# Flamingo

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

J-B. Alayrac, J. Donahue, P. Luc, A. Miech, I. Barr, Y. Hasson, K. Lenc, A. Mensch, K. Millican, M. Reynolds, R. Ring, E. Rutherford, S. Cabi, T. Han, Z. Gong, S. Samangooei, M. Monteiro, J. Menick, S. Borgeaud, A. Brock, A. Nematzadeh, S. Sharifzadeh, M. Binkowski, R. Barreira, O. Vinyals, A. Zisserman, K. Simonyan arxiv (2022)

**Digest** by Samuel Albanie, May 2022





# Outline

- Motivation and challenges for multimodal generative modelling
- Related Work
- Flamingo model
- Evaluation datasets
- Comparison to state-of-the-art
- Ablations
- Qualitative Examples
- Discussion (limitations, trade-offs, opportunities, benefits, risks)

# Motivation

# The few-shot dream

Aspect of intelligence: ability to quickly learn a task given short instruction

- Fast acquisition of categories in children (Markman et al., 1989)
- Model learning environment to make better use of data (Griffiths et al. 2019)

We'd like multimodal systems (vision and language) that achieve this property

### **Dominant computer vision paradigm:**

large-scale pretraining

+

task-specific fine-tuning

But current <u>fine-tuning</u> approaches often require:

- thousands of training samples
- careful per-task hyperparameter tuning
- significant computational resource

Can we train a multimodal model to work well in a "few-shot" regime?

### References

E. Markman, "Categorization and naming in children: Problems of induction", MIT Press (1989)
T. L. Griffiths, et al., "Doing more with less: meta-reasoning and meta-learning in humans and machines", Current Opinion in Behavioral Sciences (2019) (CLIP) A. Radford et al., "Learning transferable visual models from natural language supervision", ICML (2021) (ALIGN) C. Jia et al., "Scaling up visual and vision-language representation learning with noisy text supervision", ICML (2021) (VL-T5) J. Cho et al., "Unifying vision-and-language tasks via text generation", ICML (2021)

# Open-ended task abilities

Multimodal models like CLIP and ALIGN have shown promising zero-shot performance, but they are <u>inflexible</u>: they lack the ability to <u>generate language</u> Flexible models for visually-conditioned language generation like VL-T5 exist But these have not demonstrated strong few-shot performance

**Inspiration from NLP:** large language models like GPT-3 are flexible few-shot learners Given a few examples of a task as a prompt + query input, the language model generates a continuation to produce a predicted output.

A key factor of their success is large-scale pretraining.

In principle: image/video understanding tasks (e.g. classification, captioning, questionanswering) are text prediction problems with visual input conditioning.

Can we we learn a model capable of open-ended multimodal tasks via pretraining?



# Challenges for multimodal generative modelling

# Unifying strong unimodal models

Training large language models is extremely

computationally expensive

We'd like to save computational resources by starting from a pretrained language model But a text-only model has no built-in way to incorporate input from other modalities. We want to enable this while retaining the knowledge of the original language model

**Proposed approach:** interleave cross-attention layers with language-only self-attention layers (frozen)

# Supporting images and videos

**Goal:** enable both images and video inputs These are high-dimensional, so flattening to 1D sequences (as used in text-generation) is costly Exacerbated by quadratic cost of self-attention

Secondary goal: would also like a unified treatment of images and video

**Proposed approach:** Perceiver-based architecture with a fixed number of visual tokens

### References

A. Jaegle et al., "Perceiver: General perception with iterative attention", ICML (2021) (CLIP) A. Radford et al., "Learning transferable visual models from natural language supervision", ICML (2021) (ALIGN) C. Jia et al., "Scaling up visual and vision-language representation learning with noisy text supervision", ICML (2021)

## Heterogeneous training data

Large models require vast training datasets.

Existing (image, text) datasets used by (e.g. used by CLIP and ALIGN) my not be general enough to reach GPT-3 style few-shot learning.

Large internet-based text-only datasets exist, but not for multimodal data.

One scalable approach: scrape web pages with interleaved images and text.

But such images and text are often only weakly related **Proposed approach:** combine web scraping with

existing paired (image, text) and (video, text) datasets



# Outline

- Motivation and challenges for multimodal generative modelling
- Related Work
- Flamingo model
- Evaluation datasets
- Comparison to state-of-the-art
- Ablations
- Qualitative Examples
- Discussion (limitations, trade-offs, opportunities, benefits, risks)

# **Related Work**

# Brown et al., (2020)



# Tsimpoukelli et al., (2021)



### **References/Image credits**

T. Brown et al., "Language models are few-shot learners", NeurIPS (2020)

R. Mokady et al., "Clipcap: Clip prefix for image captioning", arxiv (2021)

(CLIP) A. Radford et al., "Learning transferable visual models from natural language supervision", ICML (2021)

# Mokady et al., (2021)

Key idea: Leverage large trained models (CLIP, GPT-2) for captioning

![](_page_5_Figure_11.jpeg)

Achieves solid captioning performance at low computational cost

ClipCap

## Aghajanyan et al., (2022)

![](_page_5_Figure_15.jpeg)

(GPT-2) A. Radford et al., "Language models are unsupervised multitask learners", (2019) M. Tsimpoukelli et al., "Multimodal few-shot learning with frozen language models", NeurIPS (2021) A. Aghajanyan et al., "CM3: A Causal Masked Multimodal Model of the Internet", arxiv (2022)

# Outline

- Motivation and challenges for multimodal generative modelling
- Related Work
- Flamingo model
- Evaluation datasets
- Comparison to state-of-the-art
- Ablations
- Qualitative Examples
- Discussion (limitations, trade-offs, opportunities, benefits, risks)

# Flamingo Model

![](_page_7_Figure_1.jpeg)

### **References/Image credits**

J-B. Alayrac et al., "Flamingo: a Visual Language Model for Few-Shot Learning", arxiv (2022) (CLIP) A. Radford et al., "Learning transferable visual models from natural language supervision", ICML (2021) (Chinchilla) J. Hoffmann et al., "Training Compute-Optimal Large Language Models", arxiv (2022)

![](_page_7_Figure_5.jpeg)

# Vision encoder: pixels to features

![](_page_8_Figure_1.jpeg)

### **References/Image credits**

J-B. Alayrac et al., "Flamingo: a Visual Language Model for Few-Shot Learning", arxiv (2022) (NFNet) A. Brock et al., "High-performance large-scale image recognition without normalization", ICML (2021) (CLIP) A. Radford et al., "Learning transferable visual models from natural language supervision", ICML (2021) (BERT) J. Devlin et al. "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding", NAACL-HLT (2019)

