



Лаборатория компьютерной  
графики и мультимедиа  
ВМК МГУ имени М.В. Ломоносова

*Курс «Компьютерное зрение»*

# **«Само-обучение и фундаментальные модели»**

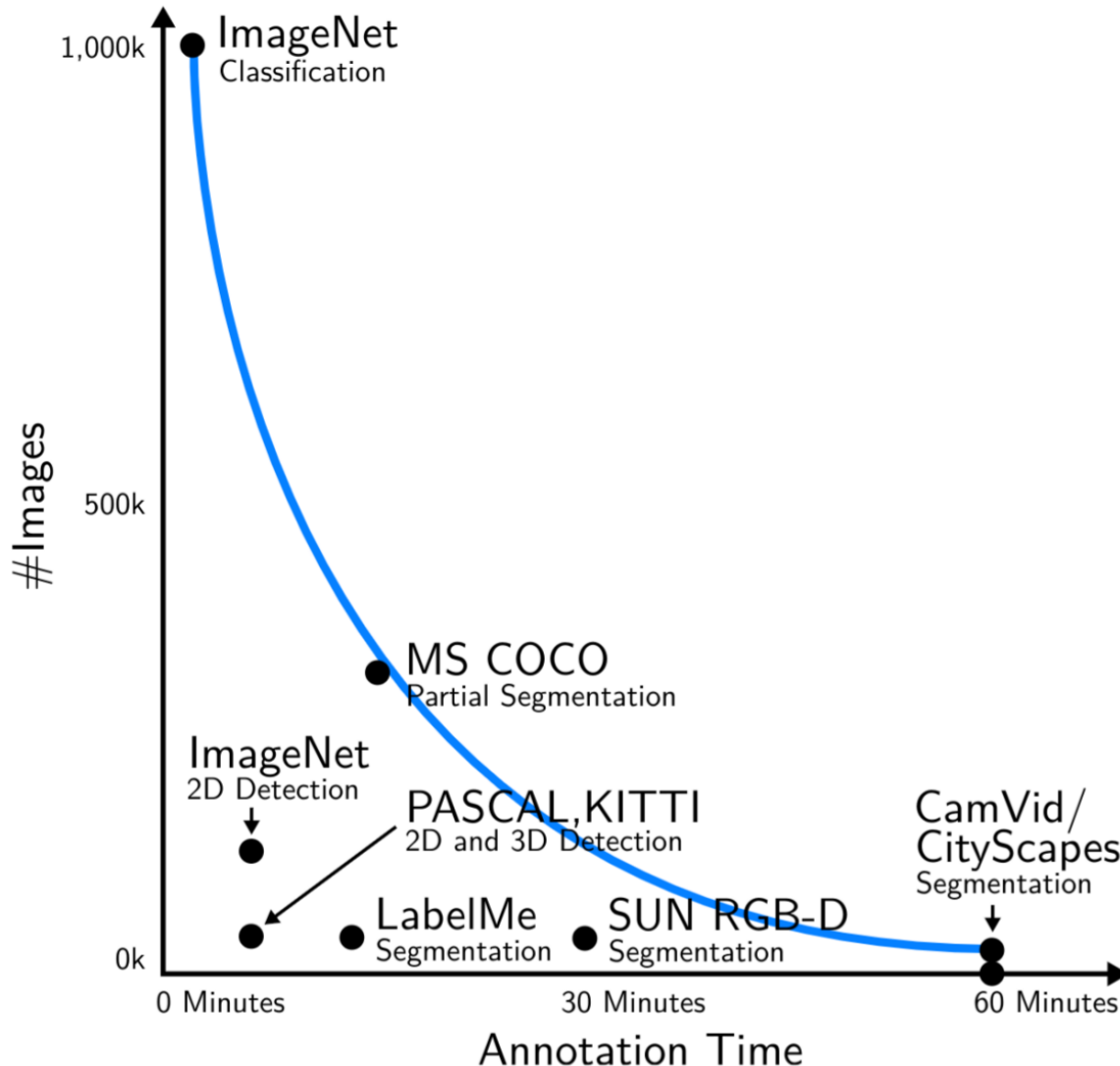
Антон Конушин и Тимур Мамедов

2025 год



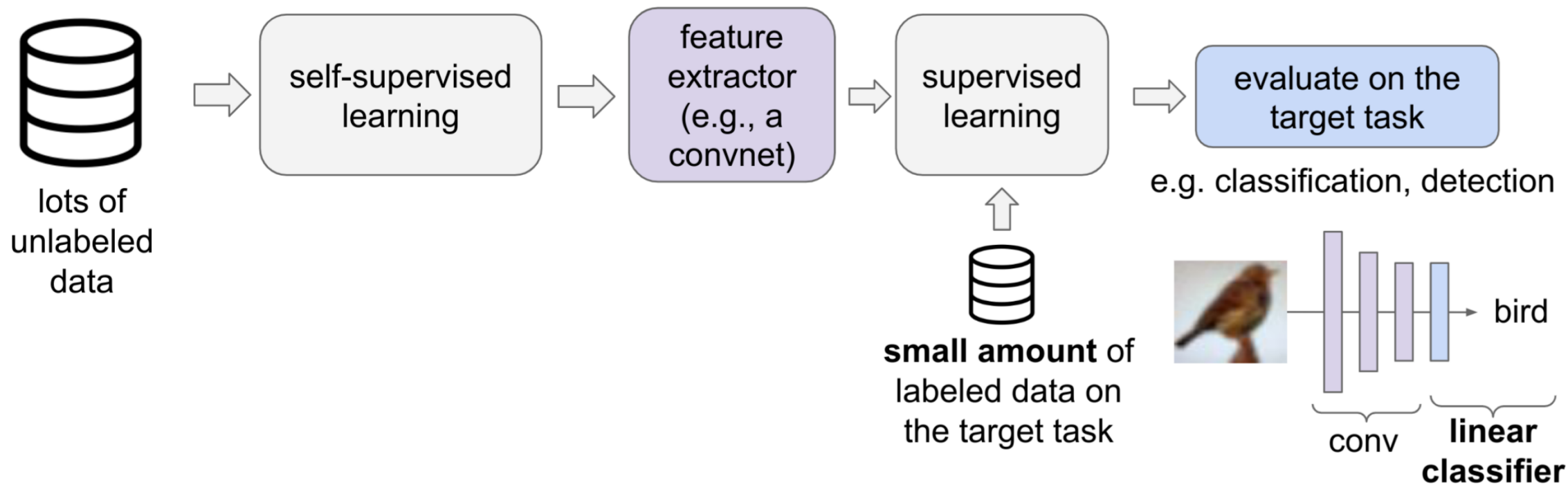
1. Введение
2. Прокси-задачи
3. Констранстное обучение и маскирование
4. Фундаментальные модели

# Дополнительные plug-in модули



- Разметка данных дорогая
- Метки всегда содержат ошибки
- Обучать людей размечать эффективно и без ошибок очень сложно
- Существует много неразмеченных данных, м.б. мы можем как-то их использовать?

# Предобучение через само-обучение



1. Обучить полезный нейро-признаки (или представления / representation) с помощью прокси-задачи
2. Добавить поверх нейропризнаков сколько-то слоёв и до-обучить на размеченных датасетах





1. Введение
2. Прокси-задачи
3. Констранстное обучение и маскирование
4. Фундаментальные модели

# Предсказание контекста



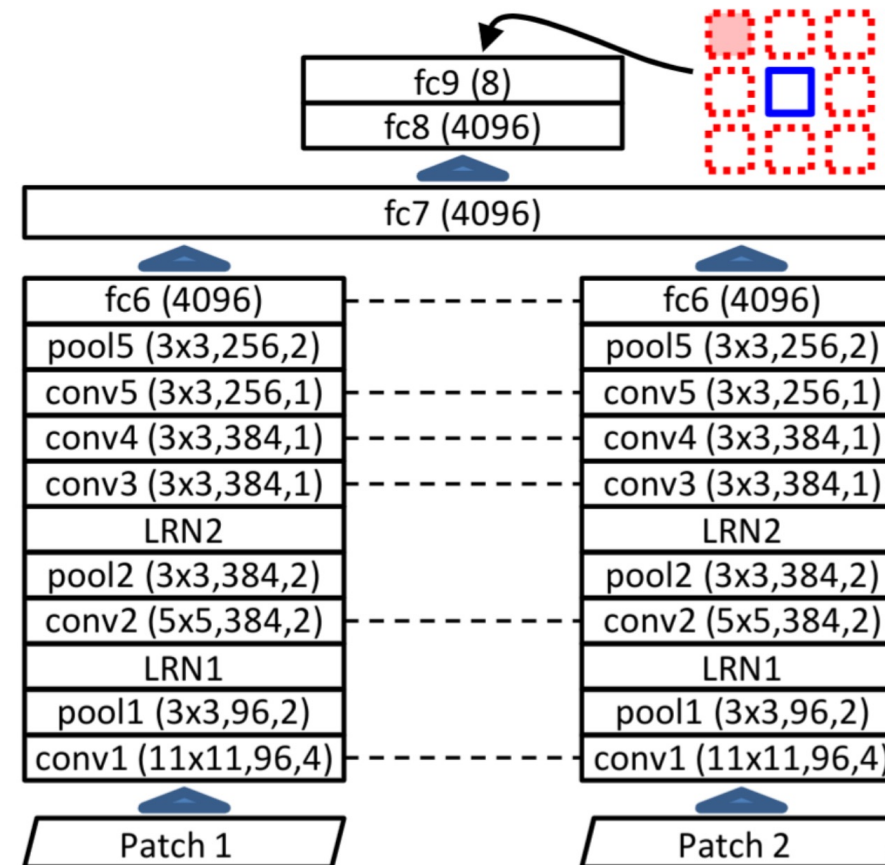
Example:



Question 1:

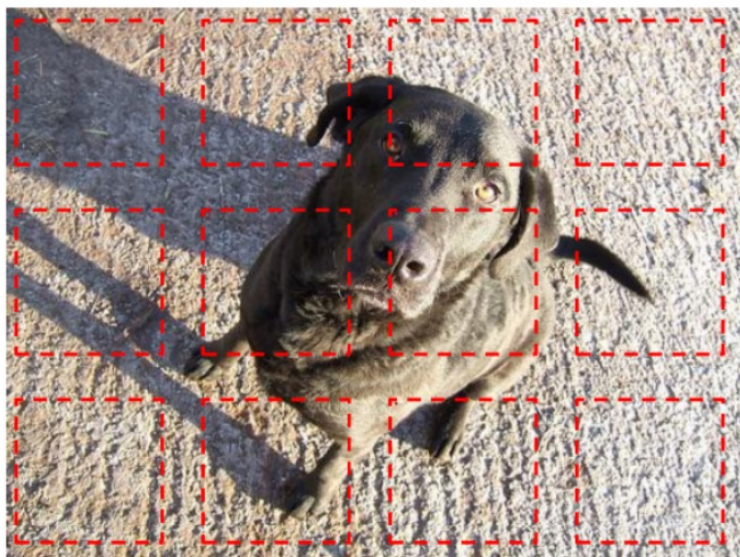


Question 2:



Предсказываем, где расположен фрагмент относительно центрального, как задача 8-и классовой классификации

# Предсказание контекста – забавный эффект



Initial layout, with sampled patches in red

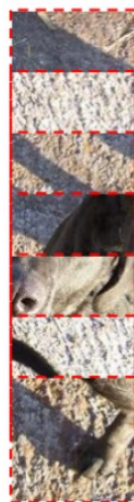
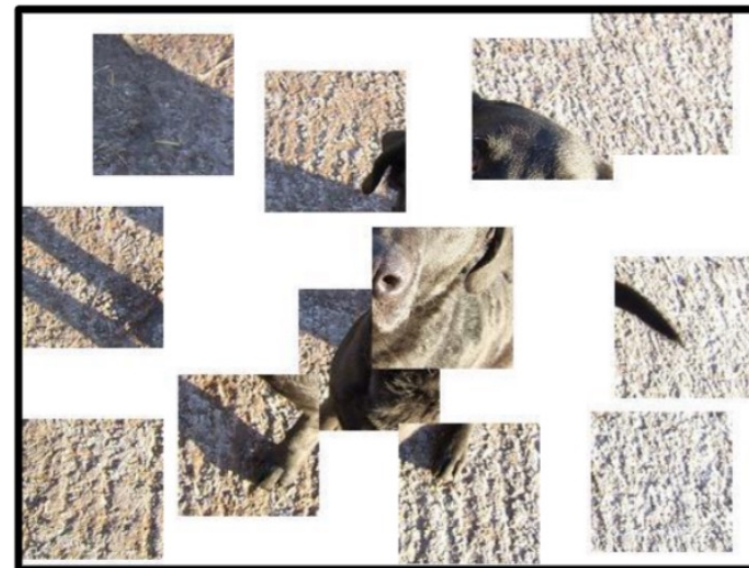
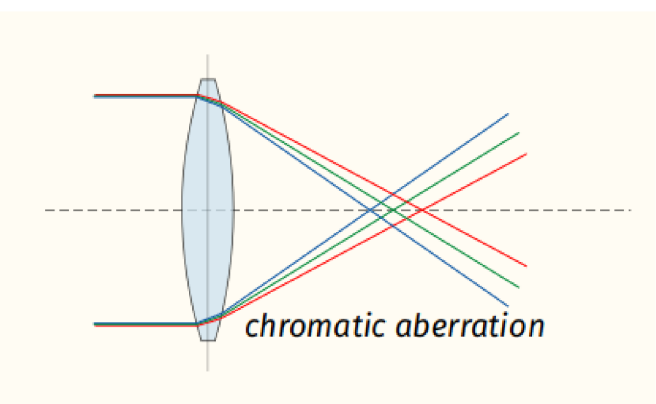


