `<image>` Image Embedding `</image>` This is WALL-E. `<image>` Image Embedding `</image>` This is EVE. `</s>`

# Kosmos - 1: Other Details

- **Vision Encoder:** CLIP ViT-L/14 that has been frozen except for the last layer.
- **Image preprocessing:** resized to 224 x 224
- **Tokenizer:** SentencePiece
- **Optimizer:** AdamW
- **Batch Size:** 1.2M tokens (0.5M tokens from text corpora, 0.5M tokens from image-caption pairs, and 0.2M tokens from interleaved data)

# Kosmos - 1 : Training Data

- **TextData:** Subset of the PILE dataset and common Crawl.

Datasets	Tokens (billion)	Weight (%)	Epochs
OpenWebText2	14.8	21.8%	1.47
CC-2021-04	82.6	17.7%	0.21
Books3	25.7	16.2%	0.63
CC-2020-50	68.7	14.7%	0.21
Pile-CC	49.8	10.6%	0.21
Realnews	21.9	10.2%	0.46
Wikipedia	4.2	5.4%	1.29
BookCorpus2	1.5	1.1%	0.75
Gutenberg (PG-19)	2.7	1.0%	0.38
CC-Stories	5.3	1.0%	0.19
NIH ExPorter	0.3	0.2%	0.75

277.5B Tokens

# Kosmos - 1 : Training Data

- **TextData:** Subset of the PILE dataset and common Crawl.
- **Image-caption pairs:** Collection of Several image caption datasets

Datasets	Tokens (billion)	Weight (%)	Epochs
OpenWebText2	14.8	21.8%	1.47
CC-2021-04	82.6	17.7%	0.21
Books3	25.7	16.2%	0.63
CC-2020-50	68.7	14.7%	0.21
Pile-CC	49.8	10.6%	0.21
Realnews	21.9	10.2%	0.46
Wikipedia	4.2	5.4%	1.29
BookCorpus2	1.5	1.1%	0.75
Gutenberg (PG-19)	2.7	1.0%	0.38
CC-Stories	5.3	1.0%	0.19
NIH ExPorter	0.3	0.2%	0.75

277.5B Tokens



LAION 2B

LAION 400M

COYO 700M

Conceptual Captions

3.1B pairs



# Kosmos - 1 : Training Data

- **TextData:** Subset of the PILE dataset and common Crawl.
- **Image-caption pairs:** Collection of Several image caption datasets
- **Interleaved Data:** documents containing images with text.

Datasets	Tokens (billion)	Weight (%)	Epochs
OpenWebText2	14.8	21.8%	1.47
CC-2021-04	82.6	17.7%	0.21
Books3	25.7	16.2%	0.63
CC-2020-50	68.7	14.7%	0.21
Pile-CC	49.8	10.6%	0.21
Realnews	21.9	10.2%	0.46
Wikipedia	4.2	5.4%	1.29
BookCorpus2	1.5	1.1%	0.75
Gutenberg (PG-19)	2.7	1.0%	0.38
CC-Stories	5.3	1.0%	0.19
NIH ExPorter	0.3	0.2%	0.75

277.5B Tokens



LAION 2B

LAION 400M

COYO 700M

Conceptual Captions

3.1B pairs

Interleaved data

71M documents

# Kosmos - 1 : Training Data

- **TextData:** Subset of the PILE dataset and common Crawl.
- **Image-caption pairs:** Collection of Several image caption datasets
- **Interleaved Data:** documents containing images with text.
- **Language only instruction data:** Unnatural Instructions and FLANv2

Datasets	Tokens (billion)	Weight (%)	Epochs
OpenWebText2	14.8	21.8%	1.47
CC-2021-04	82.6	17.7%	0.21
Books3	25.7	16.2%	0.63
CC-2020-50	68.7	14.7%	0.21
Pile-CC	49.8	10.6%	0.21
Realnews	21.9	10.2%	0.46
Wikipedia	4.2	5.4%	1.29
BookCorpus2	1.5	1.1%	0.75
Gutenberg (PG-19)	2.7	1.0%	0.38
CC-Stories	5.3	1.0%	0.19
NIH ExPorter	0.3	0.2%	0.75