### Flamingo vision encoder

Vision encoder: F6 Normalizer-Free ResNet (NFNet) backbonePretrained as dual encoder using contrastive loss employed by CLIPBERT is used for the text encoder (discarded after pretraining)Slight difference to CLIP: global average pooling is used to producethe vision embedding (rather than global attention pooling)Resolution288 x 288 pixelsEmbedding1376

Outputs 2D spatial grid of features which is flattened to 1D

For videos: frames are sampled at 1FPS (features are concatenated)

Vision encoder is frozen after pretraining

# Vision encoder details

![](_page_9_Figure_1.jpeg)

# Optimisation details

Trained on 512 TPUv4 chips using Adam optimiser

Batch size of 16,384 (fairly large)

Colour augmentation and random horizontal flips during training

Both the vision encoder and text encoder are trained from scratch

Monitor training progress on zero-shot image classification (like

CLIP, this is done with a prompt template "A photo of a {class}")

### **References/Image credits**

J-B. Alayrac et al., "Flamingo: a Visual Language Model for Few-Shot Learning", arxiv (2022)
(Adam) D. P. Kingma and J. Ba. "Adam: A Method for Stochastic Optimization", ICLR (2015)
(CLIP) A. Radford et al., "Learning transferable visual models from natural language supervision", ICML (2021)
(ALIGN) C. Jia et al., "Scaling up visual and vision-language representation learning with noisy text supervision", ICML (2021)
(NFNet) A. Brock et al., "High-performance large-scale image recognition without normalization", ICML (2021)
(BERT) J. Devlin et al. "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding", NAACL-HLT (2019)

# Pretraining data

Trained on a combination of two internal (image, text) datasets:

ALIGN (1.8 billion) - noisy

LTIP (312 million) - cleaner, longer descriptions

The manner of combination is important for performance

(Ablation study) small NFNet-F0 with BERT-mini for different regimes:

Dataset	Combination		nageNet			COO	20		
	strategy	a	ccuracy	in	image-to-text text-to-in				
			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

Accumulation: compute gradient on batch from each dataset, combine via weighted sum

Data merged: merge examples from each dataset into each batch

Round-robin: alternate batches from each dataset, update parameters each batch

![](_page_9_Picture_21.jpeg)

# Perceiver resampler

### From large, variable-size features to fixed # tokens

A variable number of input frames are processed (for videos) The vision encoder thus produces a variable number of features It outputs a fixed number of visual tokens (64) to limit complexity Temporal encodings are added to visual inputs (spatial grid position encodings are not use, since they did not help) The results are then flattened to form a 1D sequence These are combined with a fixed set of learned latent queries (64) Both are processed by attention and feed-forward layers. **Note:** differently to DETR and Perceiver, keys and values for latent queries are concatenated to those from the visual embeddings.

### **References/Image credits**

J-B. Alayrac et al., "Flamingo: a Visual Language Model for Few-Shot Learning", arxiv (2022) (DETR) N. Carion et al., "End-to-end object detection with transformers", ECCV (2020) (Perceiver) A. Jaegle et al., "Perceiver: General perception with iterative attention", ICML (2021)

![](_page_10_Picture_5.jpeg)

## Perceiver resampler module

![](_page_10_Figure_7.jpeg)

# Conditioning the language model

![](_page_11_Figure_1.jpeg)

![](_page_11_Figure_2.jpeg)

## Interleaving gated cross-attention layers

Language models: frozen Chinchillas (trained on MassiveText)

Gated xattn dense blocks (trained from scratch) are inserted between layers

Each block includes cross attention

feed-forward

Layer norm is applied to all attention inputs and the feed-forward layers (GPT-2 style)

Use tanh gates to preserve original language model behaviour at initialisation

Each  $tanh(\alpha)$  gate controlled via a layer-specific learnable scalar  $\alpha$  (initialised to zero)

### Architecture integration

Gated xattn-dense insertions are chosen according to the selected language model

Differing insertion frequencies represent different compute-performance trade-offs

### **References/Image credits**

J-B. Alayrac et al., "Flamingo: a Visual Language Model for Few-Shot Learning", arxiv (2022) (MassiveText) J. Rae, et al. "Scaling language models: Methods, analysis & insights from training gopher", arxiv (2021) (Chinchilla) J. Hoffmann et al., "Training Compute-Optimal Large Language Models", arxiv (2022) (Layer norm) J. Ba et al., "Layer normalization", arxiv (2016) (GPT-2) A. Radford et al., "Language models are unsupervised multitask learners" (2019) (tanh gates) S. Hochreiter et al., "Long short-term memory", Neural computation (1997)

![](_page_11_Picture_16.jpeg)

![](_page_11_Picture_17.jpeg)

# Per-image/video attention masking

![](_page_12_Figure_1.jpeg)

Multi-image attention is implemented with the gated xattn-dense layers with causal masking over tokens from the perceiver resampler. By default, each token only allowed to attend to the visual tokens of the image that appeared immediately before it (this restriction improved performance) Note: Although direct attention is over a single image, there is still a causal dependency on previous images (due to causal self-attention in the text decoder) Experiments show that the model can train on 5 images, but generalise up to 32. Restriction may be a useful inductive bias for single image tasks.

### **References/Image credits**

J-B. Alayrac et al., "Flamingo: a Visual Language Model for Few-Shot Learning", arxiv (2022)

### Multi-image attention

![](_page_12_Picture_10.jpeg)

# Flamingo - training data

Flamingo is trained on:

- Image-Text Pairs data
- Video-Text Pairs data
- Webpage data

![](_page_13_Picture_6.jpeg)

![](_page_13_Picture_7.jpeg)

![](_page_13_Picture_8.jpeg)

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

![](_page_13_Figure_10.jpeg)

### **References/Image credits**

J-B. Alayrac et al., "Flamingo: a Visual Language Model for Few-Shot Learning", arxiv (2022) (ALIGN) C. Jia et al., "Scaling up visual and vision-language representation learning with noisy text supervision", ICML (2021) (SafeSearch) <u>https://support.google.com/websearch/answer/510</u>

(Datasheets) T. Gebru et al., "Datasheets for datasets", Communications of the ACM (2021)

![](_page_13_Picture_14.jpeg)

### Data sources

![](_page_13_Picture_16.jpeg)

A kid doing a kickflip.