Image layout is discarded



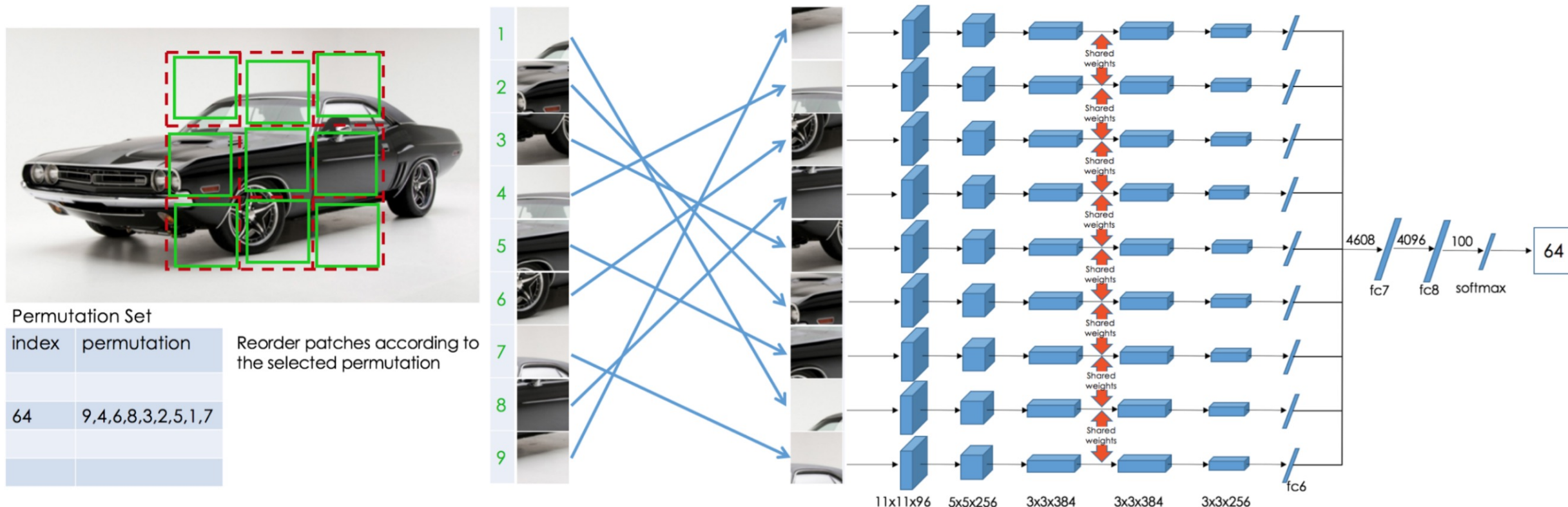
We can recover image layout automatically



Сеть может научиться «жульничать», предсказывая положения патча по эффекту хроматической абберации

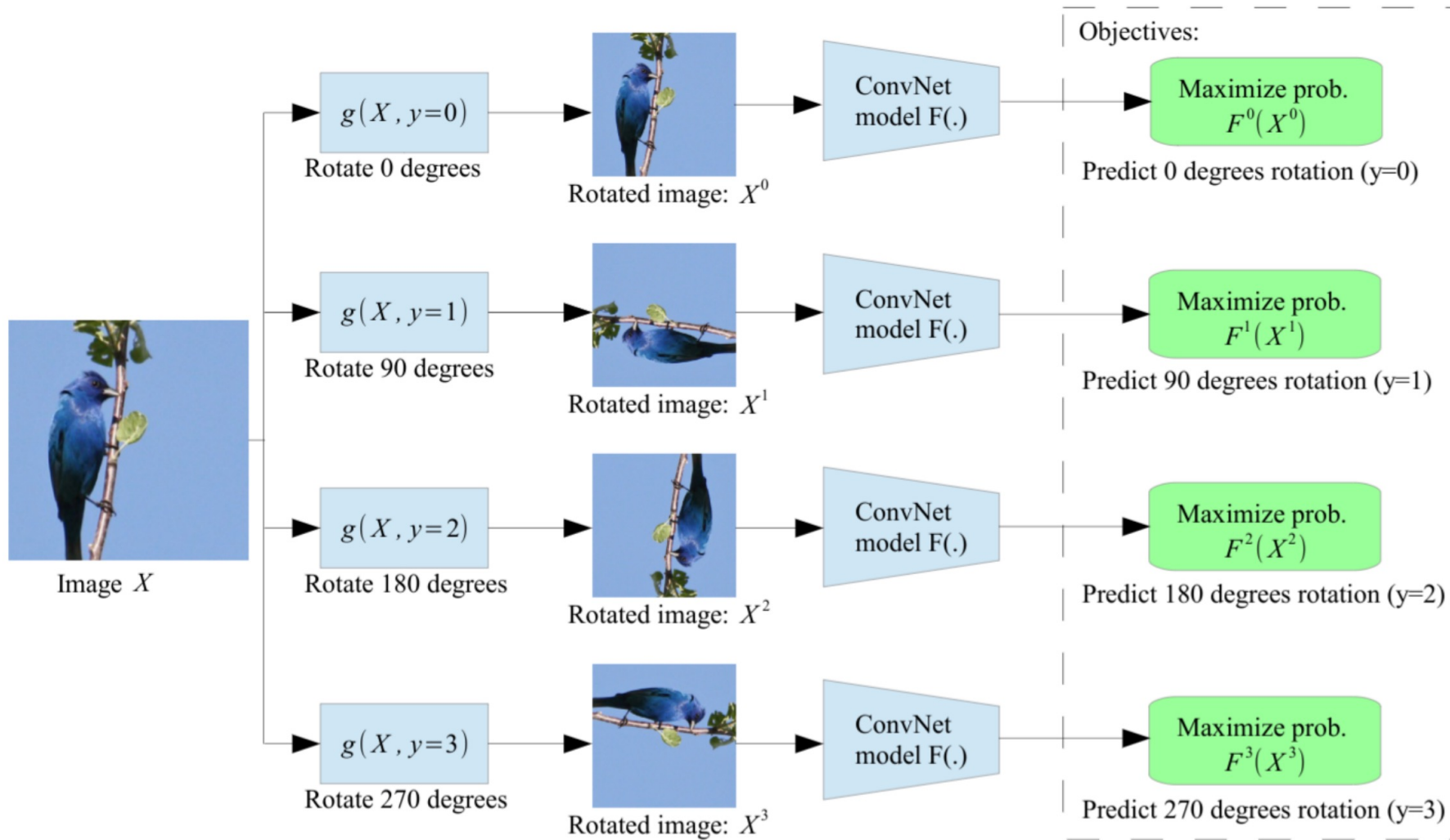


# Jigsaw puzzle



Предсказываем «перестановку» патчей – как их нужно переставить, чтобы разместить фрагменты в правильном порядке (~1000 классов)

# Повороты



Предсказываем  
поворот фрагмента  
относительно  
базового

# Повороты



Input images on the models



(a) Attention maps of supervised model

(b) Attention maps of our self-supervised model

Модель учится фокусироваться  
на ключевых фрагментах  
изображений



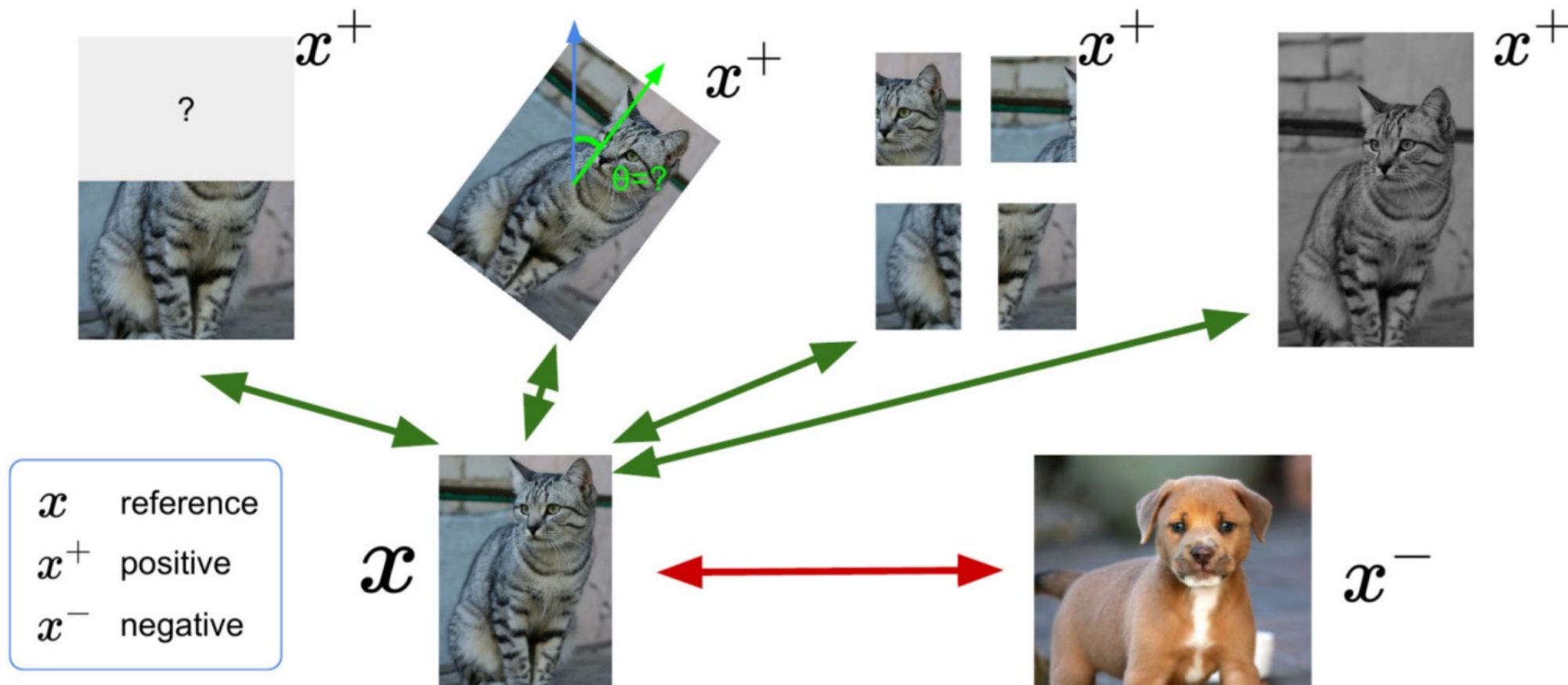
- Прокси-задачи строятся эмпирически, опираясь на какую-то априорную визуальную информацию, например, предсказание поворотов, положения в пространстве, цветов изображений
- Модели учатся понимать суть изображения для решения прокси-задачи
- Полезность обученных признаков оценивается на целевой задаче
- Придумать прокси-задачу сложно
- Сложно предсказать, насколько полезны обученные признаки и насколько модель хорошо обобщается



1. Введение
2. Прокси-задачи
3. Констранстное обучение и маскирование
4. Фундаментальные модели



# Контрастное обучение (Contrastive learning)



Given a score function  $s(\cdot, \cdot)$ , we want to learn a mapping  $f_\theta$  that yields high score for positive pairs  $(x, x^+)$  and low score for negative pairs  $(x, x^-)$ :

$$s(f_\theta(x), f_\theta(x^+)) \gg s(f_\theta(x), f_\theta(x^-))$$

# Идея контрастного обучения



Пусть у нас 1 референсное изображений ( $x$ ), 1 положительный пример ( $x^+$ ) и  $N - 1$  отрицательных примеров ( $x^-$ ). Будем использовать многоклассовую кросс-энтропию:

$$\mathcal{L} = -\mathbf{E}_x \left[ \log \frac{\exp(s(f_\theta(x), f_\theta(x^+)))}{\exp(s(f_\theta(x), f_\theta(x^+))) + \sum_{j=1}^{N-1} \exp(s(f_\theta(x), f_\theta(x_j^-)))} \right]$$

Функция известна как InfoNCE loss и её отрицание дает нижнюю границу на mutual information между  $f_\theta(x)$  и  $f_\theta(x^+)$ :

$$MI[f_\theta(x), f_\theta(x^+)] \geq \log N - \mathcal{L}$$

Максимизация совместной информации между разными «видами» на изображение позволяет извлечь высокоуровневую информацию из изображений

