

Flamingo: a Visual Language Model for Few-Shot Learning

Кириллов Даниил 213

План

- 1) Коротко о Flamingo
- 2) Немного базовых определений
- 3) Flamingo - архитектура, обучение
- 4) Результаты

Flamingo (DeepMind, 2022)

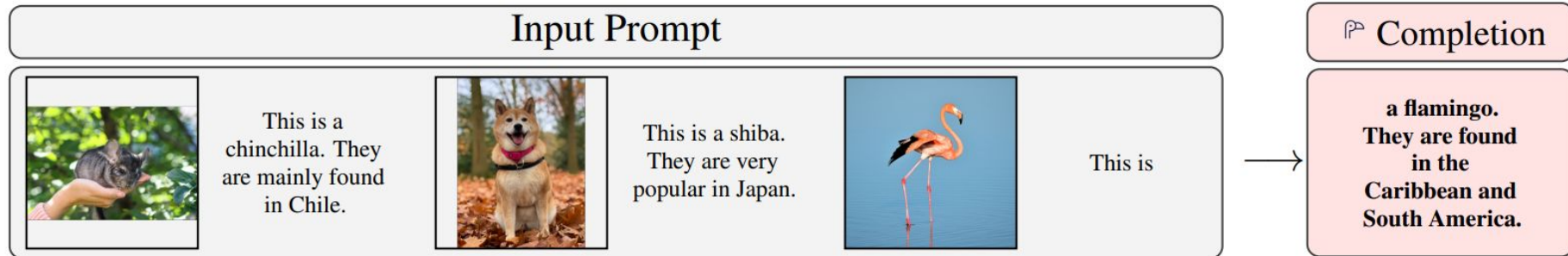
- мультимодальная сеть, получает на вход изображения и текст в произвольном порядке, выдает текстовый ответ на запрос
- few-shot learning
- Способна к обучению на мультимодальных корпусах данных с произвольным содержанием текстов, изображений, видео.
- SOTA на нескольких бенчмарках, на некоторых бенчмарках превзошла предыдущее SOTA с finetune-моделей.



Немного базовых определений

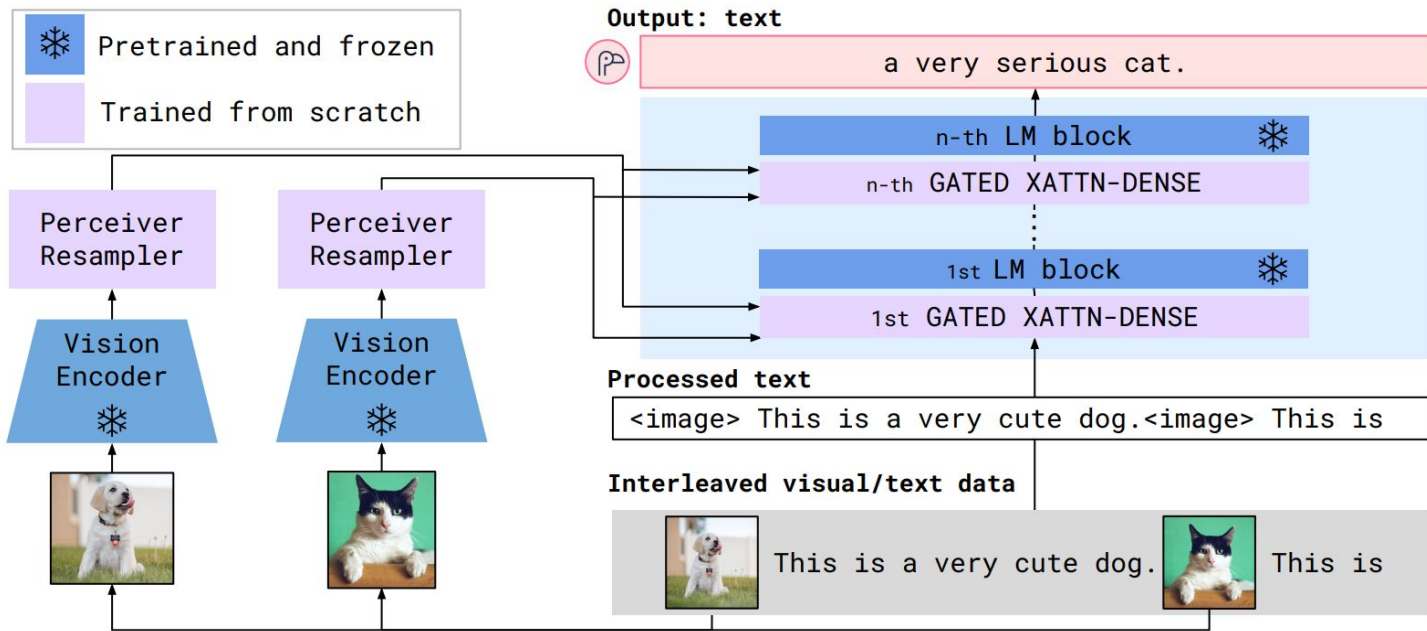
Мультимодальная сеть - сеть, способная работать с разными типами входных данных (текст, видео, картинки, ...)