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

![](_page_13_Picture_19.jpeg)

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

## MultiModal Massive Web (M3W)

Extract text and images from 43 million webpages Use the DOM to determine the interleaving order of text and images M3W contains 185 million images and 182 GB of text Documents are filtered with Google SafeSearch filter Text filters: heuristics used to remove low quality documents & repetitions Image filters: small ( < 64 pixels), extreme aspect ratio, single-colour

> (Datasheets for LTIP, VTP, M3W) Train/val splits are randomly chosen internal datasets

Text

320 x 320 pixels

Resolution

![](_page_13_Figure_25.jpeg)

# Flamingo training objective

Models are trained with a weighted sum of dataset specific negative log likelihoods of text (conditioned on visual inputs):

![](_page_14_Figure_4.jpeg)

 $\mathcal{D}_m$  - *m*-th dataset

 $\lambda_m$  - positive scalar weight for the *m*-th dataset

Similarly to the vision encoder pretraining:

- tuning these weights is important for good performance
- the accumulation strategy is used

# Training details

$$\left[-\sum_{l=1}^{L}\log p(y_{l}|y_{< l}, x_{\le l})\right]$$

# Task adaptation with few-shot in-context learning

![](_page_15_Figure_1.jpeg)

### **References/Image credits**

(GPT-3) T. Brown et al., "Language models are few-shot learners", NeurIPS (2020) J-B. Alayrac et al., "Flamingo: a Visual Language Model for Few-Shot Learning", arxiv (2022)

# Few-shot in-context learning details

## Open-ended and close-ended evaluations

### Open-ended tasks

- the text generated following the query image is used as the prediction
- stop at the first <EOC> token
- beam search with beam size 3 is used

Close-ended tasks

- each target candidate is appended independently to the query image
- sequences are ranked by their log-likelihood

## Retrieval-based in-context example selection

It can be hard to leverage large numbers of support examples:

- it is expensive to fit all examples in the prompt
- generalisation may suffer if fewer examples were used in training

Prompt example selection (Liu et al., 2021) can address these issues

Retrieval-based In-Context Example Selection (RICES) (Yang et al., 2021)

Build prompt from top-N most similar examples to query

To avoid recency bias (Zhao et al., 2021), most similar example is put last

### **References/Image credits**

(CLIP) A. Radford et al., "Learning transferable visual models from natural language supervision", ICML (2021) J. Liu et al., "What Makes Good In-Context Examples for GPT-\$3?", arxiv (2021) (RICES) Z. Yang et al., "An empirical study of gpt-3 for few-shot knowledge-based vqa", arxiv (2021) E. Perez et al., "True few-shot learning with language models", NeurIPS (2021)

S. Min et al., "Rethinking the Role of Demonstrations: What Makes In-Context Learning Work?", arxiv (2022) Z. Zhao et al., "Calibrate before use: Improving few-shot performance of language models", ICML (2021)

## Zero-shot generalisation

If no examples available, one option is prompt engineering (CLIP) Performance is sensitive to the prompt, but validation requires examples Perez et al. (2021) shows that validation with few samples is not robust Flamingo: Build prompt from two downstream examples without images/video Using one example worked poorly (model predictions are very similar to the example) Min et al., (2022): (label space/text distribution/format matter, label correctness doesn't) For close-ended tasks, no text examples are required for the zero-shot prompt.

## Prompt ensembling

Ensembling across multiple prompts can be used to improve performance **Note:** This can be combined with **RICES** over different permutations of nearest neighbours For a given answer, log likelihoods are ensembled over six random permutations of the selected few-shot prompts

![](_page_16_Picture_24.jpeg)

# Flamingo Models

### Model architectures

Three sizes of Flamingo model are considered (building on three Chinchilla sizes)

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

The largest (80B) model requires model sharding All models use a NFNet-F6 backbone for the frozen vision encoder Gated xattn-dense layers are inserted at different frequencies (trading off memory and performance) The Perceiver Resampler remains the same across each model **Model card:** model intended for internal development.

### **References/Image credits**

(Chinchilla) J. Hoffmann et al., "Training Compute-Optimal Large Language Models", arxiv (2022) J-B. Alayrac et al., "Flamingo: a Visual Language Model for Few-Shot Learning", arxiv (2022) (NFNet) A. Brock et al., "High-performance large-scale image recognition without normalization", ICML (2021) (Perceiver) A. Jaegle et al., "Perceiver: General perception with iterative attention", ICML (2021) (Model card) M. Mitchell et al., "Model cards for model reporting", ACM FAccT (2019)

# Outline

- Motivation and challenges for multimodal generative modelling
- Related Work
- Flamingo model
- Evaluation datasets
- Comparison to state-of-the-art
- Ablations
- Qualitative Examples
- Discussion (limitations, trade-offs, opportunities, benefits, risks)

# Few-shot evaluation benchmarks

		Benc
	image benchmarks	
ImageNet-1K	object classification	unimodal
MS-COCO	scene description	generative
VQAv2	scene understanding QA	generative
OKVQA	external knowledge QA	generative
Flickr30k	scene description	generative
VizWiz	scene understanding QA	generative
TextVQA	text reading QA	generative
VisDial	visual dialogue	
HatefulMemes	meme classification	custom prompt

2 datasets are unimodal

12 benchmarks require open-ended generative sampling

![](_page_19_Figure_5.jpeg)

### hmark datasets

	video benchmarks	
Kinetics700	action classification	unimodal
VATEX	event description	generative
MSVDQA	event understanding QA	generative
YouCook2	event description	generative
MSRVTTQA	event understanding QA	generative
iVQA	event understanding QA	generative
RareAct	composite action retrieval	custom prompt
NextQA	temporal/causal QA	generative
STAR	multiple choice QA	

11 used for a less biased few-shot evaluation (includes less explored capabilities):

# Flamingo: dataset deduplication

### Data deduplication against evaluation tasks

Flamingo uses large-scale web-based pretraining Necessary investigate the possibility of evaluation dataset contamination CLIP - did not deduplicate (instead performed analysis) ALIGN did perform deduplication Flamingo uses an internal Google tool for deduplication Nearest neighbour search via visual embeddings to retrieve duplicates

### References

(CLIP) A. Radford et al., "Learning transferable visual models from natural language supervision", ICML (2021) (ALIGN) C. Jia et al., "Scaling up visual and vision-language representation learning with noisy text supervision", ICML (2021)