# SimCLR

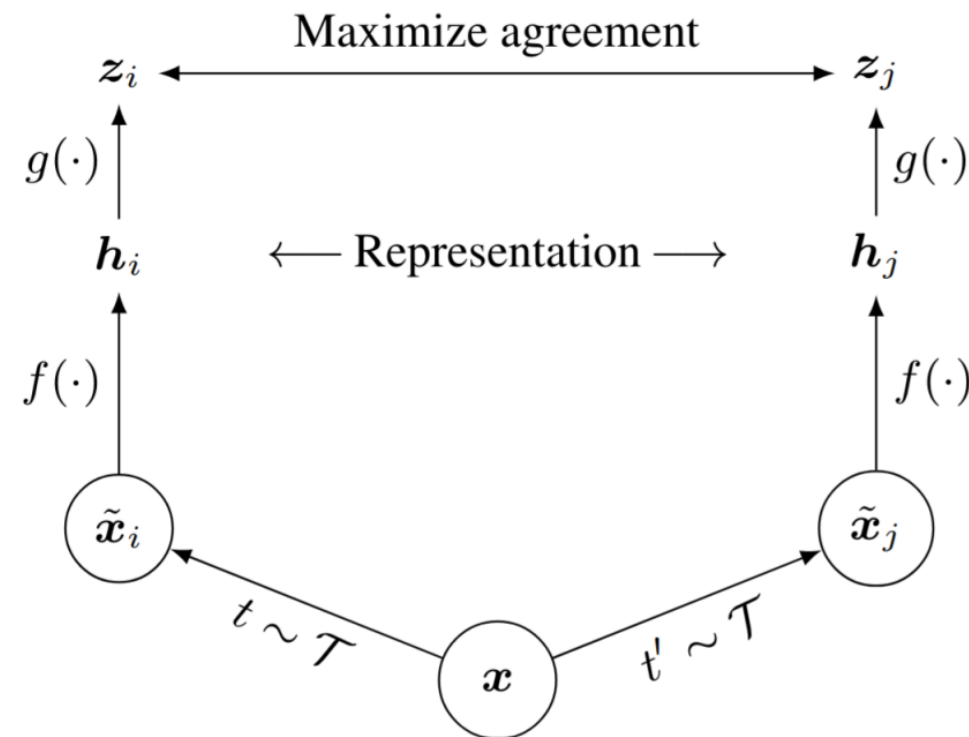


Функция косинусного расстояния как целевая функция:

$$s(\mathbf{z}_i, \mathbf{z}_j) = \frac{\mathbf{z}_i^T \mathbf{z}_j}{\|\mathbf{z}_i\| \|\mathbf{z}_j\|}$$

Используем вспомогательную сеть-проектор  $g(\cdot)$  для проецирования признаков в пространство, в котором применяем контрастное обучение

Такой проектор улучшает обучение, т.к. более важная информация сохраняется в  $h$

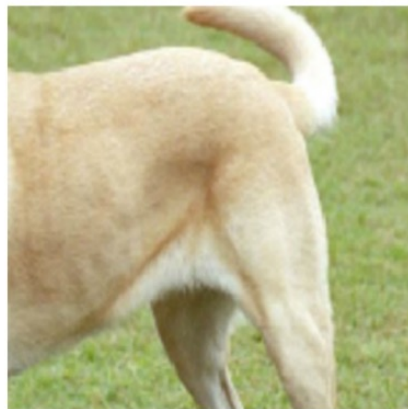




# SimCLR



(a) Original



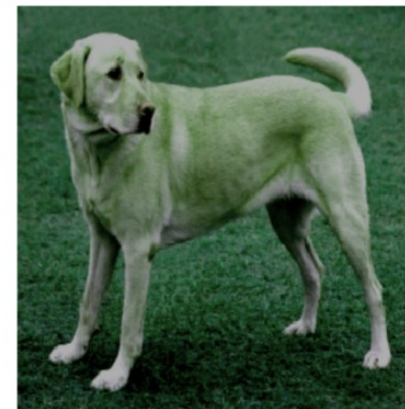
(b) Crop and resize



(c) Crop, resize (and flip)



(d) Color distort. (drop)



(e) Color distort. (jitter)



(f) Rotate  $\{90^\circ, 180^\circ, 270^\circ\}$



(g) Cutout



(h) Gaussian noise



(i) Gaussian blur



(j) Sobel filtering

Примеры аугментаций изображений при получении пар изображений

```

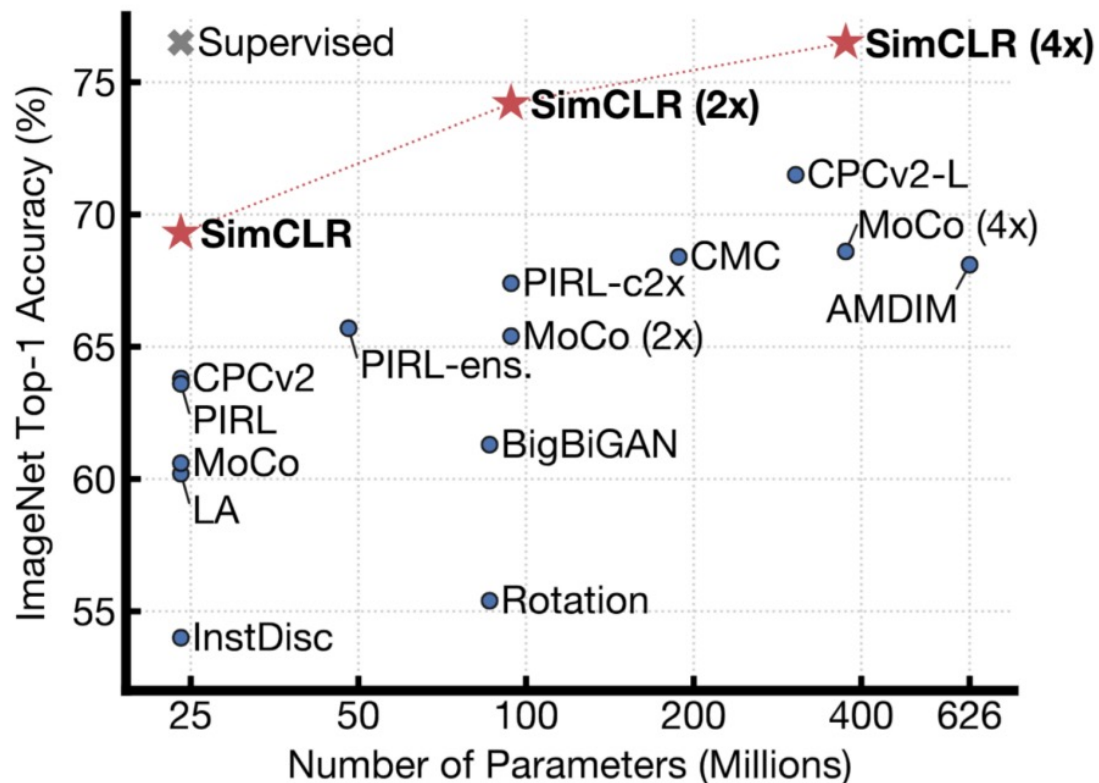
input: batch size  $N$ , constant  $\tau$ , structure of  $f, g, \mathcal{T}$ .
for sampled minibatch  $\{\mathbf{x}_k\}_{k=1}^N$  do
  for all  $k \in \{1, \dots, N\}$  do
    draw two augmentation functions  $t \sim \mathcal{T}, t' \sim \mathcal{T}$ 
    # the first augmentation
     $\tilde{\mathbf{x}}_{2k-1} = t(\mathbf{x}_k)$ 
     $\mathbf{h}_{2k-1} = f(\tilde{\mathbf{x}}_{2k-1})$  # representation
     $\mathbf{z}_{2k-1} = g(\mathbf{h}_{2k-1})$  # projection
    # the second augmentation
     $\tilde{\mathbf{x}}_{2k} = t'(\mathbf{x}_k)$ 
     $\mathbf{h}_{2k} = f(\tilde{\mathbf{x}}_{2k})$  # representation
     $\mathbf{z}_{2k} = g(\mathbf{h}_{2k})$  # projection
  end for
  for all  $i \in \{1, \dots, 2N\}$  and  $j \in \{1, \dots, 2N\}$  do
     $s_{i,j} = \mathbf{z}_i^\top \mathbf{z}_j / (\|\mathbf{z}_i\| \|\mathbf{z}_j\|)$  # pairwise similarity
  end for
  define  $\ell(i, j)$  as  $\ell(i, j) = -\log \frac{\exp(s_{i,j}/\tau)}{\sum_{k=1}^{2N} \mathbb{1}_{[k \neq i]} \exp(s_{i,k}/\tau)}$ 
   $\mathcal{L} = \frac{1}{2N} \sum_{k=1}^N [\ell(2k-1, 2k) + \ell(2k, 2k-1)]$ 
  update networks  $f$  and  $g$  to minimize  $\mathcal{L}$ 
end for
return encoder network  $f(\cdot)$ , and throw away  $g(\cdot)$ 

```

Для каждого примера в минибатче  
сэплируем две функции  
аугментации, применяем их, и  
считаем расстояние между  
признаками



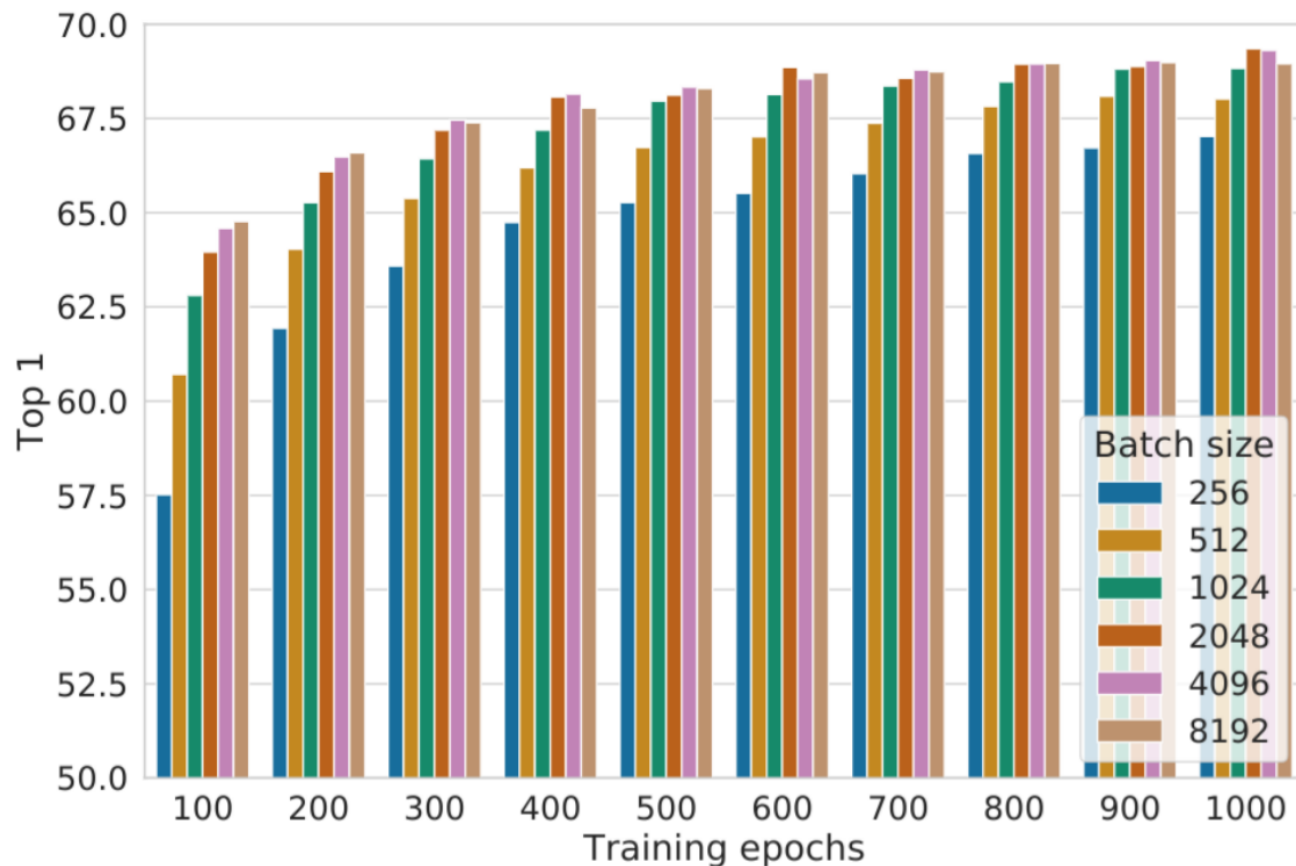
# SimCLR



Method	Architecture	Label fraction	
		1%	10%
		Top 5	
Supervised baseline	ResNet-50	48.4	80.4
<i>Methods using other label-propagation:</i>			
Pseudo-label	ResNet-50	51.6	82.4
VAT+Entropy Min.	ResNet-50	47.0	83.4
UDA (w. RandAug)	ResNet-50	-	88.5
FixMatch (w. RandAug)	ResNet-50	-	89.1
S4L (Rot+VAT+En. M.)	ResNet-50 (4×)	-	91.2
<i>Methods using representation learning only:</i>			
InstDisc	ResNet-50	39.2	77.4
BigBiGAN	RevNet-50 (4×)	55.2	78.8
PIRL	ResNet-50	57.2	83.8
CPC v2	ResNet-161(*)	77.9	91.2
SimCLR (ours)	ResNet-50	75.5	87.8
SimCLR (ours)	ResNet-50 (2×)	83.0	91.2
SimCLR (ours)	ResNet-50 (4×)	<b>85.8</b>	<b>92.6</b>

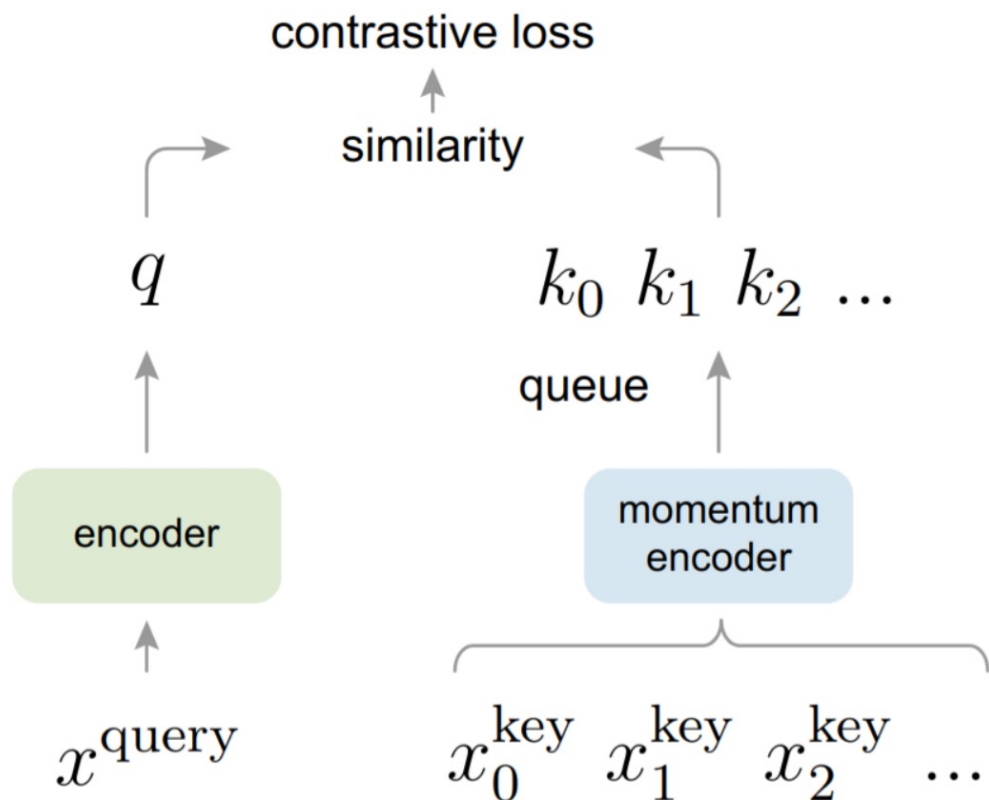
Table 7. ImageNet accuracy of models trained with few labels.