277.5B Tokens



LAION 2B

LAION 400M

COYO 700M

Conceptual Captions

3.1B pairs

Interleaved data

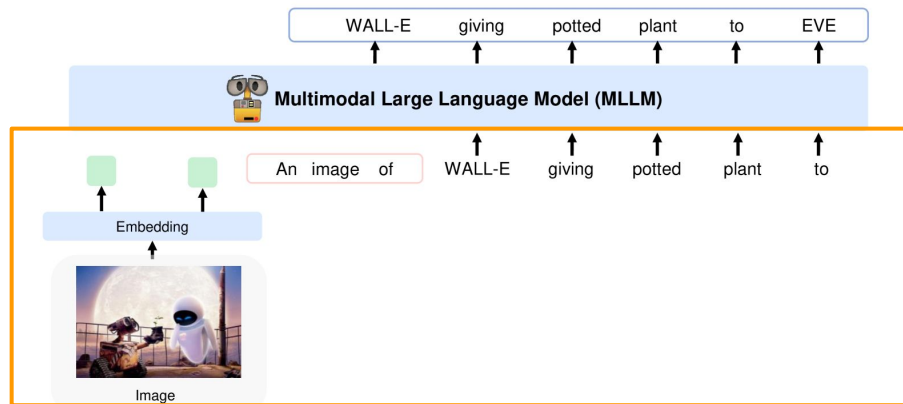
71M documents

Unnatural  
Instructions

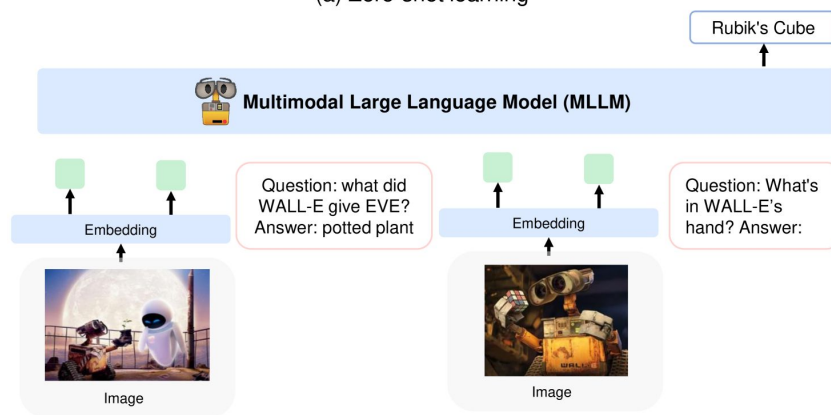
FLANv2

122K pairs

# Kosmos-1: Evaluation Format



(a) Zero-shot learning



(b) Few-shot learning

# Evaluation Tasks

Dataset	Task description	Metric	Zero-shot	Few-shot
<i>Language tasks</i>				
StoryCloze [34]	Commonsense reasoning	Accuracy	✓	✓
HellaSwag [61]	Commonsense NLI	Accuracy	✓	✓
Winograd [28]	Word ambiguity	Accuracy	✓	✓
Winogrande [40]	Word ambiguity	Accuracy	✓	✓
PIQA [8]	Physical commonsense	Accuracy	✓	✓
BoolQ [11]	Question answering	Accuracy	✓	✓
CB [16]	Textual entailment	Accuracy	✓	✓
COPA [37]	Causal reasoning	Accuracy	✓	✓
Rendered SST-2 [38]	OCR-free sentiment classification	Accuracy	✓	
HatefulMemes [25]	OCR-free meme classification	ROC AUC	✓	
<i>Cross-modal transfer</i>				
RelativeSize [5]	Commonsense reasoning (object size)	Accuracy	✓	
MemoryColor [36]	Commonsense reasoning (object color)	Accuracy	✓	
ColorTerms [4]	Commonsense reasoning (object color)	Accuracy	✓	
<i>Nonverbal reasoning tasks</i>				
IQ Test	Raven's Progressive Matrices	Accuracy	✓	
<i>Perception-language tasks</i>				
COCO Caption [32]	Image captioning	CIDEr, etc.	✓	✓
Flicker30k [60]	Image captioning	CIDEr, etc.	✓	✓
VQAv2 [18]	Visual question answering	VQA acc.	✓	✓
VizWiz [19]	Visual question answering	VQA acc.	✓	✓
WebSRC [14]	Web page question answering	F1 score	✓	
<i>Vision tasks</i>				
ImageNet [15]	Zero-shot image classification	Top-1 acc.	✓	
CUB [51]	Zero-shot image classification with descriptions	Accuracy	✓	