# Deduplication details

(Image-text) Deduplicated LTIP and ALIGN training images against:

![](_page_20_Figure_7.jpeg)

![](_page_20_Picture_9.jpeg)

# Flamingo: nuts and bolts training details

### Data augmentation and pre-processing

During training, 50% of text samples are prepended with a space character Effectiveness likely due subword tokenizer (tokens depend on preceding space) Augmentation enforces invariance to this artifact without degrading punctuation Visual inputs processed at 320 pixels (rather than 288 pixels used in pretraining) Inspired by FixRes - not too expensive since vision encoder is frozen On interleaved datasets, image indices  $\phi$  are also perturbed (next/prev prob. 0.5) For videos: clips of 8 frames (at 1 fps) are sampled from each training video In inference: 30 video frames processed at 3 fps (interpolating pos. embeddings)

### Loss and optimisation

All models are trained with AdamW

Optimisation is done with linear warmup followed by a flat learning rate Dataset mixing weights  $(\lambda_m)$ :

M3W
Align

### References

H. Touvron et al., "Fixing the train-test resolution discrepancy: FixEfficientNet", arxiv (2020)
D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization", ICLR (2015)
I. Loshchilov and F. Hutter, "Decoupled weight decay regularization", arxiv (2017)

(JAX) (Haiku (Mega (ZeRC Infrastructure/implementation

The model is trained using JAX and Haiku Training used TPUv4 instances Largest (80B) model trained for 15 days for 1536 chips over 16 devices Megatron sharding used for Embedding/S-Attention/X-Attenion/FFW ZeRO stage 1 is used to shard the optimiser state Activations + gradients: bfloat16, params + optim. accumulators: float32

- (JAX) J. Bradbury et al., "JAX: composable transformations of Python+ NumPy programs" (2018)
- (Haiku) T. Hennigan et al., "Haiku: Sonnet for Jax" (2020)
- (Megatron) M. Shoeybi et al., "Megatron-Im: Training multi-billion parameter language models using model parallelism", arxiv (2019) (ZeRO) S. Rajbhandari et al., "Zero: Memory optimizations toward training trillion parameter models", SC (2020)

![](_page_21_Picture_15.jpeg)

# Outline

- Motivation and challenges for multimodal generative modelling
- Related Work
- Flamingo model
- Evaluation datasets
- Comparison to state-of-the-art
- Ablations
- Qualitative Examples
- Discussion (limitations, trade-offs, opportunities, benefits, risks)

# Few-shot: comparison to SotA

![](_page_23_Figure_2.jpeg)

### **References/Image credits**

# Few-shot: further analysis

![](_page_24_Figure_1.jpeg)

### **References/Image credits**

J-B. Alayrac et al., "Flamingo: a Visual Language Model for Few-Shot Learning", arxiv (2022) (GPT-3) T. Brown et al., "Language models are few-shot learners", NeurIPS (2020) (Model soup) M. Wortsman et al., "Model soups: averaging weights of multiple fine-tuned models improves accuracy without increasing inference time", arxiv (2022) (MTV) S. Yan et al., "Multiview Transformers for Video Recognition", arxiv (2022) (BASIC) Pham et al., "Combined Scaling for Zero-shot Transfer Learning", arxiv (2021) (CLIP) A. Radford et al., "Learning transferable visual models from natural language supervision", ICML (2021)

# Few-shot classification on classification tasks

Model	Method	Prompt size	shots/class	ImageNet top 1	Kinetics700 avg top1/5
SotA	Fine-tuned	-	full	91.0 soup	89.0 MTV
SotA	Contrastive	-	0	85.7 BASIC	69.6 CLIP
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-tuned Sot	soup: ViT-G/14 (mo	del soup) MTV	/: Multiview Tran	sformer for Vide	eo Recognition
Zero-shot SotA	BASIC: 3B params	, 6.6B image-text	pairs	CLIP: ViT-L/14@	236 px

![](_page_24_Picture_7.jpeg)

# **Contrastive pretraining: zero-shot retrieval**

## **Retrieval benchmarks**

Evaluation of the pretrained dual encoder for zero-shot retrieval

			Flick	r30K			COCO					
	in	image-to-text text-to-image					in	nage-to-	text	text-to-image		
	R@1	R@5	R@10	R@1	R@5	R@10	R@1	R@5	R@10	R@1	R@5	R@10
Florence	90.9	99.1	-	76.7	93.6	-	64.7	85.9	-	47.2	71.4	-
ALIGN	88.6	98.7	<b>99.</b> 7	75.7	93.8	96.8	58.6	83.0	89.7	45.6	69.8	78.6
CLIP	88.0	98.7	99.4	68.7	90.6	95.2	58.4	81.5	88.1	37.7	62.4	72.2
Flamingo	89.3	98.8	99.7	79.5	95.3	97.9	65.9	87.3	92.9	48.0	73.3	82.1

**Observation:** training on short text descriptions improves ImageNet classification but harms text-image retrieval

Flamingo pretraining optimises for retrieval rather than classification to capture the whole scene in images.

### **References/Image credits**

J-B. Alayrac et al., "Flamingo: a Visual Language Model for Few-Shot Learning", arxiv (2022) (Florence) L. Yuan et al., "Florence: A New Foundation Model for Computer Vision", arxiv (2021) (ALIGN) C. Jia et al., "Scaling up visual and vision-language representation learning with noisy text supervision", ICML (2021) (CLIP) A. Radford et al., "Learning transferable visual models from natural language supervision", ICML (2021)

# Fine-tuning Flamingo

## Low-level details

During fine-tuning, Flamingo keeps the language model layers frozen Input image resolution is increased from 320 x 320 pixels to 480 x 480 pixels The base vision encoder is also fine-tuned (unlike Flamingo pretraining) Hyperparameters are set by grid search on validation subsets of the training sets Search over: learning rate, decay schedule, training steps, batch size, augmentation

Method

*Flamingo -* 32 shots SimVLM OFA Florence Flamingo Fine-tuned

Restricted SotA<sup>†</sup>

Unrestricted SotA

### **References/Image credits**

J-B. Alayrac et al., "Flamingo: a Visual Language Model for Few-Shot Learning", arxiv (2022) (SimVLM) Z. Wang et al., "Simvlm: Simple visual language model pretraining with weak supervision", arxiv (2021) (OFA) P. Wang et al., "Unifying Architectures, Tasks, and Modalities Through a Simple Sequence-to-Sequence Learning Framework", arxiv (2022) (Florence) L. Yuan et al., "Florence: A New Foundation Model for Computer Vision", arxiv (2021)