# SimCLR



- SimCLR требует обучения с большими батчами
- Это возможно только для распределённого обучения  
requires training with large
- SimCLR v2 – большие модели, больше слоёв в проекторе

# Momentum Contrast (MoCo)



$$\mathcal{L}_q = -\log \frac{\exp(q \cdot k^+ / \tau)}{\exp(q \cdot k^+ / \tau) + \sum_{k^-} \exp(q \cdot k^- / \tau)}$$

- Цель – уменьшить размер батча для работы на меньшем числе GPU
- Поддерживаем «очередь» из примеров (ключей), которые используем для оценки loss, но не используем их для расчёта градиента
- Для стабилизации «ключи» прогоняем через momentum encoder
- Кодировщик обновляется с большим моментом ( $m = 0.999$ )



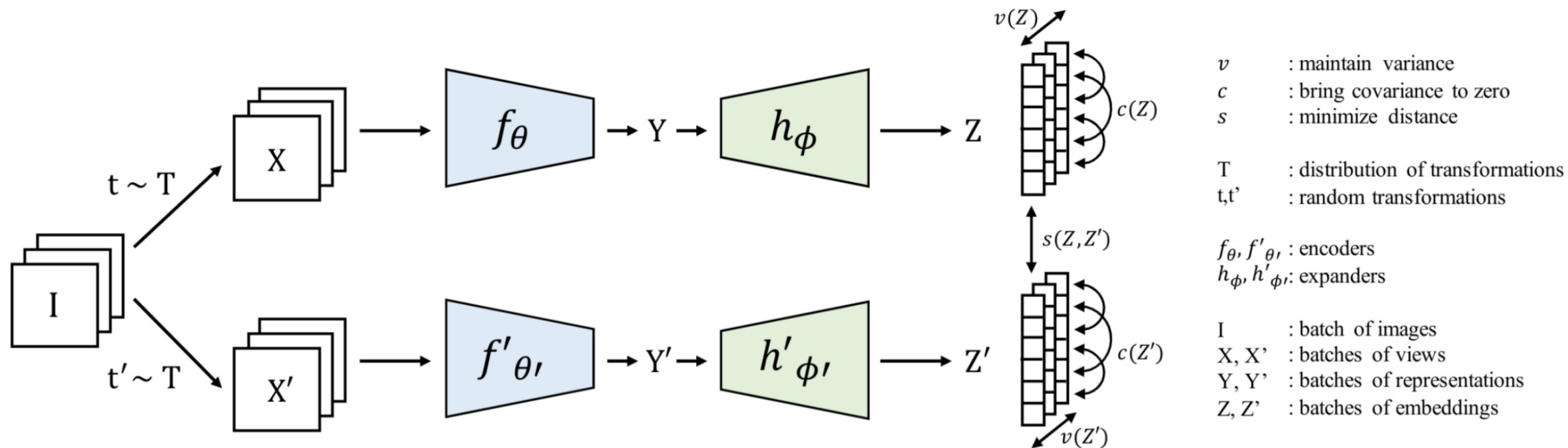
# MoCo v2



case	unsup. pre-train					ImageNet acc.
	MLP	aug+	cos	epochs	batch	
MoCo v1 [6]				200	256	60.6
SimCLR [2]	✓	✓	✓	200	256	61.9
SimCLR [2]	✓	✓	✓	200	8192	66.6
<b>MoCo v2</b>	✓	✓	✓	200	256	<b>67.5</b>
<i>results of longer unsupervised training follow:</i>						
SimCLR [2]	✓	✓	✓	1000	4096	69.3
<b>MoCo v2</b>	✓	✓	✓	800	256	<b>71.1</b>

- Нелинейный проектор и сильные аугментации необходимы для получения хорошего качества
- MoCo обгоняет SimCLR при заметно меньших батчах
- MoCo v2 мы можем обучать на узле с 8×V100 GPUs)

# VICReg



$$v(Z) = \frac{1}{d} \sum_{i=1}^d \max \left( 0, 1 - \sqrt{\text{Var}(z^i) + \varepsilon} \right)$$

$$c(Z) = \frac{1}{d} \sum_{i \neq j} [C(Z)]_{ij}^2, \quad C(Z) = \frac{1}{n-1} \sum_{i=1}^n (z_i - \bar{z})(z_i - \bar{z})^T$$

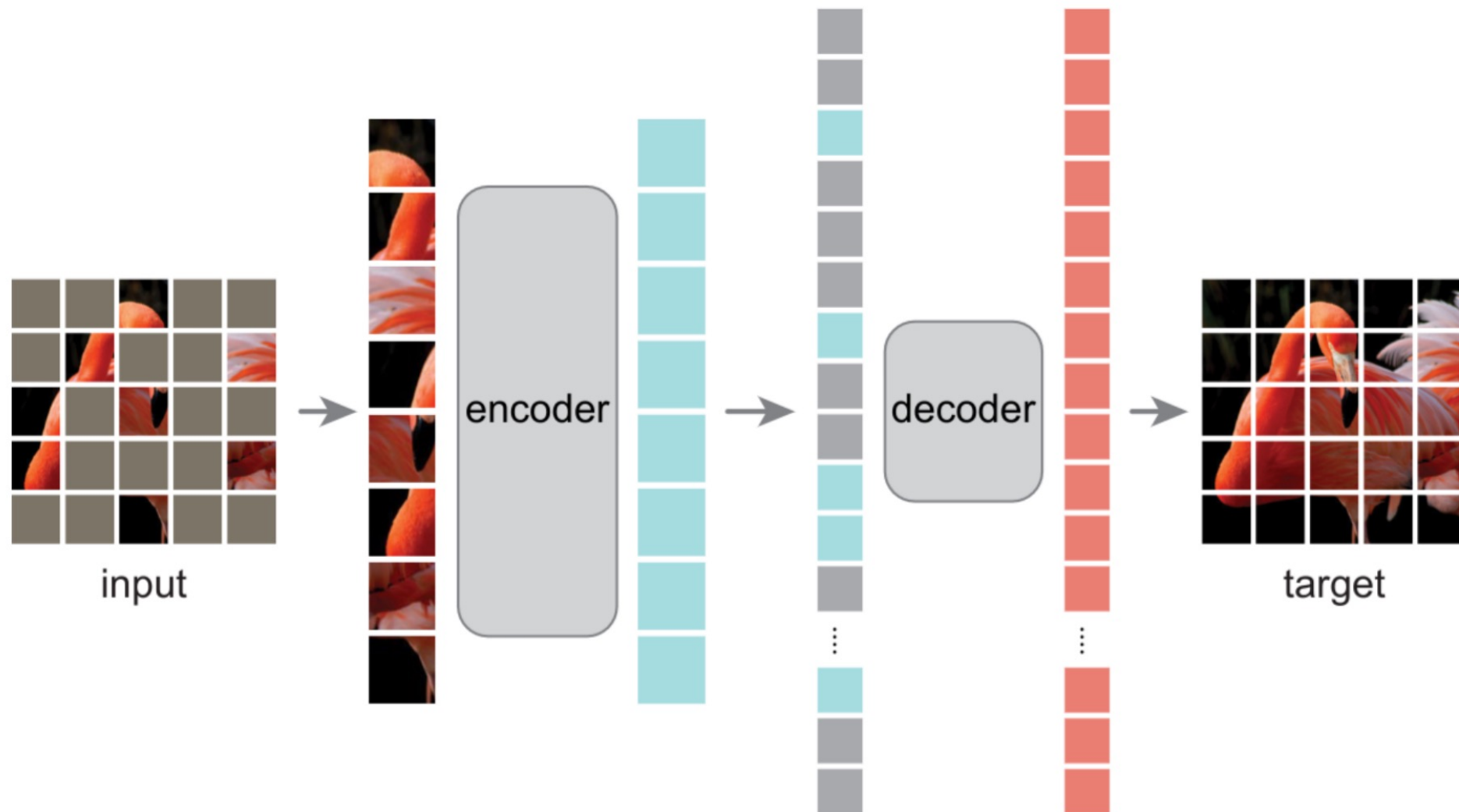
$$s(Z, Z') = \text{MSE}(Z, Z')$$



Method	Linear		Semi-supervised			
	Top-1	Top-5	Top-1		Top-5	
			1%	10%	1%	10%
Supervised	76.5	-	25.4	56.4	48.4	80.4
MoCo <a href="#">He et al. (2020)</a>	60.6	-	-	-	-	-
PIRL <a href="#">Misra &amp; Maaten (2020)</a>	63.6	-	-	-	57.2	83.8
CPC v2 <a href="#">Hénaff et al. (2019)</a>	63.8	-	-	-	-	-
CMC <a href="#">Tian et al. (2019)</a>	66.2	-	-	-	-	-
SimCLR <a href="#">Chen et al. (2020a)</a>	69.3	89.0	48.3	65.6	75.5	87.8
MoCo v2 <a href="#">Chen et al. (2020c)</a>	71.1	-	-	-	-	-
SimSiam <a href="#">Chen &amp; He (2020)</a>	71.3	-	-	-	-	-
SwAV <a href="#">Caron et al. (2020)</a>	71.8	-	-	-	-	-
InfoMin Aug <a href="#">Tian et al. (2020)</a>	73.0	<u>91.1</u>	-	-	-	-
OBoW <a href="#">Gidaris et al. (2021)</a>	<u>73.8</u>	-	-	-	<u>82.9</u>	<u>90.7</u>
BYOL <a href="#">Grill et al. (2020)</a>	<u>74.3</u>	<u>91.6</u>	53.2	68.8	78.4	89.0
SwAV (w/ multi-crop) <a href="#">Caron et al. (2020)</a>	<u>75.3</u>	-	<u>53.9</u>	<u>70.2</u>	78.5	<u>89.9</u>
Barlow Twins <a href="#">Zbontar et al. (2021)</a>	73.2	91.0	<u>55.0</u>	<u>69.7</u>	<u>79.2</u>	89.3
VICReg (ours)	73.2	<u>91.1</u>	<u>54.8</u>	<u>69.5</u>	<u>79.4</u>	<u>89.5</u>

Обучение выполняется на 32×V100 GPUs

# Masked Autoencoders



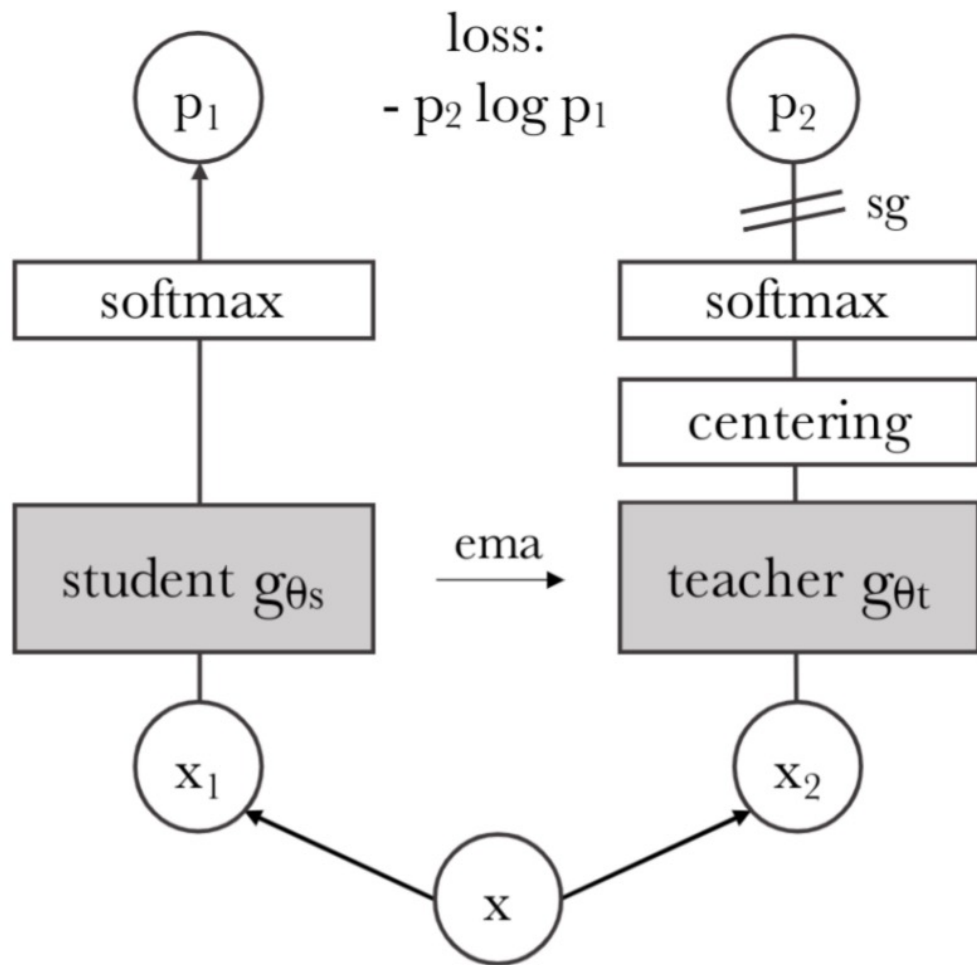


# Masked Autoencoders



method	pre-train data	ViT-B	ViT-L	ViT-H	ViT-H <sub>448</sub>
scratch, our impl.	-	82.3	82.6	83.1	-
DINO [5]	IN1K	82.8	-	-	-
MoCo v3 [9]	IN1K	83.2	84.1	-	-
BEiT [2]	IN1K+DALLE	83.2	85.2	-	-
MAE	IN1K	<u>83.6</u>	<u>85.9</u>	<u>86.9</u>	<b>87.8</b>