# Image Captioning

Kosmos-1 is able to outperform both 3B and 9B Flamingo models while being only 1.6B

Model	COCO		Flickr30k	
	CIDEr	SPICE	CIDEr	SPICE
ZeroCap	14.6	5.5	-	-
VLKD	58.3	13.4	-	-
FewVLM	-	-	31.0	10.0
META LM	82.2	15.7	43.4	11.7
Flamingo-3B*	73.0	-	60.6	-
Flamingo-9B*	79.4	-	61.5	-
<b>KOSMOS-1 (1.6B)</b>	<b>84.7</b>	<b>16.8</b>	<b>67.1</b>	<b>14.5</b>
Flamingo (80B)	84.3	-	<b>67.2</b>	-

Zero Shot results on COCO caption karpathy split

# Image Captioning

Kosmos-1 Performs degrades when number of shots are increased from 4

Model	COCO			Flickr30k		
	$k = 2$	$k = 4$	$k = 8$	$k = 2$	$k = 4$	$k = 8$
Flamingo-3B	-	85.0	90.6	-	72.0	71.7
Flamingo-9B	-	93.1	<b>99.0</b>	-	72.6	<b>73.4</b>
<b>KOSMOS-1 (1.6B)</b>	<b>99.6</b>	<b>101.7</b>	96.7	<b>70.0</b>	<b>75.3</b>	68.0

	$K = 32$	$k = 32$
Flamingo - 3B	99.0	71.2
Flamingo - 9B	106.3	72.8

Zero Shot results on COCO caption  
karpathy split

# Visual Question answering

Kosmos-1 outperforms Flamingo on  
VizWiz

<b>Model</b>	<b>VQAv2</b>	<b>VizWiz</b>
Frozen	29.5	-
VLKDViT-B/16	38.6	-
METALM	41.1	-
Flamingo-3B*	49.2	28.9
Flamingo-9B*	<b>51.8</b>	28.8
<b>KOSMOS-1 (1.6B)</b>	51.0	<b>29.2</b>
Flamingo (80B)	<b>56.3</b>	<b>31.6</b>

# Visual Question Answering

Kosmos - 1 does better in k=2,4 but does poorly with higher K

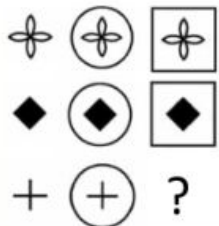
Model	VQAv2			VizWiz		
	$k = 2$	$k = 4$	$k = 8$	$k = 2$	$k = 4$	$k = 8$
Frozen	-	38.2	-	-	-	-
METALM	-	45.3	-	-	-	-
Flamingo-3B	-	53.2	55.4	-	34.4	38.4
Flamingo-9B	-	<b>56.3</b>	<b>58.0</b>	-	34.9	<b>39.4</b>
KOSMOS-1 (1.6B)	<b>51.4</b>	51.8	51.4	<b>31.4</b>	<b>35.3</b>	39.0

	$k = 32$	$k = 32$
Flamingo - 3B	57.1	<b>45.5</b>
Flamingo - 9B	<b>60.4</b>	44.0



# IQ- Test

Example of IQ Test



Which option can complete the matrix?