![](_page_26_Picture_9.jpeg)

### Fine-tuning comparison on 9 benchmarks

VOAV?		COCO	VATEX	VizWiz		MSRVTTQA		VisDial	YouCook2		TextVQA	
test-dev	test-std	test	test	test-dev	test-std	test	valid	test-std	valid	valid	test-std	test
67.6	-	113.8	65.1	49.8	-	31.0	56.8	-	86.8	36.0	-	70
80.0	80.3	143.3	-	-	-	-	-	-	-	-	-	
79.9	80.0	149.6	-	-	-	-	-	-	-	-	-	
80.2	80.4	-	-	-	-	-	-	-	-	-	-	
82.0	<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</u>
80.2	80.4	143.3	76.3	-	-	46.8	75.2	74.5	138.7	54.7	<u>73.7</u>	75
[ <b>150</b> ]	[ <b>150</b> ]	[ <b>134</b> ]	[ <mark>165</mark> ]	-	-	[ <mark>57</mark> ]	[ <mark>87</mark> ]	[ <mark>87</mark> ]	[ <mark>142</mark> ]	[147]	[ <mark>92</mark> ]	[6
81.3	81.3	149.6	81.4	57.2	60.6	-	-	75.4	-	-	-	84
[143]	[143]	[ <b>129</b> ]	[ <b>165</b> ]	[ <b>70</b> ]	[ <b>70</b> ]	-	-	[ <mark>133</mark> ]	-	-	-	[10

**Summary:** fine-tuning, while expensive, brings significant gains in performance

![](_page_26_Figure_13.jpeg)

# Outline

- Motivation and challenges for multimodal generative modelling
- Related Work
- Flamingo model
- Evaluation datasets
- Comparison to state-of-the-art
- Ablations
- Qualitative Examples
- Discussion (limitations, trade-offs, opportunities, benefits, risks)

![](_page_28_Figure_0.jpeg)

a V

\_\_\_\_ 20

\_\_\_\_\_ 21

1.0 — 23

22

24

Ablated

setting

(iv)

Cross-attention

architecture

Flamingo 3B

### **References/Image credits**

Attention

J-B. Alayrac et al., "Flamingo: a Visual Language Model for Few-Shot Learning", arxiv (2022) (Round Robin) J. Cho et al., "Unifying vision-and-language tasks via text generation", ICML (2021)

0.4

Training progress

0.2

(Vanilla XAttn) A. Vaswani et al., "Attention is all you need", NeurIPS (2017) (Grafting) Z. Luo et al., "VC-GPT: Visual Conditioned GPT for End-to-End Generative Vision-and-Language Pre-training", arxiv (2022)

## Influence of training data mixture

ngo 3B	Changed value	Param. count↓	Step time↓	COCO CIDEr↑	OKVQA top1↑	VQAv2 top1↑	ImageNet top1↑	MSVDQA top1↑	VATEX CIDEr↑	Kinetics top1-top5↑	Ove: scor
model (sh	ort training)	3.2B	1.74s	86.5	42.1	55.8	59.9	36.3	53.4	49.4	68
	M3W	3.2B	0.68s	58.0	37.2	48.6	35.7	29.5	33.6	34.0	50
to	w/o VTP	3.2B	1.42s	84.2	43.0	53.9	59.6	34.5	46.0	45.8	65.
la	w/o LTIP/ALIGN	3.2B	0.95s	66.3	39.2	51.6	41.4	32.0	41.6	38.2	56
	w/o M3W	3.2B	1.02s	54.1	36.5	52.7	24.9	31.4	23.5	28.3	46

For some dataset combinations, no <EOC> token is produced (instead additional prompts are predicted)

For these cases, the prediction is trimmed to the text preceding the prompt keywords.

## Optimisation strategy for mixing datasets

ngo 3B	Changed value	Param. count↓	Step time↓	COCO CIDEr↑	OKVQA top1↑	VQAv2 top1↑	ImageNet top1↑	MSVDQA top1↑	VATEX CIDEr↑	Kinetics top1-top5↑	Ove sco
model (short t	raining)	3.2B	1.74s	86.5	42.1	55.8	59.9	36.3	53.4	49.4	68
accumulation	Round Robin	3.2B	1.68s	76.1	39.8	52.1	50.7	33.2	40.8	39.7	59

## Architecture: tanh cross-attention gating

ngo 3B Changed	Param.	Step COCO	OKVQA	VQAv2	ImageNet	MSVDQA	VATEX	Kinetics	Over
value	count↓	time↓   CIDEr↑	top1↑	top1↑	top1↑	top1↑	CIDEr↑	top1-top5↑	scor
model (short training)	3.2B	1.74s 86.5	42.1	55.8	59.9	36.3	53.4	49.4	68.
×	3.2B	1.74s   78.4	40.5	52.9	54.0	35.9	47.5	46.4	64.

## Conditioning architectures for the frozen language model

Flamingo 3B	Changed	Param.	Step	COCO	OKVQA	VQAv2	ImageNet	MSVDQA	VATEX	Kinetics	Ove
value	value	count↓	time↓	CIDEr↑	top1↑	top1↑	top1↑	top1↑	CIDEr↑	top1-top5↑	sco
go 3B model (short training)		3.2B	1.74s	86.5	42.1	55.8	59.9	36.3	53.4	49.4	68
GATED	VANILLA XATTN	2.4B	1.16s	80.6	41.5	53.4	59.0	32.9	50.7	46.8	65
XATTN-DENSE	GRAFTING	3.3B	1.74s	79.2	36.1	50.8	47.5	32.2	47.8	27.9	57

rall	
re†	
.4	
.7	
.4	
.5	
.9	

rall	
reî	
.4	
7	
• /	

![](_page_28_Picture_16.jpeg)

![](_page_28_Picture_17.jpeg)

![](_page_28_Picture_18.jpeg)

# Ablation studies cont.

		Со
	Ablated	Flaming
	setting	value
	Flamir	ngo 3B m
v)	Cross-attention frequency	Every