# DINO



*Supervised*



*DINO*



	Random	Supervised	DINO
ViT-S/16	22.0	27.3	45.9
ViT-S/8	21.8	23.7	44.7

# DINO



Method	Arch.	Param.	im/s	Linear	$k$ -NN
Supervised	RN50	23	1237	79.3	79.3
SCLR [12]	RN50	23	1237	69.1	60.7
MoCov2 [15]	RN50	23	1237	71.1	61.9
InfoMin [67]	RN50	23	1237	73.0	65.3
BarlowT [81]	RN50	23	1237	73.2	66.0
OBoW [27]	RN50	23	1237	73.8	61.9
BYOL [30]	RN50	23	1237	74.4	64.8
DCv2 [10]	RN50	23	1237	75.2	67.1
SwAV [10]	RN50	23	1237	<b>75.3</b>	65.7
DINO	RN50	23	1237	<b>75.3</b>	<b>67.5</b>
Supervised	ViT-S	21	1007	79.8	79.8
BYOL* [30]	ViT-S	21	1007	71.4	66.6
MoCov2* [15]	ViT-S	21	1007	72.7	64.4
SwAV* [10]	ViT-S	21	1007	73.5	66.3
DINO	ViT-S	21	1007	<b>77.0</b>	<b>74.5</b>
<i>Comparison across architectures</i>					
SCLR [12]	RN50w4	375	117	76.8	69.3
SwAV [10]	RN50w2	93	384	77.3	67.3
BYOL [30]	RN50w2	93	384	77.4	–
DINO	ViT-B/16	85	312	78.2	76.1
SwAV [10]	RN50w5	586	76	78.5	67.1
BYOL [30]	RN50w4	375	117	78.6	–
BYOL [30]	RN200w2	250	123	79.6	73.9
DINO	ViT-S/8	21	180	79.7	<b>78.3</b>
SCLRv2 [13]	RN152w3+SK	794	46	79.8	73.1
DINO	ViT-B/8	85	63	<b>80.1</b>	77.4

Обучение на одном узле 8×V100 GPUs



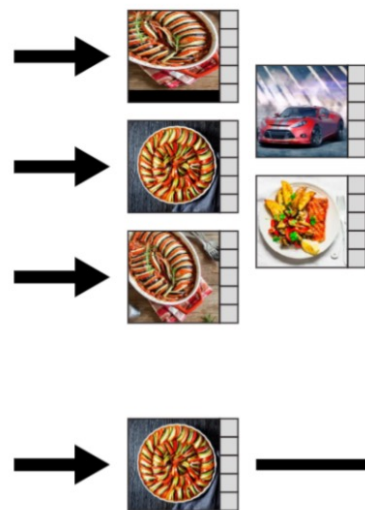
# DINOv2



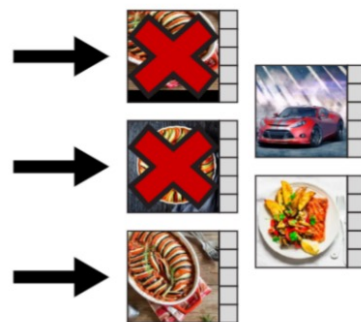
Uncurated Data



Curated Data



Embedding



Deduplication



Retrieval

Augmented Curated Data



Данные + большая модель (ViT-g, 1.1B params) с дистилляцией + несколько функций потерь и регуляризации + эффективное выполнение. Код для обучения и веса выложены в опенсор с коммерческой лицензией



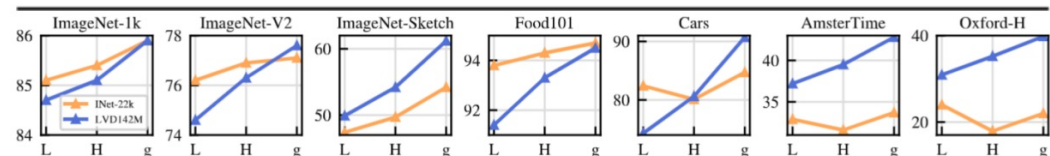
# DINOv2



Task	Dataset / Split	Images	Retrieval	Retrieved	Final
classification	ImageNet-22k / –	14,197,086	as is	–	14,197,086
classification	ImageNet-22k / –	14,197,086	sample	56,788,344	56,788,344
classification	ImageNet-1k / train	1,281,167	sample	40,997,344	40,997,344
fine-grained classif.	Caltech 101 / train	3,030	cluster	2,630,000	1,000,000
fine-grained classif.	CUB-200-2011 / train	5,994	cluster	1,300,000	1,000,000
fine-grained classif.	DTD / train1	1,880	cluster	1,580,000	1,000,000
fine-grained classif.	FGVC-Aircraft / train	3,334	cluster	1,170,000	1,000,000
fine-grained classif.	Flowers-102 / train	1,020	cluster	1,060,000	1,000,000
fine-grained classif.	Food-101 / train	75,750	cluster	21,670,000	1,000,000
fine-grained classif.	Oxford-IIIT Pet / trainval	3,680	cluster	2,750,000	1,000,000
fine-grained classif.	Stanford Cars / train	8,144	cluster	7,220,000	1,000,000
fine-grained classif.	SUN397 / train1	19,850	cluster	18,950,000	1,000,000
fine-grained classif.	Pascal VOC 2007 / train	2,501	cluster	1,010,000	1,000,000
segmentation	ADE20K / train	20,210	cluster	20,720,000	1,000,000
segmentation	Cityscapes / train	2,975	cluster	1,390,000	1,000,000
segmentation	Pascal VOC 2012 (seg.) / trainaug	1,464	cluster	10,140,000	1,000,000
depth estimation	Mapillary SLS / train	1,434,262	as is	–	1,434,262
depth estimation	KITTI / train (Eigen)	23,158	cluster	3,700,000	1,000,000
depth estimation	NYU Depth V2 / train	24,231	cluster	10,850,000	1,000,000
depth estimation	SUN RGB-D / train	4,829	cluster	4,870,000	1,000,000
retrieval	Google Landmarks v2 / train (clean)	1,580,470	as is	–	1,580,470
retrieval	Google Landmarks v2 / train (clean)	1,580,470	sample	6,321,880	6,321,880
retrieval	AmsterTime / new	1,231	cluster	960,000	960,000
retrieval	AmsterTime / old	1,231	cluster	830,000	830,000
retrieval	Met / train	397,121	cluster	62,860,000	1,000,000
retrieval	Revisiting Oxford / base	4,993	cluster	3,680,000	1,000,000
retrieval	Revisiting Paris / base	6,322	cluster	3,660,000	1,000,000

142,109,386

Training Data	INet-1k	Im-A	ADE-20k	Oxford-M
INet-22k	85.9	73.5	46.6	62.5
INet-22k \ INet-1k	85.3	70.3	46.2	58.7
Uncurated data	83.3	59.4	48.5	54.3
LVD-142M	85.8	73.9	47.7	64.6



# DINOv2



Method	Arch.	Data	Text sup.	kNN	linear		
				val	val	ReaL	V2
Weakly supervised							
CLIP	ViT-L/14	WIT-400M	✓	79.8	84.3	88.1	75.3
CLIP	ViT-L/14 <sub>336</sub>	WIT-400M	✓	80.5	85.3	88.8	75.8
SWAG	ViT-H/14	IG3.6B	✓	82.6	85.7	88.7	77.6
OpenCLIP	ViT-H/14	LAION	✓	81.7	84.4	88.4	75.5
OpenCLIP	ViT-G/14	LAION	✓	83.2	86.2	89.4	77.2
EVA-CLIP	ViT-g/14	custom*	✓	<b>83.5</b>	86.4	89.3	77.4
Self-supervised							
MAE	ViT-H/14	INet-1k	✗	49.4	76.6	83.3	64.8
DINO	ViT-S/8	INet-1k	✗	78.6	79.2	85.5	68.2
SEERv2	RG10B	IG2B	✗	—	79.8	—	—
MSN	ViT-L/7	INet-1k	✗	79.2	80.7	86.0	69.7
EsViT	Swin-B/W=14	INet-1k	✗	79.4	81.3	87.0	70.4
Mugs	ViT-L/16	INet-1k	✗	80.2	82.1	86.9	70.8
iBOT	ViT-L/16	INet-22k	✗	72.9	82.3	87.5	72.4
DINOv2	ViT-S/14	LVD-142M	✗	79.0	81.1	86.6	70.9
	ViT-B/14	LVD-142M	✗	82.1	84.5	88.3	75.1
	ViT-L/14	LVD-142M	✗	<b>83.5</b>	86.3	89.5	78.0
	ViT-g/14	LVD-142M	✗	<b>83.5</b>	<b>86.5</b>	<b>89.6</b>	<b>78.4</b>



1. Введение
2. Прокси-задачи
3. Констранстное обучение и маскирование
4. Фундаментальные модели

# Что такое фундаментальная модель?

---



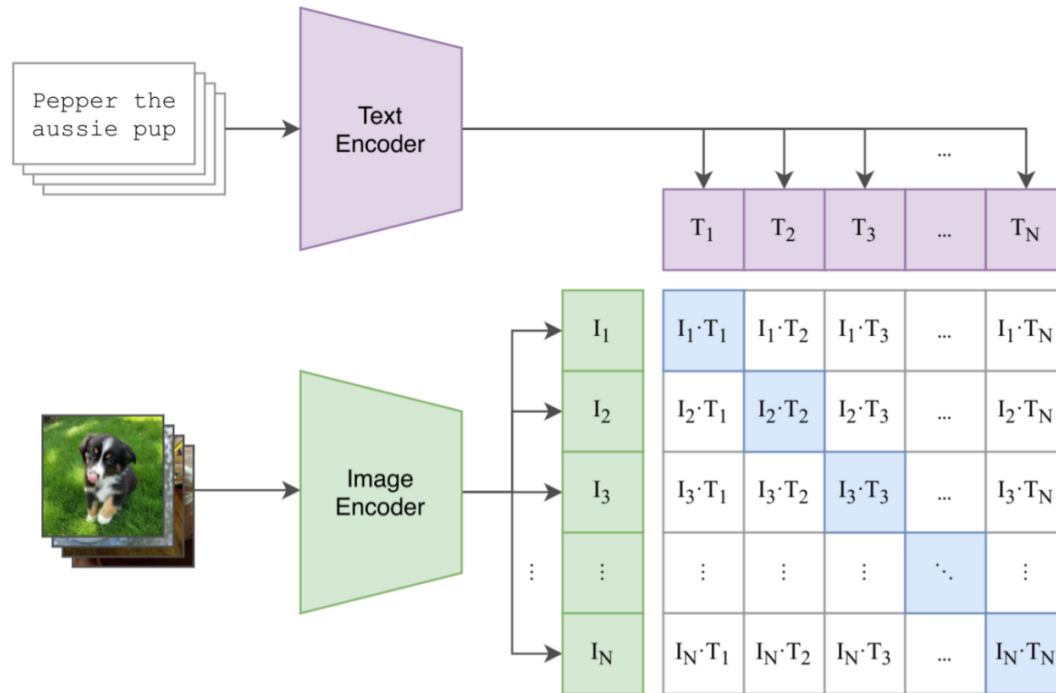
Модель называют «фундаментальной», если:

- Работает с несколькими модальностями, например, текстом и изображениями
- Решают несколько задач в разных доменах, например, классификацию, сегментацию, описание, ответы на вопросы
- Работает с запросами, т.е. поддерживает несколько вариантов запросов относительно анализируемой информации
- Хорошо работает без тонкой настройки на целевой бенчмарк