# IQ-Test

Example of IQ Test

		?

Which option can complete the matrix?

A B C D E F

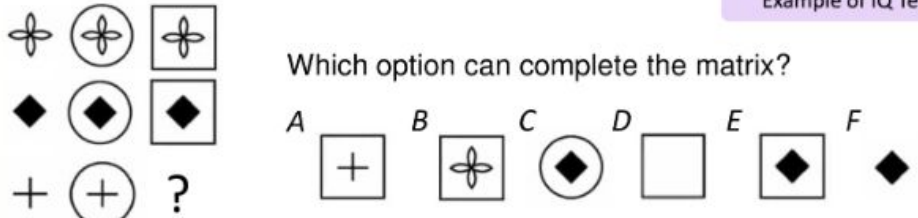
Input Prompt

Here are eight images: The following image is:







Is it correct?	Is it correct?	Is it correct?	Is it correct?	Is it correct?	Is it correct?
Yes	Yes	Yes	Yes	Yes	Yes

# IQ-Test









Example of IQ Test









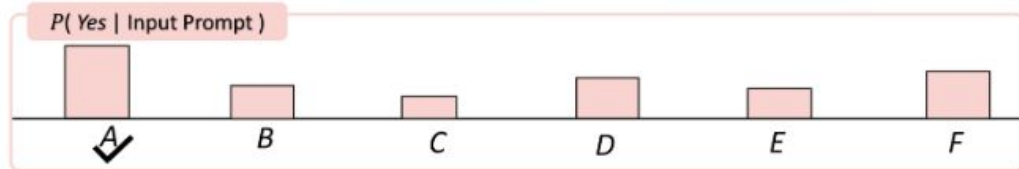
Which option can complete the matrix?

A  B  C  D  E  F 

Input Prompt

Here are eight images:         The following image is:

					
Is it correct?	Is it correct?	Is it correct?	Is it correct?	Is it correct?	Is it correct?
Yes	Yes	Yes	Yes	Yes	Yes



# IQ-Test

Example of IQ Test

Which option can complete the matrix?

A B C D E F

Input Prompt

Here are eight images: The following image is:

Is it correct?  
Yes

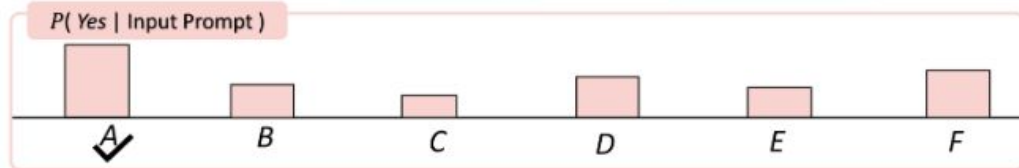
Is it correct?  
Yes

Is it correct?  
Yes

Is it correct?  
Yes

Is it correct?  
Yes

Is it correct?  
Yes



Method	Accuracy
Random Choice	17%
KOSMOS-1	22%
w/o language-only instruction tuning	26%

Results

# OCR Free Language understanding



“Question: does this picture contain real hate speech? Answer: {answer}”

It's clear the filmmakers weren't sure where they wanted their story to go, and even more clear that they lack the skills to get us to this undetermined destination.

“Question: what is the sentiment of the opinion? Answer: {answer}”

# OCR Free Language understanding

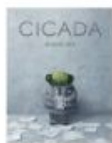
Kosmos -1 models dont uses OCR text  
whereas flamingo models do use  
them

<b>Model</b>	<b>HatefulMemes</b>	<b>Rendered SST-2</b>
CLIP ViT-B/32	57.6	59.6
CLIP ViT-B/16	61.7	59.8
CLIP ViT-L/14	63.3	64.0
Flamingo-3B	53.7	-
Flamingo-9B	57.0	-
<b>Kosmos-1 (1.6B)</b>	<b>63.9</b>	<b>67.1</b>
Flamingo (80B)	46.4	-

# WEb page Question answering

NEW RELEASES

[View all](#)



Cicada

Shaun Tan

\$26.99



Siege

Deborah Snow

\$32.99



The Other Wife

Michael Robotham

\$32.99



Kings of the Road: 50 Cars...