	Resampler architecture and size												
	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↑	Ove sco
	Flamingo 3B model (short training)			3.2B	1.74s	86.5	42.1	55.8	59.9	36.3	53.4	49.4	68
(vi)	Resampler	Perceiver	MLP Transformer	3.2B 3.2B	1.85s 1.81s	78.6 83.2	42.2 41.7	54.7 55.6	53.6 59.0	35.2 31.5	44.7 48.3	42.1 47.4	63 65
(vii)	Resampler size	Medium	<mark>Small</mark> Large	3.1B 3.4B	1.58s 1.87s	81.1 84.4	40.4 42.2	54.1 54.4	60.2 60.4	36.0 35.1	50.2 51.4	48.9 49.4	66 67

	Number of images attended to												
	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↑	Ove scor
	Flamingo 3B model (short training)			3.2B	1.74s	86.5	42.1	55.8	59.9	36.3	53.4	49.4	68
(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

	Ablated setting	Flami value
		Flamingo 3B
(ix)	<i>p</i> <sub>next</sub>	0.5

### **References/Image credits**

J-B. Alayrac et al., "Flamingo: a Visual Language Model for Few-Shot Learning", arxiv (2022)

# mpute/capacity vs. performance trade-off for cross-attention

ngo 3B	Changed value	Param. count↓	Step time↓	COCO CIDEr↑	OKVQA top1↑	VQAv2 top1↑	ImageNet top1↑	MSVDQA top1↑	VATEX CIDEr↑	Kinetics top1-top5↑	Ove scor
model (short training)		3.2B	1.74s	86.5	42.1	55.8	59.9	36.3	53.4	49.4	68
	Single in middle	2.0B	0.87s	71.5	38.1	50.2	44.0	29.1	42.3	28.3	54
	Every 4th	2.3B	1.02s	82.3	42.7	55.1	57.1	34.6	50.8	45.5	65
	Every 2nd	2.6B	1.24s	83.7	41.0	55.8	59.6	34.5	49.7	47.4	66

# M3W image placement data augmentation

ngo 3B	Changed	Param.	Step	COCO	OKVQA	VQAv2	ImageNet	MSVDQA	VATEX	Kinetics	Over
	value	count ↓	time ↓	CIDEr↑	top1↑	top1↑	top1↑	top1↑	CIDEr↑	top1-top5↑	scor
model (sho	ort training)	3.2B	1.74s	86.5	42.1	55.8	59.9	36.3	53.4	49.4	68.
	0.0	3.2B	1.74s	85.0	41.6	55.2	60.3	36.7	50.6	49.9	67.
	1.0	3.2B	1.74s	81.3	43.3	55.6	57.8	36.8	52.7	47.8	67.

![](_page_29_Figure_13.jpeg)

![](_page_29_Figure_14.jpeg)

·all	-
e↑	
4	-
8	-
	-

all	
e↑	
4	
8	
6	

# Ablation studies cont.

	Ablated	Flami
	Flamin	ngo 3B
(x)	Vision encoder	NFNe
	Ablated	Flami
	setting	value
	Flami	1go 3B
(xi)	LM pretraining	Massi
		Fr
	Ablated	Flami
	setting	value
	Flamii	1go 3B
(xii)	Freezing Vision	$\checkmark$
(xiii)	Freezing LM	$\checkmark$
	Ablated	Flami
	setting	value
		180 3B
(xiv)	Co-train LM on MassiveText	X
Ma	ssiveText is o	adde

### **References/Image credits**

J-B. Alayrac et al., "Flamingo: a Visual Language Model for Few-Shot Learning", arxiv (2022) (CLIP) A. Radford et al., "Learning transferable visual models from natural language supervision", ICML (2021) (MassiveText) J. Rae, et al. "Scaling language models: Methods, analysis & insights from training gopher", arxiv (2021) C. Raffel, et al., "Exploring the limits of transfer learning with a unified text-to-text transformer", JMLR (2019)

### Vision encoder pretraining

ngo 3B	Changed	Param.	Step	COCO	OKVQA	VQAv2	ImageNet	MSVDQA	VATEX	Kinetics	Ove:
	value	count↓	time↓	CIDEr↑	top1↑	top1↑	top1↑	top1↑	CIDEr↑	top1-top5↑	scor
model (short t	raining)	3.2B	1.74s	86.5	42.1	55.8	59.9	36.3	53.4	49.4	68
t-F6	CLIP ViT-L/14	3.1B	1.58s	76.5	41.6	53.4	49.5	33.2	44.5	42.3	61.
	NFNet-F0	2.9B	1.45s	73.8	40.5	52.8	49.8	31.1	42.9	36.6	58

### Language model pretraining

ngo 3B	Changed value	Param. count↓	Step time↓	COCO CIDEr↑	OKVQA top1↑	VQAv2 top1↑	ImageNet top1↑	MSVDQA top1↑	VATEX CIDEr↑	Kinetics top1-top5↑	Ove: scor
model (sho	ort training)	3.2B	1.74s	86.5	42.1	55.8	59.9	36.3	53.4	49.4	68
veText	C4	3.2B	1.74s	81.3	34.4	47.1	60.6	30.9	53.9	46.9	62.

## reezing model components to prevent catastrophic forgetting

ngo 3B	Changed value	Param. count↓	Step time↓	COCO CIDEr↑	OKVQA top1↑	VQAv2 top1↑	ImageNet top1↑	MSVDQA top1↑	VATEX CIDEr↑	Kinetics top1-top5↑	Over scor
model (short t	training)	3.2B	1.74s	86.5	42.1	55.8	59.9	36.3	53.4	49.4	68
	<ul><li>(random init)</li><li>(pretrained)</li></ul>	3.2B 3.2B	4.70s* 4.70s*	74.5 83.5	41.6 40.6	52.7 55.1	45.2 55.6	31.4 34.6	35.8 50.7	32.6 41.2	56. 64.
	<ul><li>(random init)</li><li>(pretrained)</li></ul>	3.2B 3.2B	2.42s 2.42s	74.8 81.2	31.5 33.7	45.6 47.4	59.5 60.7	26.9 31.0	50.1 53.9	43.4 49.9	58. 62.