Few-shot - способность модели дообучиться при подаче нескольких примеров работы.



Flamingo - архитектура

Содержит две предобученные модели - для извлечения признаков из изображений и языковой.

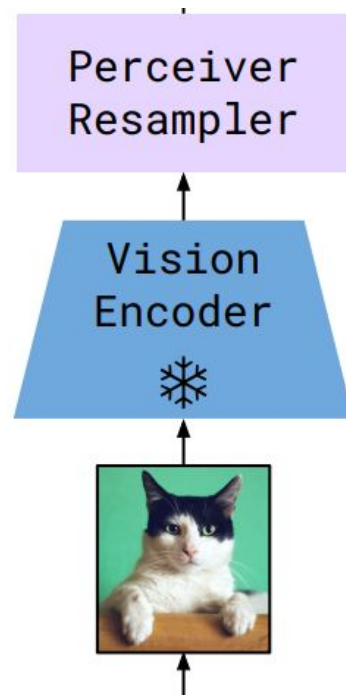


Flamingo - Visual features

Vision Encoder - получение признаков из изображения.

Normalized-Free ResNet (NFNet), учится по парам изображение-текст на contrastive loss (как CLIP).

В случае работы с видео разбиваем поток на кадры по 1 в секунду (1 fps).



Flamingo - Visual features

Perceiver Resampler -

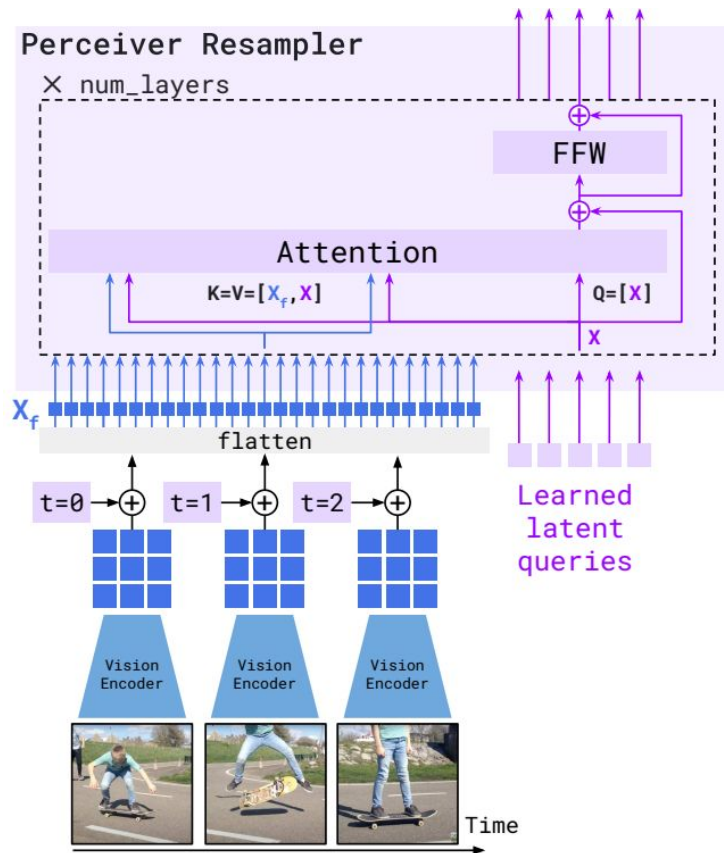
приведение полученным признаков

к вектору размера 64.