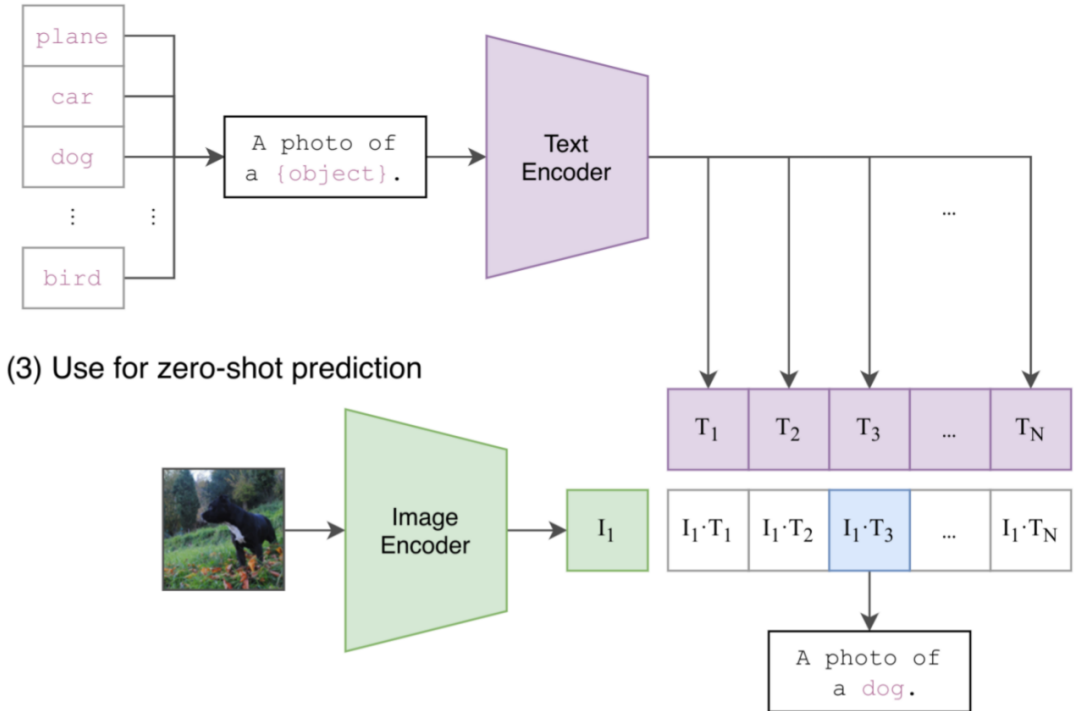
# CLIP



(1) Contrastive pre-training



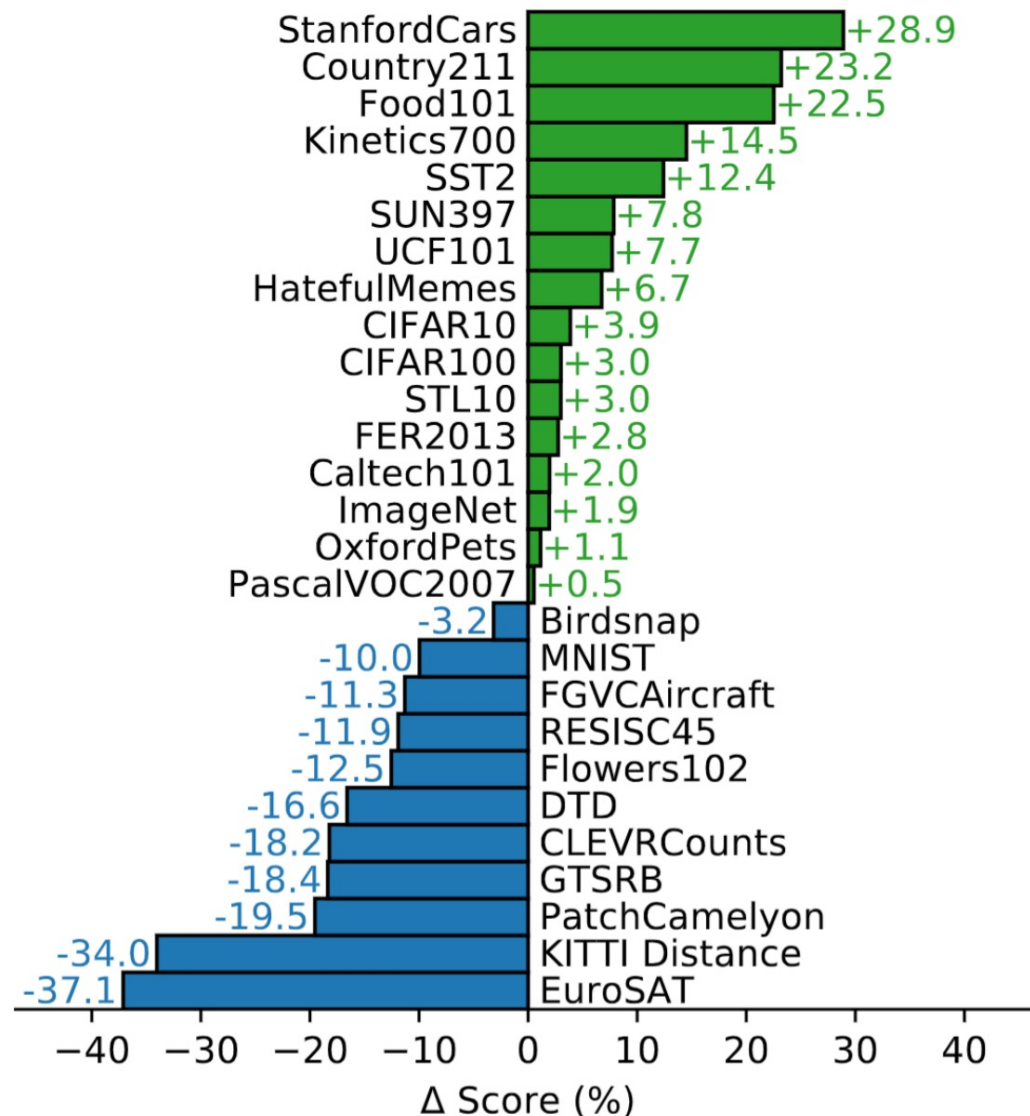
(2) Create dataset classifier from label text



(3) Use for zero-shot prediction

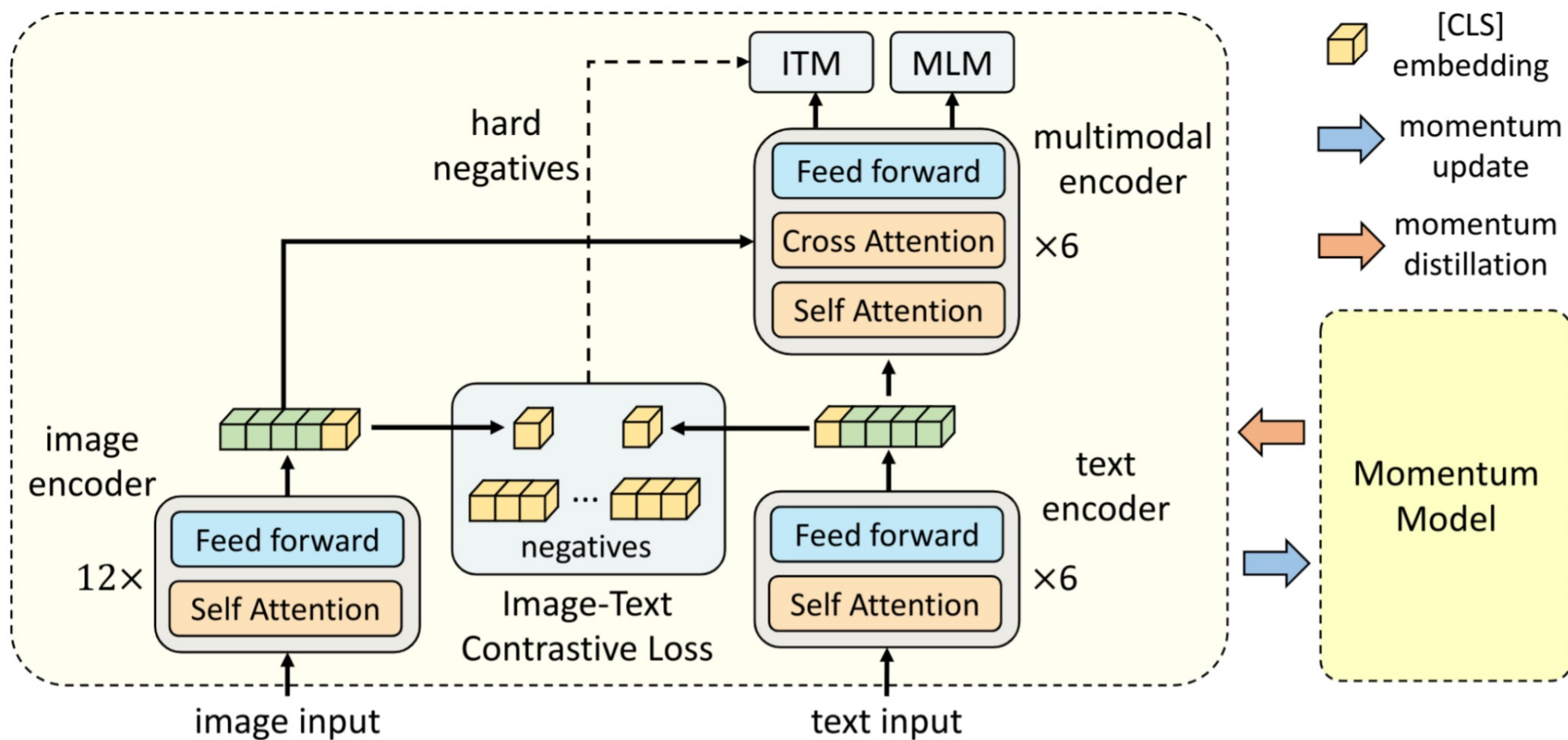
400M (image, text) пар, 500×V100 GPUs для предобучения

# CLIP zero-shot results



Zero-Shot CLIP vs. Linear Probe on ResNet50

# ALBEF



Loss: contrastive loss + image-text matching + masked language modelling



# ALBEF pseudo-targets



## masked language modelling ↓

“polar bear in the [MASK]”

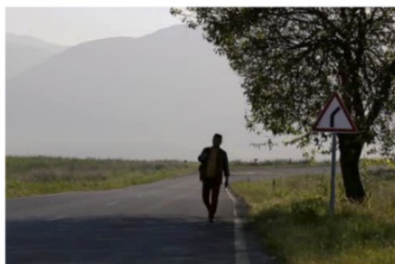


GT: wild

Top-5 pseudo-targets:

1. zoo
2. pool
3. water
4. pond
5. wild

“a man [MASK] along a road in front of nature in summer”

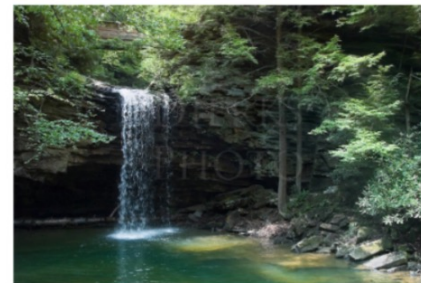


GT: standing

Top-5 pseudo-targets:

1. walks
2. walking
3. runs
4. running
5. goes

“a [MASK] waterfall in the deep woods”



GT: remote

Top-5 pseudo-targets:

1. small
2. beautiful
3. little
4. secret
5. secluded



GT: breakdown of the car on the road

Top-5 pseudo-targets:

1. young woman get out of the car near the road
2. a woman inspects her damaged car under a tree
3. a woman looking into a car after locking her keys inside
4. young woman with a broken car calling for help
5. breakdown of the car on the road



GT: the harbor a small village

Top-5 pseudo-targets:

1. the harbour with boats and houses
2. replica of the sailing ship in the harbour
3. ships in the harbor of the town
4. the harbor a small village
5. boats lined up alongside the geographical feature category in the village

## image-text matching ↑



# Сравнение ALBEF с CLIP



Method	# Pre-train Images	Flickr30K (1K test set)					
		TR			IR		
		R@1	R@5	R@10	R@1	R@5	R@10
UNITER [2]	4M	83.6	95.7	97.7	68.7	89.2	93.9
CLIP [6]	400M	88.0	98.7	99.4	68.7	90.6	95.2
ALIGN [7]	1.2B	88.6	98.7	99.7	75.7	93.8	96.8
ALBEF	4M	90.5	98.8	99.7	76.8	93.7	96.7
ALBEF	14M	<b>94.1</b>	<b>99.5</b>	<b>99.7</b>	<b>82.8</b>	<b>96.3</b>	<b>98.1</b>

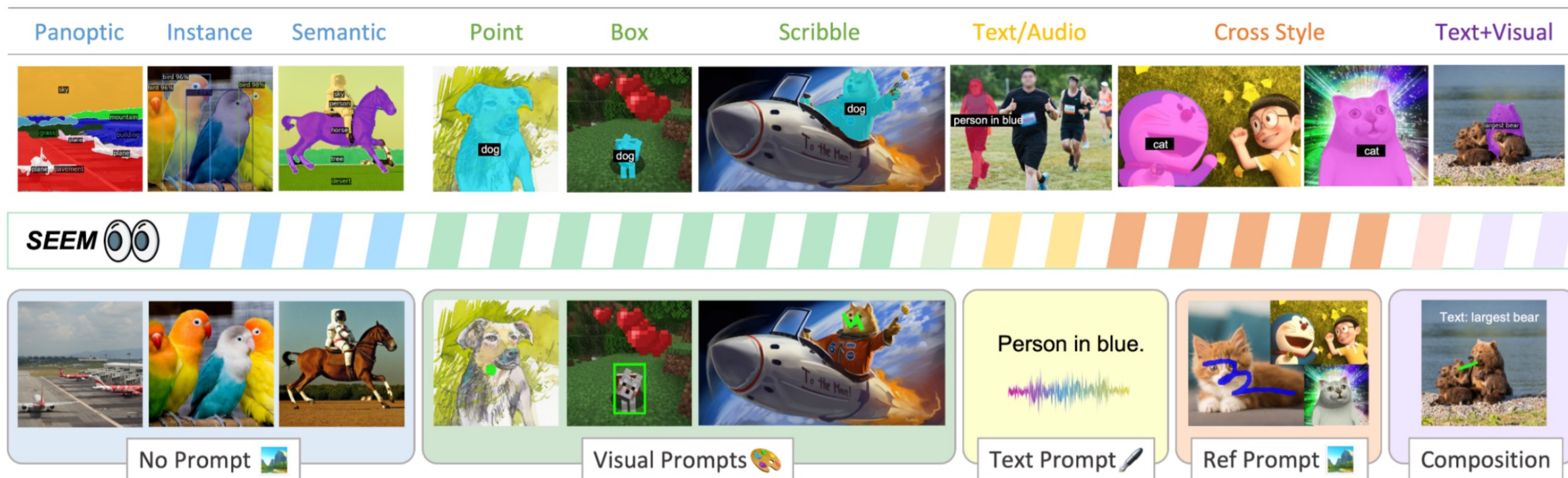
Table 3: Zero-shot image-text retrieval results on Flickr30K.