Toby Hagon

\$49.99



The Hate U Give

Angie Thomas

\$17.99



“Given the context below from web page, extract the answer from the given text like this:

Question: Who is the publisher of this book?

Answer: Penguin Books Ltd. Context: {WebText}

Q: {question} A: {answer} ”

# WEb page Question answering

NEW RELEASES

[View all](#)



Cicada

Shaun Tan

~~\$26.99~~



Siege

Deborah Snow

~~\$32.99~~



The Other Wife

Michael Robotham

~~\$32.99~~



Kings of the Road: 50 Cars...

Toby Hagon

~~\$49.99~~



The Hate U Give

Angie Thomas

~~\$17.99~~



“Given the context below from web page, extract the answer from the given text like this:

Question: Who is the publisher of this book?

Answer: Penguin Books Ltd. Context: {WebText}

Q: {question} A: {answer} ”

Models	EM	F1
<i>Using extracted text</i>		
LLM	7.6	17.9
KOSMOS-1	<b>15.8</b>	<b>31.3</b>
<i>Without using extracted text</i>		
KOSMOS-1	3.8	10.6



# Multimodal chain of thought prompting

Input Prompt



Question: which studio created the movie characters in the image?  
Answer:

Completion


Disney



Standard Prompting

# Multimodal chain of thought prompting

**Input Prompt**



Question: which studio created the movie characters in the image?  
Answer:

**Completion**

Disney

**X**

Standard Prompting

## Chain-of-Thought Prompting

### Model Input

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: Roger started with 5 balls. 2 cans of 3 tennis balls each is 6 tennis balls.  $5 + 6 = 11$ . The answer is 11.

Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?


### Model Output

A: The cafeteria had 23 apples originally. They used 20 to make lunch. So they had  $23 - 20 = 3$ . They bought 6 more apples, so they have  $3 + 6 = 9$ . The answer is 9. ✓

Chain of thought prompting for language

# Multimodal chain of thought prompting

**Input Prompt**



Question: which studio created the movie characters in the image?  
Answer:


**Completion**

Disney

**X**

Standard Prompting

**Input Prompt**




Introduce this picture in detail:

**Completion**

It is a robot that is holding up a cube that is the same color as the wall. The image is from the film, WALL-E, by Pixar Animation Studios. It is the story of a small robot named WALL-E and his quest to build the ultimate life form. WALL-E is created by a group...

# Multimodal chain of thought prompting

**Input Prompt**



Question: which studio created the movie characters in the image?  
Answer:


**Completion**

Disney

✗

Standard Prompting

**Input Prompt**




Introduce this picture in detail:

**Completion**

It is a robot that is holding up a cube that is the same color as the wall. The image is from the film, WALL-E, by Pixar Animation Studios. It is the story of a small robot named WALL-E and his quest to build the ultimate life form. WALL-E is created by a group...

**Input Prompt**



It is a robot that is holding up a cube that is the same color as the wall. The image is from the film, WALL-E, by Pixar Animation Studios. It is the story...  
Question: which studio created the movie characters in the image?  
Answer:

**Completion**

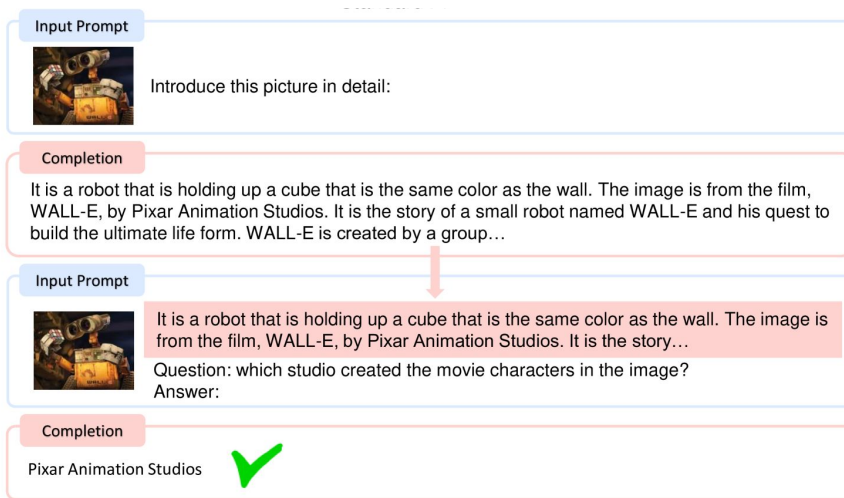
Pixar Animation Studios

✓

Multimodal chain of thought

# Multimodal chain of thought prompting

Models	Accuracy
CLIP ViT-B/32	59.6
CLIP ViT-B/16	59.8
CLIP ViT-L/14	64.0
Kosmos-1	67.1
w/ multimodal CoT prompting	<b>72.9</b>



Multimodal chain of thought

# Zero shot image classification



Input is "The photo of the" and  
output is constrained to 1K classes

# Zero shot image classification



Input is "The photo of the" and  
output is constrained to 1K classes

Model	Without Constraints	With Constraints
GIT [58]	1.9	33.5
KOSMOS-1	<b>4.0</b>	<b>38.1</b>

# Zero shot Image classification with Descriptions

## Input Prompt



Question: what is the name of the woodpecker in the picture, three toed or downy?  
Answer:

## Completion

downy woodpecker



## Input Prompt

Description of three toed woodpecker: It has black and white stripes throughout the body and a yellow crown.  
Description of downy woodpecker: It has white spots on its black wings and some red on its crown.



Question: what is the name of the woodpecker in the picture?  
Answer:

## Completion

three toe woodpecker





# Zero shot Image classification with Descriptions

## Input Prompt



Question: what is the name of the woodpecker in the picture, three toed or downy?  
Answer:

## Completion

downy woodpecker



## Input Prompt

Description of three toed woodpecker: It has black and white stripes throughout the body and a yellow crown.  
Description of downy woodpecker: It has white spots on its black wings and some red on its crown.



Question: what is the name of the woodpecker in the picture?  
Answer:

## Completion

three toe woodpecker



Settings	Accuracy
Without Descriptions	61.7
With Descriptions	90.0

# Language Tasks

MLLM trained on the text corpora as an LLM is the baseline

Task	Zero-shot		One-shot		Few-shot ( $k = 4$ )	
	LLM	KOSMOS-1	LLM	KOSMOS-1	LLM	KOSMOS-1
StoryCloze	<b>72.9</b>	72.1	<b>72.9</b>	72.2	<b>73.1</b>	72.3
HellaSwag	<b>50.4</b>	50.0	<b>50.2</b>	50.0	<b>50.4</b>	50.3
Winograd	<b>71.6</b>	69.8	<b>71.2</b>	68.4	<b>70.9</b>	69.8
Winogrande	<b>56.7</b>	54.8	<b>56.7</b>	54.5	<b>57.0</b>	55.7
PIQA	<b>73.2</b>	72.9	<b>73.0</b>	72.5	<b>72.6</b>	72.3
BoolQ	<b>56.4</b>	<b>56.4</b>	55.1	<b>57.2</b>	58.7	<b>59.2</b>
CB	39.3	<b>44.6</b>	41.1	<b>48.2</b>	42.9	<b>53.6</b>
COPA	<b>68.0</b>	63.0	<b>69.0</b>	64.0	<b>69.0</b>	64.0
<b>Average</b>	61.1	60.5	61.2	60.9	61.8	62.2

# Cross Modal Transfer

Comparing Kosmos-1 trained with and without language-only instruction tuning on V+L tasks.