## Co-training the language model on MassiveText

ngo 3B	Changed	Param.	Step	COCO	OKVQA	VQAv2	ImageNet	MSVDQA	VATEX	Kinetics	Ove
	value	count ↓	time $\downarrow$	CIDEr↑	top1↑	top1↑	top1↑	top1↑	CIDEr↑	top1-top5↑	sco
model (short	training)	3.2B	1.74s	86.5	42.1	55.8	59.9	36.3	53.4	49.4	68
	✓ (random init)	3.2B	5.34s*	69.3	29.9	46.1	59.9	28.1	45.5	46.9	57
	✓ (pretrained)	3.2B	5.34s*	83.0	42.5	53.3	60.9	35.1	51.1	50.1	67

ed to the mixture with  $\lambda_m = 1.0$  (determined by grid search)

rall re↑	
.4	
.4 .9	

rall re↑	
.4	
.5	1

rall	
re†	
.4	
.6	
.5	
.2	
.9	

	_
rall	
re↑	
.4	
.4	
.2	

# Outline

- Motivation and challenges for multimodal generative modelling
- Related Work
- Flamingo model
- Evaluation datasets
- Comparison to state-of-the-art
- Ablations
- Qualitative Examples
- Discussion (limitations, trade-offs, opportunities, benefits, risks)

# Low-level details

To generate qualitative examples, greedy decoding is used, rather than beam search (as used for benchmarks) The examples are selected (rather than random)

![](_page_32_Picture_3.jpeg)

### **References/Image credits**

![](_page_32_Picture_7.jpeg)

			Text completion of ima	ge and text prompts		
Input Prompt	We stateWe stateStateCharacterizationWe stateState <td>We have a state of the state</td> <td>Westion: What season is this? Explain why. Answer:</td> <td>Question: What is the title of the book? Answer:</td> <td>The ingredients of this dish are</td> <td><image/><image/></td>	We have a state of the state	Westion: What season is this? Explain why. Answer:	Question: What is the title of the book? Answer:	The ingredients of this dish are	<image/> <image/>
Prompletion	Vitamin C, potassium, and fiber.	To keep the mosquitoes out of the bed.	This is winter season. The Christmas tree and the presents are the main reasons why it is winter.	The House Book.	: beef, potatoes, carrots, broccoli, and lemon.	the songs "Let It Go" and the First Time in Forever Tony Award® winner Ic Menzel, who plays Elsa

### References/Image credits

![](_page_33_Figure_5.jpeg)

## Text completion of image and text prompts

		Inpu	t Prompt
	This is a chinchilla. They are mainly found in Chile.		This is a shiba are very popu Japan.
	What is the title of this painting? Answer: The Hallucinogenic Toreador.		Where is this p displayed? A Louvres Museu
UNDERGROUND	Output: "Underground"	CONGRESS 400	Output: "Con
2+1	2+1=3	5+6	5+6=1

### References/Image credits

![](_page_34_Figure_5.jpeg)

## Text completion of image and text prompts

![](_page_35_Figure_2.jpeg)

### References/Image credits

# Qualitative results - dialogue

### Prompt

This is a conversation between a human, User, and an intelligent visual AI, Flamingo. User sends images, and Flamingo describes them.

User: <a cat image>

Flamingo: That is a cat. It's a tiny kitten with really cute big ears.

User: <a dinner image>

Flamingo: This is a picture of a group of people having dinner. They are having a great time!

User: Can you guess what are they celebrating?

Flamingo: They might be celebrating the end of a successful project or maybe a birthday?

User: <a graph image>

Flamingo: This is a graph, it looks like a cumulative density function graph.

![](_page_36_Picture_11.jpeg)

### **References/Image credits**

J-B. Alayrac et al., "Flamingo: a Visual Language Model for Few-Shot Learning", arxiv (2022) https://www.rome2rio.com/map/Grand-Central-42nd-Street-Station/590-Madison-Avenue

## Selected dialogue samples

![](_page_36_Figure_16.jpeg)

	Radison Are East 42nd S	6
	Which city is this?	$\odot$
D	This is a street sign in New York City.	
V	Which street corner is it?	$\odot$
Д	It is the intersection of Madison Avenue and East 42nd Street.	
V n	Thich subway station is earby?	•
Д	The closest subway sta- tion is Grand Central Sta- tion.	
H p	low about the nearest ark?	<b>I</b>
D	The nearest park is Bryant Park.	
T W	here is another sign. Vhat does it say?	<b>•</b>
D	It says "One way".	

I would like you to read the color and then say the color it is written in. For example:	6
PURPLE	•
Color is "Purple" and it is written in red.	6
GREEN	©
Color is "Green" and it is written in blue.	©
YELLOW	©
Color is "Yellow" and it is written in green.	
BLACK	6
Color is "Black" and it is written in yellow.	
Well done! Do you know the name of the test these images come from?	6
I think it is called the Stroop test.	]
Can you explain how hu- mans perform in this test?	©
Humans are slower when the color of the word and the color of the word are different.	
How about you?	6
I am not affected by this difference.	

![](_page_36_Picture_19.jpeg)

# Qualitative results - dialogue

![](_page_37_Figure_2.jpeg)

### **References/Image credits**

# **Qualitative results - video**

![](_page_38_Picture_2.jpeg)

### **References/Image credits**

J-B. Alayrac et al., "Flamingo: a Visual Language Model for Few-Shot Learning", arxiv (2022)

# selected video samples

![](_page_38_Picture_6.jpeg)

# **Qualitative results - more videos**

![](_page_39_Figure_2.jpeg)

### **References/Image credits**

J-B. Alayrac et al., "Flamingo: a Visual Language Model for Few-Shot Learning", arxiv (2022)

# selected video samples

![](_page_39_Picture_6.jpeg)

# Outline

- Motivation and challenges for multimodal generative modelling
- Related Work
- Flamingo model
- Evaluation datasets
- Comparison to state-of-the-art
- Ablations
- Qualitative Examples
- Discussion (limitations, trade-offs, opportunities, benefits, risks)

# Flamingo limitations

## Classification

Flamingo models lag contrastive models for classification Possibly contrastive training optimises for retrieval classification can be viewed as special case of retrieval (CLIP) By contrast, language models do have this objective alignment The work of Zhao et al. showed that language models are sensitive to prompt sample selection and their ordering It is possible to calibrate to minimise these effects, but this requires assumptions on the label space (restrictive) **Future work:** bridging the performance gap

![](_page_41_Figure_3.jpeg)

### References

(CLIP) A. Radford et al., "Learning transferable visual models from natural language supervision", ICML (2021) Z. Zhao et al., "Calibrate before use: Improving few-shot performance of language models", ICML (2021) T. Wang et al., "What Language Model Architecture and Pretraining Objective Work Best for Zero-Shot Generalization?", arxiv (2022) T. Brown et al., "Language models are few-shot learners", NeurIPS (2020)

### Legacies of language models

Flamingo builds on pretrained language models, inheriting their weaknesses Causal modelling is strictly less expressive than bidirectional modelling Recent work by Wang et al. (2022) suggests non-causal masked language modelling with multi-task fine-tuning may be a better strategy for zero-shot generalisation Challenge: if the expected output text is long, it is difficult to leverage enough shots E.g. for VisDial, 1 shot is an image with 21 sentences, so 32 shots = 672 sentences Results in 4096+ tokens - longer than the max training sequence length (2048) This may explain why performance drops with more shots (16 vs 32) on VisDial Language modelling suffers from poor sample efficiency (Brown et al., 2020) Language model priors may also cause hallucinations and ungrounded guesses.