Method	VQA		NLVR <sup>2</sup>		SNLI-VE	
	test-dev	test-std	dev	test-P	val	test
VisualBERT [13]	70.80	71.00	67.40	67.00	-	-
VL-BERT [10]	71.16	-	-	-	-	-
LXMERT [1]	72.42	72.54	74.90	74.50	-	-
12-in-1 [12]	73.15	-	-	78.87	-	76.95
UNITER [2]	72.70	72.91	77.18	77.85	78.59	78.28
VL-BART/T5 [54]	-	71.3	-	73.6	-	-
ViLT [21]	70.94	-	75.24	76.21	-	-
OSCAR [3]	73.16	73.44	78.07	78.36	-	-
VILLA [8]	73.59	73.67	78.39	79.30	79.47	79.03
ALBEF (4M)	74.54	74.70	80.24	80.50	80.14	80.30
ALBEF (14M)	<b>75.84</b>	<b>76.04</b>	<b>82.55</b>	<b>83.14</b>	<b>80.80</b>	<b>80.91</b>

Table 4: Comparison with state-of-the-art methods on downstream vision-language tasks.

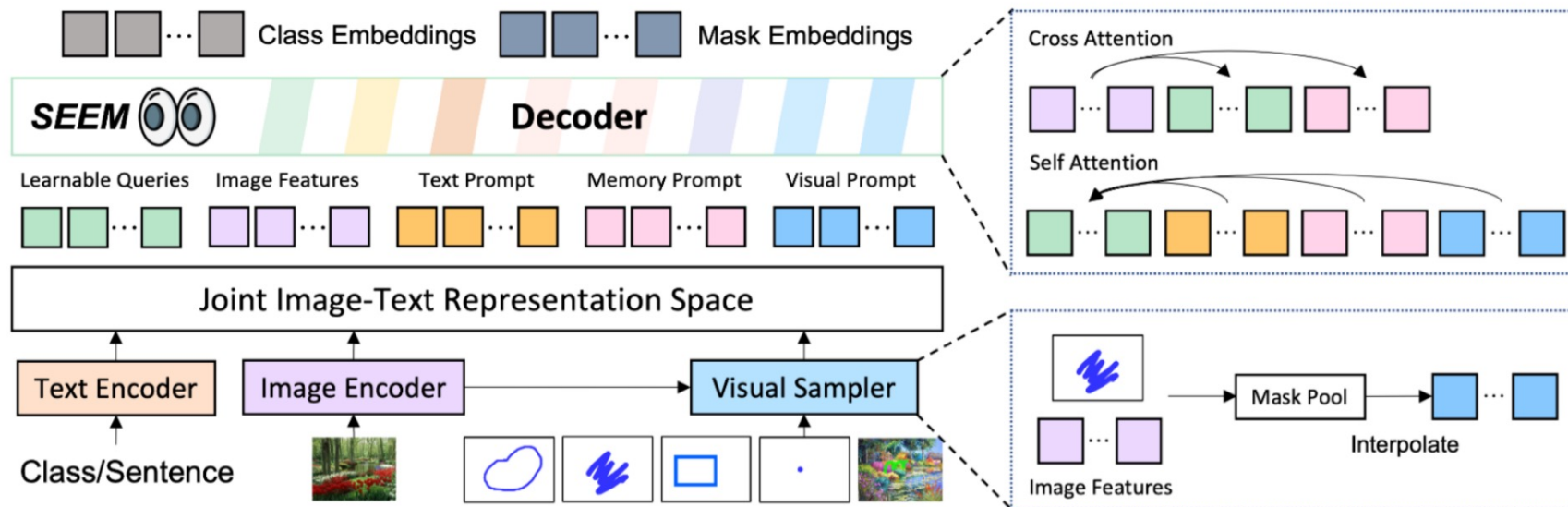
8×A100 GPUs для  
предобучения

# SEEM

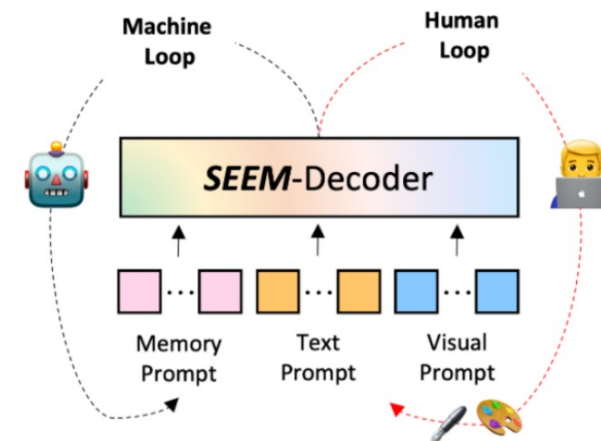


Panoptic + referring + interactive segmentation

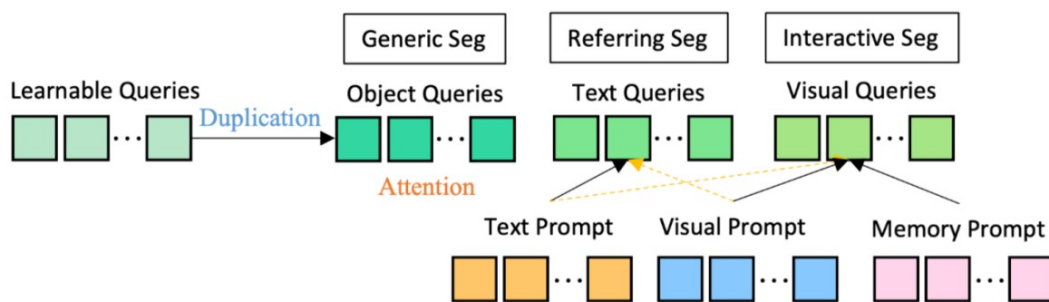
# SEEM



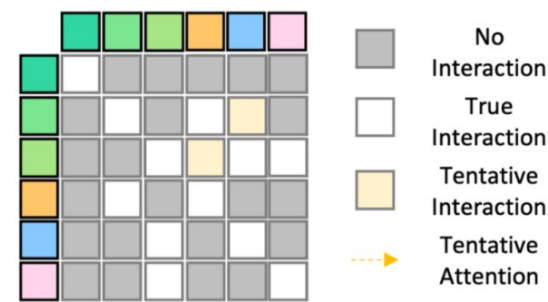
(a) Model Architecture



(b) Human-Model Interaction



(a) Queries and Prompt Interaction



(b) Self-Attention Mask



# SEEM



Method	Segmentation Data	Type	Generic Segmentation COCO			Referring Segmentation RefCOCOg			Interactive Segmentation PascalVOC					
			PQ	mAP	mIoU	cIoU	mIoU	AP50	5-NoC85	10-NoC85	20-NoC85	5-NoC90	10-NoC90	20-NoC90
Mask2Former (T) <a href="#">[6]</a>	COCO (0.12M)	Segmentation	53.2	43.3	63.2	-	-	-	-	-	-	-	-	-
Mask2Former (B) <a href="#">[6]</a>	COCO (0.12M)		56.4	46.3	67.1	-	-	-	-	-	-	-	-	-
Mask2Former (L) <a href="#">[6]</a>	COCO (0.12M)		57.8	48.6	67.4	-	-	-	-	-	-	-	-	-
Pano/SegFormer (B) <a href="#">[45]</a>	COCO (0.12M)		55.4	*	*	-	-	-	-	-	-	-	-	-
LAVT (B) <a href="#">[53]</a>	Ref-COCO (0.03M)		-	-	-	61.2	*	*	-	-	-	-	-	-
PolyFormer (B) <a href="#">[17]</a>	Ref-COCO+VG+... (0.16M)		-	-	-	69.3	*	*	-	-	-	-	-	-
PolyFormer (L) <a href="#">[17]</a>	Ref-COCO+VG+... (0.16M)		-	-	-	71.1	*	*	-	-	-	-	-	-
RITM (<T) <a href="#">[18]</a>	COCO+LVIS (0.12M)	Interactive	-	-	-	-	-	-	*	*	2.19	*	*	2.57
PseudoClick (<T) <a href="#">[54]</a>	COCO (0.12M)		-	-	-	-	-	-	*	*	1.94	*	*	2.25
FocalClick (T) <a href="#">[21]</a>	COCO (0.12M)		-	-	-	-	-	-	*	*	2.97	*	*	3.52
FocalClick (B) <a href="#">[21]</a>	COCO (0.12M)		-	-	-	-	-	-	*	*	2.46	*	*	2.88
SimpleClick (B) <a href="#">[20]</a>	COCO+LVIS (0.12M)		-	-	-	-	-	-	1.75	1.93	2.06	1.94	2.19	2.38
SimpleClick (L) <a href="#">[20]</a>	COCO+LVIS (0.12M)		-	-	-	-	-	-	1.52	1.64	1.72	1.67	1.84	1.96
SimpleClick (H) <a href="#">[20]</a>	COCO+LVIS (0.12M)		-	-	-	-	-	-	1.51	1.64	1.76	1.64	1.83	1.98
UViM (L) <a href="#">[55]</a>	COCO (0.12M)	Generalist	45.8	*	*	-	-	-	-	-	-	-	-	-
Pix2Seq v2 (B) <a href="#">[56]</a>	COCO (0.12M)		-	38.2	-	-	-	-	-	-	-	-	-	-
X-Decoder (T) <a href="#">[11]</a>	COCO (0.12M)		52.6	41.3	62.4	59.8	*	*	-	-	-	-	-	-
X-Decoder (B) <a href="#">[11]</a>	COCO (0.12M)		56.2	45.8	66.0	64.5	*	*	-	-	-	-	-	-
X-Decoder (L) <a href="#">[11]</a>	COCO (0.12M)		56.9	46.7	67.5	64.6	*	*	-	-	-	-	-	-
UNINEXT (T) <a href="#">[48]</a>	Image+Video (3M)		-	44.9	-	70.0	*	*	-	-	-	-	-	-
UNINEXT (L) <a href="#">[48]</a>	Image+Video (3M)		-	49.6	-	73.4	*	*	-	-	-	-	-	-
Painter (L) <a href="#">[57]</a>	COCO+ADE+NYUv2 (0.16M)		43.4	*	*	-	-	-	-	-	-	-	-	-
#SegGPT (L) <a href="#">[50]</a>	COCO+ADE+NYUv2 (0.16M)		34.4	*	*	-	-	-	-	-	-	-	-	-
#SAM (B) <a href="#">[36]</a>	SAM (11M)		-	-	-	-	-	-	2.47	2.65	3.28	2.23	3.13	4.12
#SAM (L) <a href="#">[36]</a>	SAM (11M)		-	-	-	-	-	-	1.85	2.15	2.60	2.01	2.46	3.12
#SAM (H) <a href="#">[36]</a>	SAM (11M)		-	-	-	-	-	-	1.82	2.13	2.55	1.98	2.43	3.11
SEEM (T)	COCO+LVIS (0.12M)		50.8	39.7	62.2	60.9	65.7	74.8	1.72	2.30	3.37	1.97	2.83	4.41
SEEM (B)	COCO+LVIS (0.12M)		56.1	46.4	66.3	65.0	69.6	78.2	1.56	2.04	2.93	1.77	2.47	3.79
SEEM (L)	COCO+LVIS (0.12M)		57.5	47.7	67.6	65.6	70.3	78.9	1.51	1.95	2.77	1.71	2.36	3.61
SEEM (T)	COCO+LVIS (0.12M)	Composition	-	-	-	70.4	71.7	82.1	1.72	2.28	3.32	1.97	2.82	4.37
SEEM (B)	COCO+LVIS (0.12M)		-	-	-	76.2	77.8	87.8	1.56	2.03	2.91	1.77	2.46	3.76
SEEM (L)	COCO+LVIS (0.12M)		-	-	-	75.1	76.9	86.8	1.52	1.97	2.81	1.72	2.38	3.64