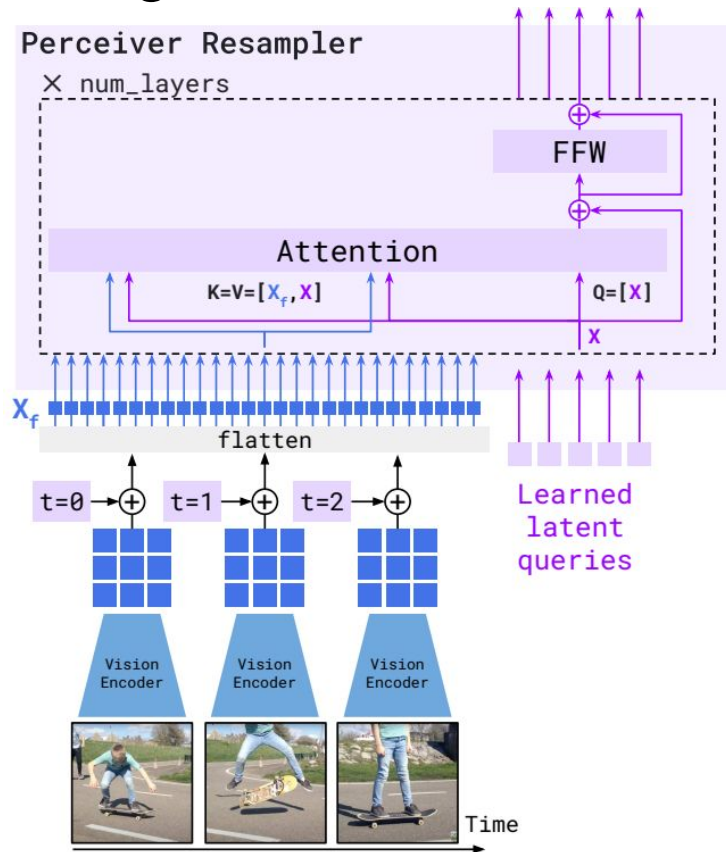
queries на первом слое - обучаемые параметры

keys, values получаются из конкатенации

признаковы изображения и queries.



Flamingo - Visual features



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

    # Add the time position embeddings and flatten.
    x_f = x_f + time_embeddings
    x_f = flatten(x_f) # [T, S, d] -> [T * S, d]

    # Apply the Perceiver Resampler layers.
    for i in range(num_layers):
        # Attention.
        x = x + attention_i(q=x, kv=concat([x_f, x]))

        # Feed forward.
        x = x + ffw_i(x)

    return x
```

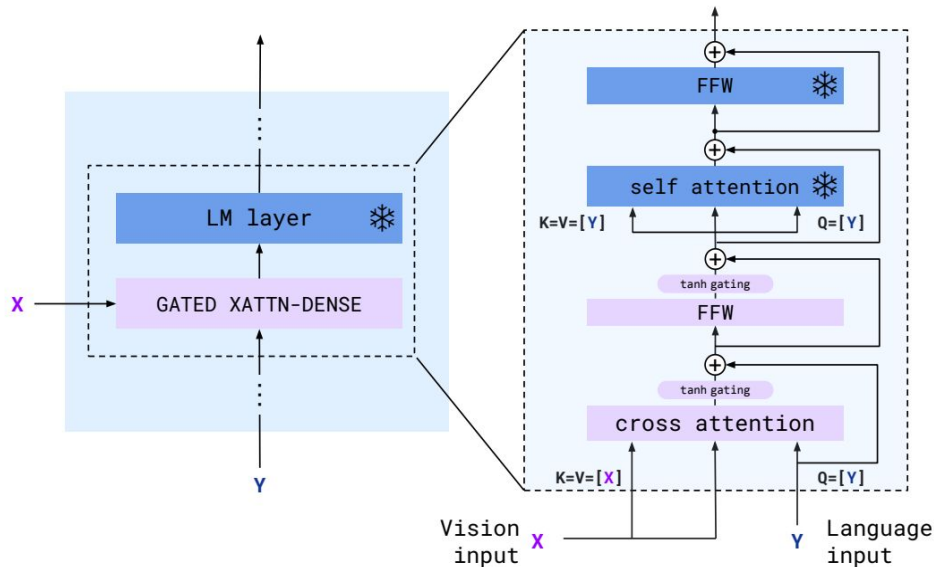

Flamingo - Gated XATTN-Dense

Учимся на правдоподобие:

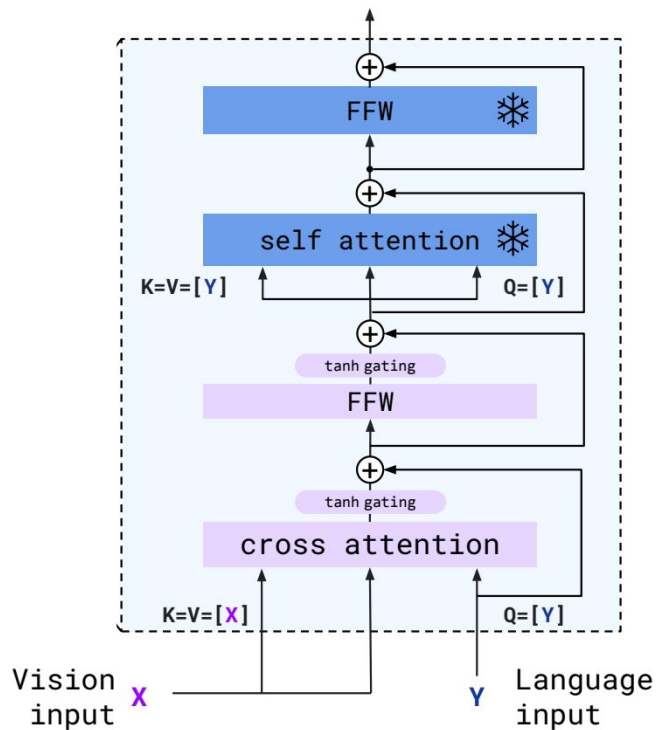
у - текстовые токены

х - токены изображений.

Пытаемся передать
признаки об изображении
в языковую модель.



Flamingo - Gated XATTN-Dense