<b>Model</b>	<b>COCO</b>	<b>Flickr30k</b>	<b>VQAv2</b>	<b>VizWiz</b>
KOSMOS-1	84.7	<b>67.1</b>	<b>51.0</b>	<b>29.2</b>
w/o language-only instruction tuning	<b>87.6</b>	65.2	46.7	27.9

# Cross Modal Transfer

Comparing Kosmos-1 trained with and without language-only instruction tuning on V+L tasks.

Model	COCO	Flickr30k	VQAv2	VizWiz
KOSMOS-1	84.7	<b>67.1</b>	<b>51.0</b>	<b>29.2</b>
w/o language-only instruction tuning	<b>87.6</b>	65.2	46.7	27.9

Comparing Kosmos-1 and LLM baseline on common sense reasoning tasks.

Model	Size Reasoning	Color Reasoning	
	RELATIVE SIZE	MEMORY COLOR	COLOR TERMS
<i>Using retrieved images</i>			
VALM [53]	85.0	58.6	52.7
<i>Language-only zero-shot evaluation</i>			
LLM	92.7	61.4	63.4
KOSMOS-1	<b>94.2</b>	<b>76.1</b>	<b>73.1</b>

# Strengths

- Kosmos - 1 uses fewer parameters than Flamingo models but is competitive on results and often outperforms.
- Unlike Flamingo their zero shot evaluation methods dont use two mock demonstrations.
- They train their models on open source datasets.

# Weaknesses

- Kosmos - 1 seems to be scaling poorly to higher K, also it should be noted that they haven't done experiments for K = 32 a setting at which Flamingo does best and can outperform kosmos-1.
- Kosmos - 1 seems to be doing poorly on classification tests, when looking at imagenet results and it performs well only with instructions, which may not be available at all times.
- Even though kosmos - 1 can ingest interleaved multimodal input, they have not performed any experiments around video reasoning tasks.

# Future Work

- Authors can try to scale the models so that their sizes can be comparable to Flamingo for an apples to apples comparison.
- Authors can try to incorporate more modalities into their model, such as video etc.
- Authors can look into distillation methods to create smaller models when they scale up.

# Discussion

- Do you think not releasing models and datasets to the public hurts future research ?
- IQ tests seem like a effective method to learn how close a model is to human intelligence. Can we train models that do well on theses tests, have logic and pattern recognition skills?



# References

1. Brown, Tom, et al. "Language models are few-shot learners." *Advances in neural information processing systems* 33 (2020): 1877-1901.
2. Hao, Yaru, et al. "Language models are general-purpose interfaces." *arXiv preprint arXiv:2206.06336* (2022).
3. Alayrac, Jean-Baptiste, et al. "Flamingo: a visual language model for few-shot learning." *arXiv preprint arXiv:2204.14198* (2022).
4. Huang, Shaohan, et al. "Language Is Not All You Need: Aligning Perception with Language Models." *arXiv preprint arXiv:2302.14045* (2023).
5. Hoffmann, Jordan, et al. "Training compute-optimal large language models." *arXiv preprint arXiv:2203.15556* (2022).
6. Tsimpoukelli, Maria, et al. "Multimodal few-shot learning with frozen language models." *Advances in Neural Information Processing Systems* 34 (2021): 200-212.
7. Aghajanyan, Armen, et al. "Cm3: A causal masked multimodal model of the internet." *arXiv preprint arXiv:2201.07520* (2022).
8. Brock, Andy, et al. "High-performance large-scale image recognition without normalization." *International Conference on Machine Learning*. PMLR, 2021.
9. Wang, Wenhui, et al. "Image as a foreign language: Beit pretraining for all vision and vision-language tasks." *arXiv preprint arXiv:2208.10442* (2022).
10. Wang, Hongyu, et al. "Foundation transformers." *arXiv preprint arXiv:2210.06423* (2022).
11. Sun, Yutao, et al. "A Length-Extrapolatable Transformer." *arXiv preprint arXiv:2212.10554* (2022).