# Flamingo failures: hallucinations/ungrounded guesses

![](_page_42_Figure_1.jpeg)

![](_page_42_Picture_5.jpeg)

# Trade-offs of few-shot learning methods

## In-context learning: advantages

In-context learning has many advantages over fine-tuning:

- requires (almost) no hyperparameter tuning
- works reasonably in low-data regime (dozens of examples)
- only requires inference (simpler deployment)

By contrast, fine-tuning methods require:

- carefully tuned design choices (learning rates, architecture)
- more data (e.g. thousands of examples) to work well

Advantages motivate choice of in-context learning for Flamingo

### Takeaway

No "golden" few-shot method: best choice depends on the scenario

(k-shots Prompt

**Inference compute cost:** in-context learning cost scales • linearly with the number of shots if the few-shot prompt can be re-used by caching • quadratically with the number of shots if no such caching is possible By contrast: gradient-based few-shot learning has constant complexity w.r.t. shots **Prompt sensitivity:** in-context learning is sensitive to prompt order and format (Zhao et al., 2021) Leveraging more shots: In-context learning performance plateaus as shots increase (e.g. >32) By contrast: gradient-based approaches tend to continue to benefit from more examples RICES helps to some extent, but there are still issues for larger numbers of examples per class Task location: several works (Min et al., 2022; Reynolds and McDonell, 2021) suggest in-context learning may not be learning a task, but instead identifying the task to be performed A few examples may therefore help task location, but its possible the model can't do anything more than task location from them, and hence cannot scale up to usefully use more examples

### References

Z. Zhao et al., "Calibrate before use: Improving few-shot performance of language models", ICML (2021)

S. Min et al., "Rethinking the Role of Demonstrations: What Makes In-Context Learning Work?", arxiv (2022)

L. Reynolds and K. McDonell, "Prompt programming for large language models: Beyond the few-shot paradigm", CHI (2021)

## In-context learning: disadvantages

![](_page_43_Picture_30.jpeg)

![](_page_43_Picture_31.jpeg)

# Flamingo opportunities

## Extending the visual and text interface

Natural language provides a versatile interface to:

- provide descriptions of visual tasks
- generate model outputs
- estimate conditional likelihoods over candidate outputs

However, it is cumbersome for structured prediction, ill suited for:

- conditioning on/predicting structures like bounding boxes
- dense predictions (over space or time)
- continuous predictions (like optical flow)

Further modalities (like audio) could extend the interface

Flamingo scales up to 80B parameters and provides some insights about scaling behaviour Scaling laws were studied for language (Kaplan et al., 2020) and vision (Zhai et al., 2021) There is limited work understanding scaling for vision-language models Rather than focusing on perplexity, downstream task performance may be a better metric

### **References/Image credits**

J. Kaplan et al., "Scaling laws for neural language models", arxiv (2020)

X. Zhai et al., "Scaling vision transformers", arxiv (2021)

## Scaling laws for vision-language models

![](_page_44_Picture_17.jpeg)

# Flamingo benefits

# Accessibility

Flamingo can be trained with minimal examples and used through

a chat-like interface for open-ended dialogue

This could enable non-expert users to apply Flamingo in low-

resource settings

Example: Flamingo works well on VizWiz

Dialogue interface could help highlight issues with bias/toxicity

## Model recycling

Although costly to train, Flamingo demonstrates how to leverage frozen pretrained vision encoders and language models This suggests new modalities can be introduced into frozen models Could help with reducing environmental impact (Strubell et al., 2019)

# Flamingo risks and mitigation strategies

## Inherited risks of large language models

Flamingo relies heavily on a pretrained language model With no input images, it defaults to language model behaviour Result: offensive language, stereotypes, private info leakage Flamingo is based on Chinchilla (Hoffmann et al., 2020) slightly less gendered biased than prior models, but still biased Chinchilla also has relatively low toxicity, as measured using the PerspectiveAPI toxicity score on 25,000 samples Potential mitigations: social/public policy interventions (regulation and guidelines), research on AI Ethics/NLP, better benchmarks

### Toxicity when prompted with images

Some Flamingo captions were tagged as toxic by PerspectiveAPI However, on manual inspection, no clear toxicity was found Toxic outputs not observed for with "safe-for-work" imagery

### **References/Image credits**

(Chinchilla) J. Hoffmann et al., "Training Compute-Optimal Large Language Models", arxiv (2022) D. Zhao et al., "Understanding and Evaluating Racial Biases in Image Captioning", ICCV (2021) J-B. Alayrac et al., "Flamingo: a Visual Language Model for Few-Shot Learning", arxiv (2022)

## Gender and racial biases

Conduct a study to assess captioning bias following Zhao et al. (2021)

Evaluate how performance varies on COCO as a function of gender and skin colour

	CIDEr difference				
	female - male = $\Delta$	darker - lighter = $\Delta$	overall		
Flamingo, 0 shot	$0.899 - 0.870 = +0.029 \ (p = 0.52)$	$0.955 - 0.864 = +0.091 \ (p = 0.25)$	0.843		
Flamingo, 32 shots	$1.172 - 1.142 = +0.030 \ (p = 0.54)$	$1.128 - 1.152 = -0.025 \ (p = 0.76)$	1.138		

No statistically significant differences were observed.

# Flamingo for mitigation

Flamingo could be used for filtering purposes for toxic samples in the training data

During evaluation, models adapted on filtered data could be used to down-rank/

exclude outputs that do not meet desired standards

Flamingo performance on HatefulMemes suggesting it may be well-suited for this task

Could be used for "red-teaming" to identify issues in other models (Perez et al., 2022)

As shown in the qualitative examples, Flamingo can, in cases, explain its own outputs

![](_page_46_Picture_20.jpeg)

# Summary

images and video data with limited training data State-of-the-art results a variety of tasks Qualitative examples demonstrating interactive abilities In summary, a highly flexible vision and language model

## Summary

- Flamingo is a "general-purpose" family of models that can be applied to