```
def gated_xattn_dense(
    y, # input language features
    x, # input visual features
    alpha_xattn, # xattn gating parameter - init at 0.
    alpha_dense, # ffw gating parameter - init at 0.
):
    """Applies a GATED XATTN-DENSE layer."""

    # 1. Gated Cross Attention
    y = y + tanh(alpha_xattn) * attention(q=y, kv=x)

    # 2. Gated Feed Forward (dense) Layer
    y = y + tanh(alpha_dense) * ffw(y)

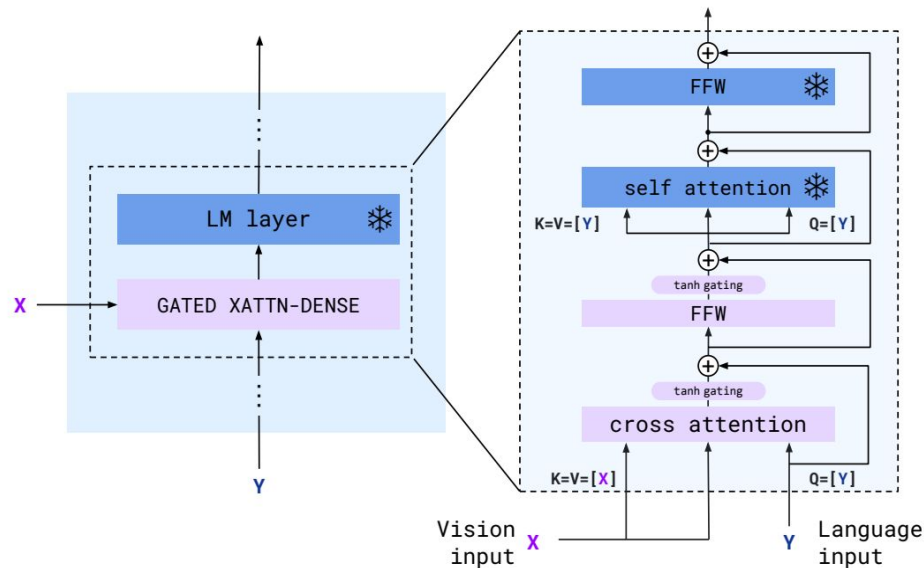
    # Regular self-attention + FFW on language
    y = y + frozen_attention(q=y, kv=y)
    y = y + frozen_ffw(y)