# Florence-2

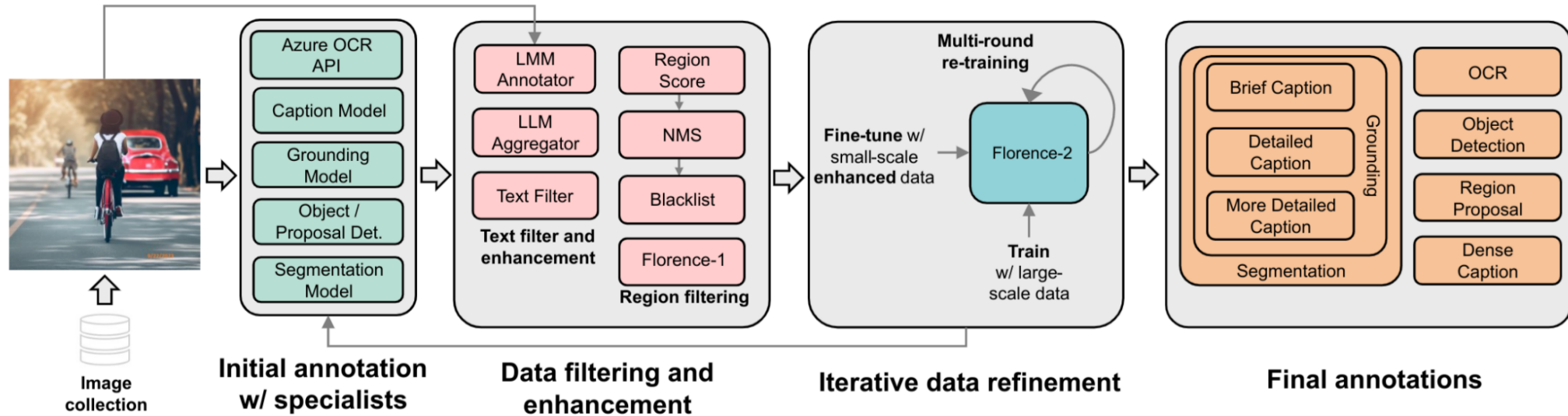
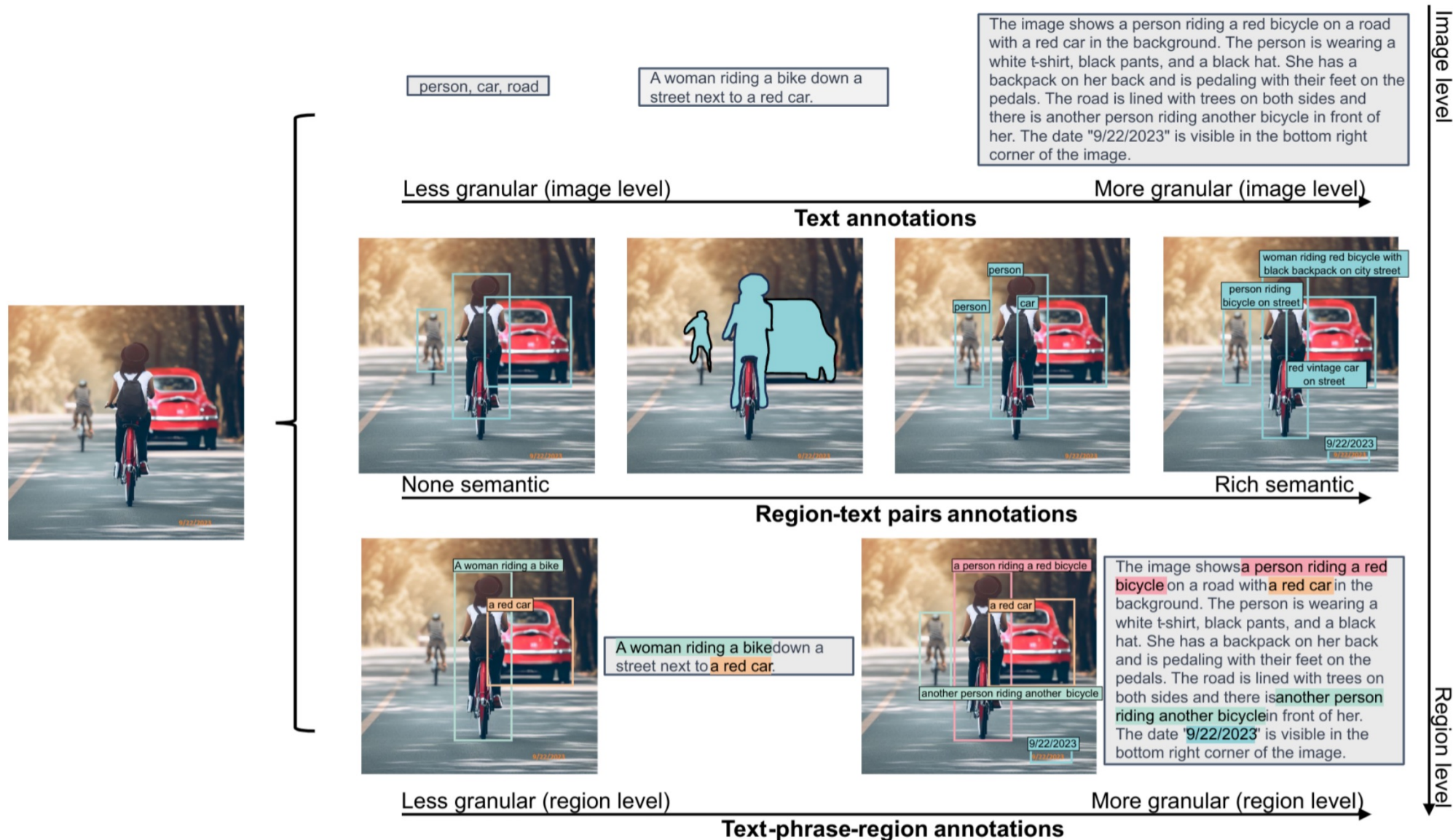


Figure 3. **Florence-2 data engine** consists of three essential phrases: (1) initial annotation employing specialist models, (2) data filtering to correct errors and remove irrelevant annotations, and (3) an iterative process for data refinement. Our final dataset (**FLD-5B**) of over **5B** annotations contains **126M** images, **500M** text annotations, **1.3B** region-text annotations, and **3.6B** text-phrase-region annotations.

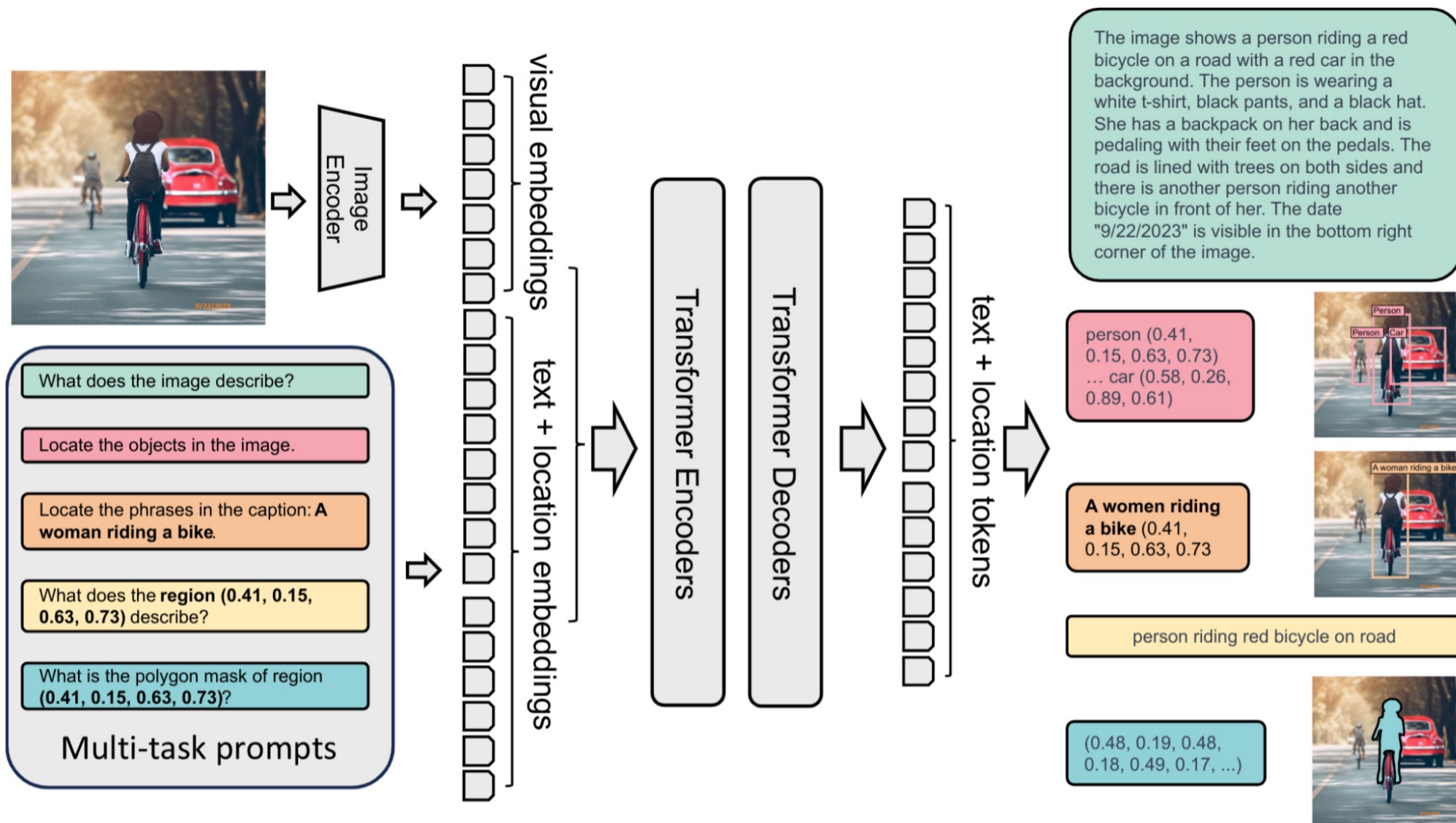
Dataset	Rep. Model	#Images	#Annotations	Spatial hierarchy	Semantics granularity
JFT300M [21]	ViT	300M	300M	Image-level	Coarse
WIT [64]	CLIP	400M	400M	Image-level	Coarse
SA-1B [32]	SAM	11M	1B	Region-level	Non-semantic
GrIT [60]	Kosmos-2	91M	137M	Image & Region-level	Fine-grained
M3W [2]	Flamingo	185M	43.3M*	Multi-image-level	Fine-grained
<b>FLD-5B (ours)</b>	<b>Florence-2 (ours)</b>	<b>126M</b>	<b>5B</b>	<b>Image &amp; Region-level</b>	<b>Coarse to fine-grained</b>



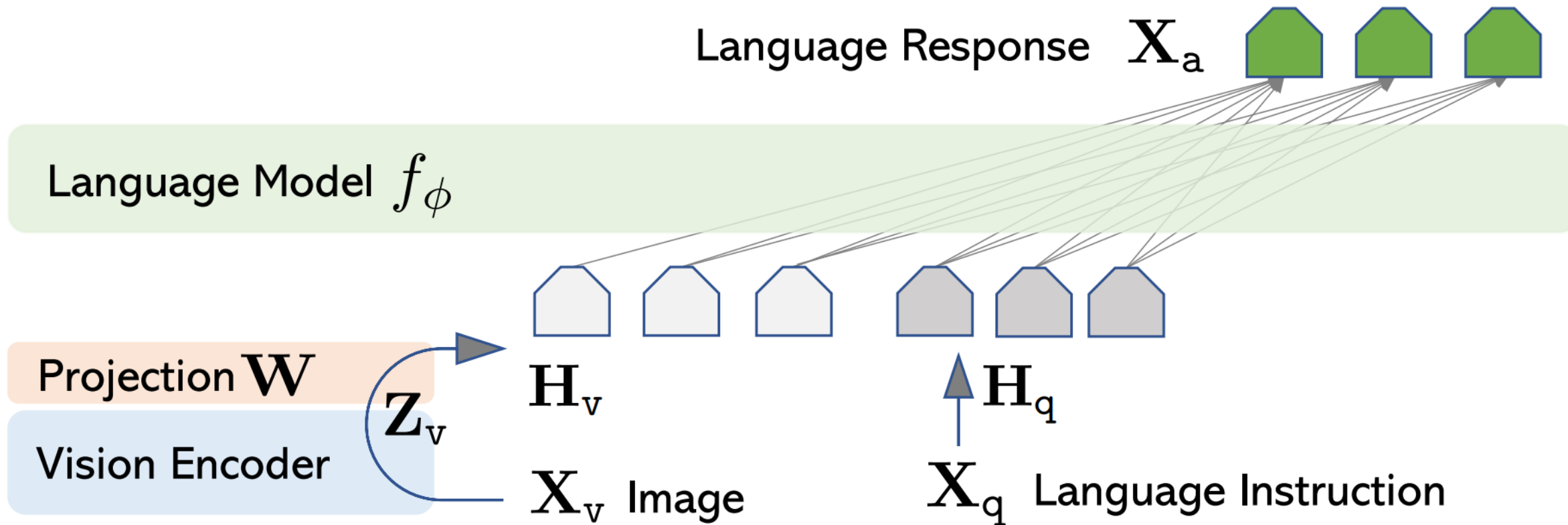
# Florence 2



# Florence 2



# LLaVA: Large Language and Vision Assistant





- Использование неразмеченных данных выгодно, т.к. Их много а размечать долго
- Первый подход к само-обучению заключался в эвристических прокси-задачах
- Контрастное обучение и маскирование – мощные подходы к само-обучению
- Мы рассмотрели определение и несколько примеров фундаментальных моделей