    return y # output visually informed language features
```

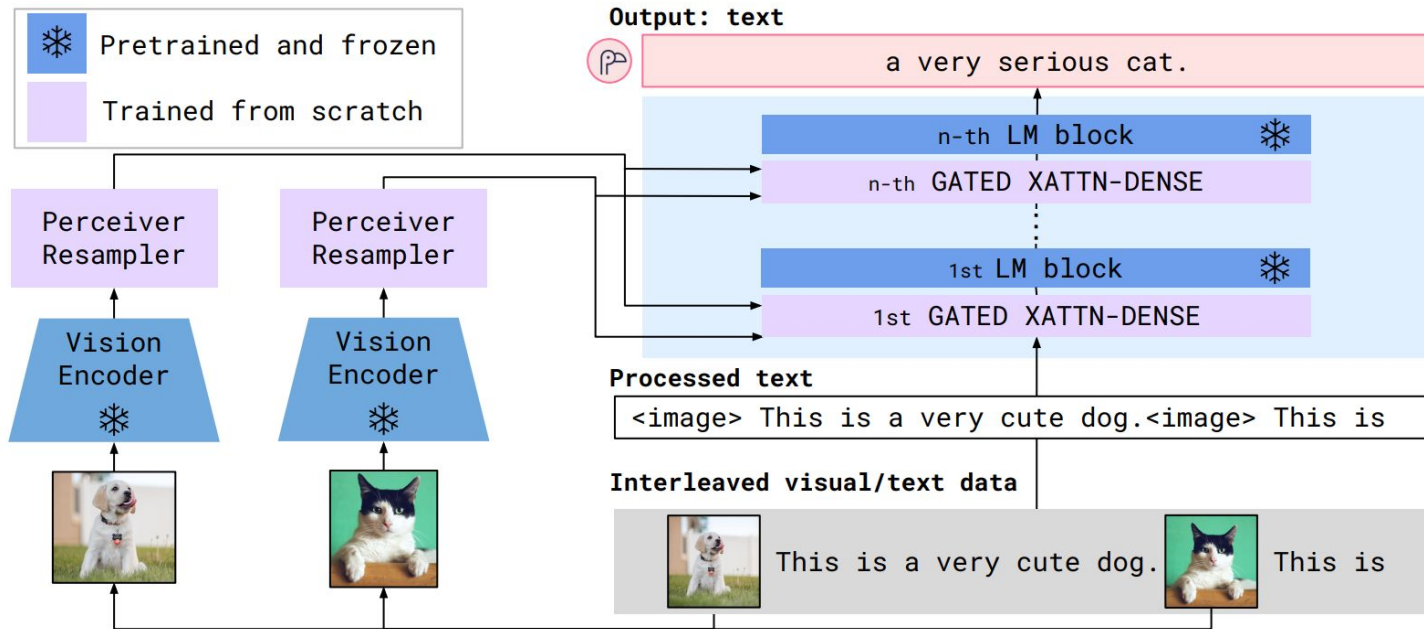
Flamingo - Gated XATTN-Dense

cross-attention маскируется
чтобы смотреть только
на последний visual token.

в функции активации
 α - обучаемый параметр,
изначально равный 0.



Flamingo



Flamingo - обучение

Авторы собрали свой датасет - MultiModal MassiveWeb (M3W), собранный из интернета - со страницы берутся 256 подряд идущих случайных токенов и до 5-ти первых изображений, попавших в этот отрезок.

Из датасета ALIGN (картинка + описание), взяли часть картинок с длинным описанием - датасет LTIP (Long Text and Image Pairs) и собрали аналогичный датасет с видео - VTP (Video and Text Pairs).

Flamingo - обучение

Multi-objective training and optimisation.

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

Считается loss по всем датасетам и суммируется с некоторыми весами.

Flamingo - результаты

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

Flamingo - результаты

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)
<i>Flamingo-3B</i>	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	-
<i>Flamingo-9B</i>	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	-
<i>Flamingo</i>	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	✓		54.4	80.2	143.3	47.9	76.3	57.2	67.4	46.8	35.4	138.7	36.7	75.2	54.7	25.2	79.1	
		(X)	[34] (10K)	[140] (444K)	[124] (500K)	[28] (27K)	[153] (500K)	[65] (20K)	[150] (30K)	[51] (130K)	[135] (6K)	[132] (10K)	[128] (46K)	[79] (123K)	[137] (20K)	[129] (38K)	[62] (9K)	-

Источники

Оригинальная статья про архитектуру - <https://arxiv.org/abs/2204.14198>