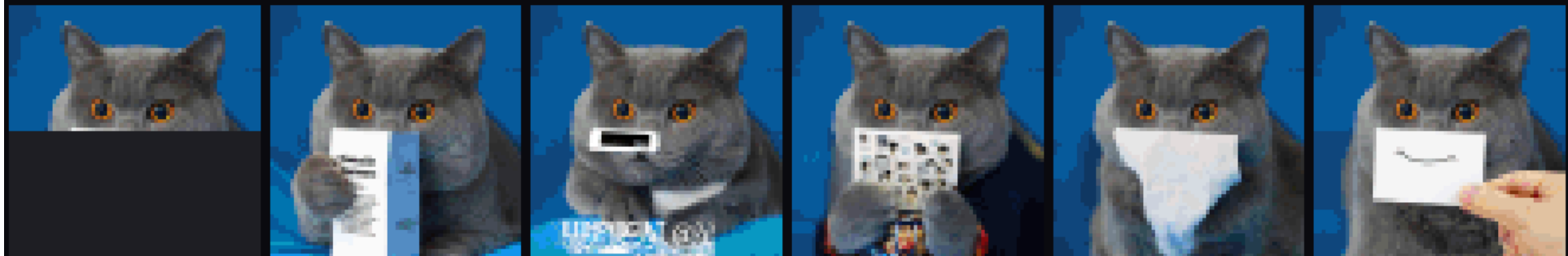


Image Pretraining from Pixels



Сеньченко Тимофей

Группа 172

Unsupervised and self-supervised learning

- Успешно используются в NLP (Авторегрессионные модели, BERT)
- Практически нет аналогичных современных применений в задачах связанных с изображениями
- Хорошие генеративные модели выучивают более хорошие представления?
- Воспользуемся domain-agnostic трансформером для авторегрессионного предсказания пикселей

Pretraining

$$x = (x_1, \dots, x_n)$$

$$p(x) = \prod_{i=1}^n p(x_{\pi_i} | x_{\pi_1}, \dots, x_{\pi_{i-1}}, \theta)$$

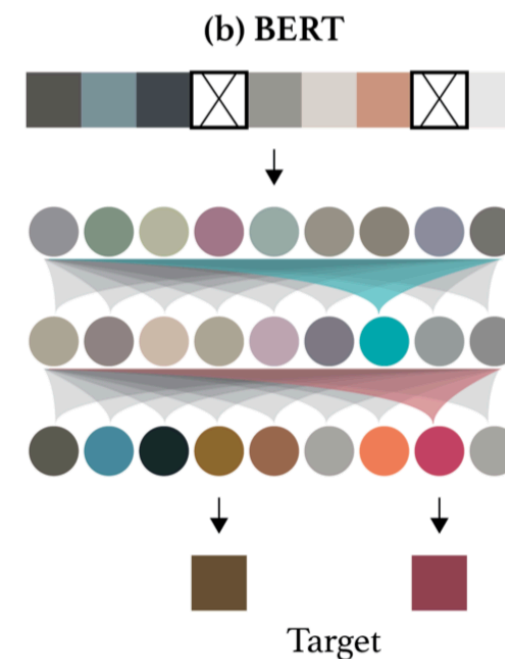
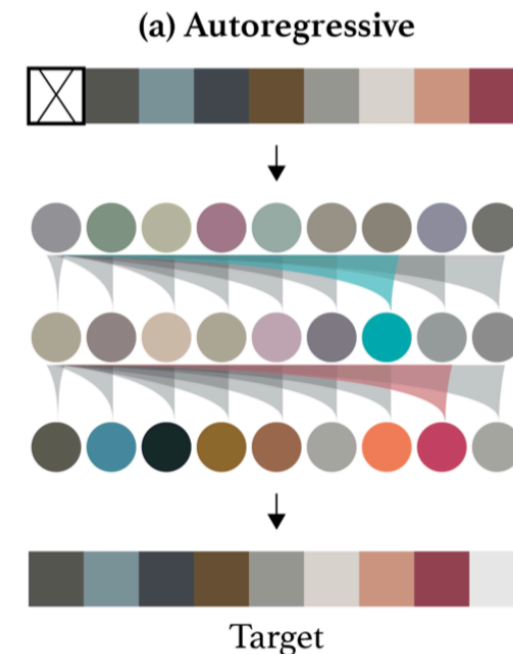
Для изображений будем минимизировать:

$$L_{AR} = \mathbb{E}_{x \sim X} [-\log p(x)]$$

Или:

$$L_{BERT} = \mathbb{E}_{x \sim X} \mathbb{E}_M \sum_{i \in M} [-\log p(x_i | x_{[1,n] \setminus M})]$$

$$M \subset [1, n]$$



Architecture

Вход: токены x_1, \dots, x_n

-> эмбединги размерности d для каждой позиции

Декодер состоит из L блоков, где l -й блок подает на выход промежуточные эмбединги h_1^l, \dots, h_n^l размерности d

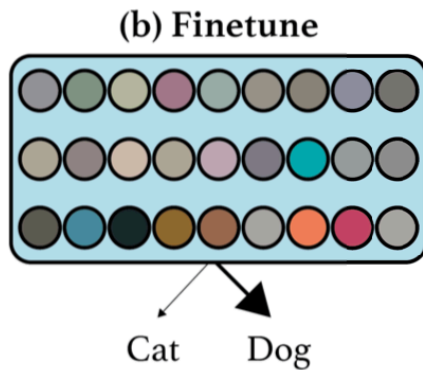
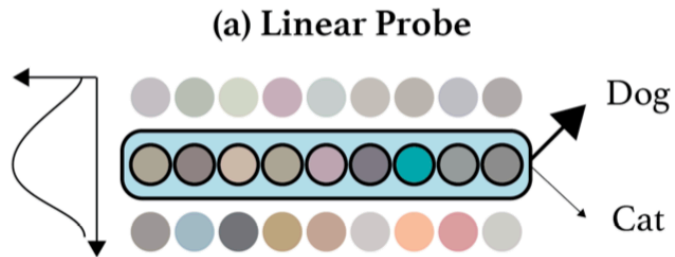
Авторы статьи используют структуру блока GPT-2: (тензор h^l - вход)

$$n^l = \text{layer_norm}(h^l)$$

$$a^l = h^l + \text{multihead_attention}(n^l)$$

$$h^{l+1} = a^l + \text{mlp}(\text{layer_norm}(a^l))$$

Fine-tuning and Linear Probing



Fine-tuning:

- Добавляем голову классификатор,
- Дообучаем на лейблах, используя выученные представления
- Loss: кросс-энтропия L_{CLF} , а лучше: $L_{GEN} + L_{CLF}$
 $L_{GEN} \in \{L_{AR}, L_{BERT}\}$

Linear Probing:

- Выбираем некоторый слой i и по его выходам учим Логистическую регрессию на лейблах

Context reduction

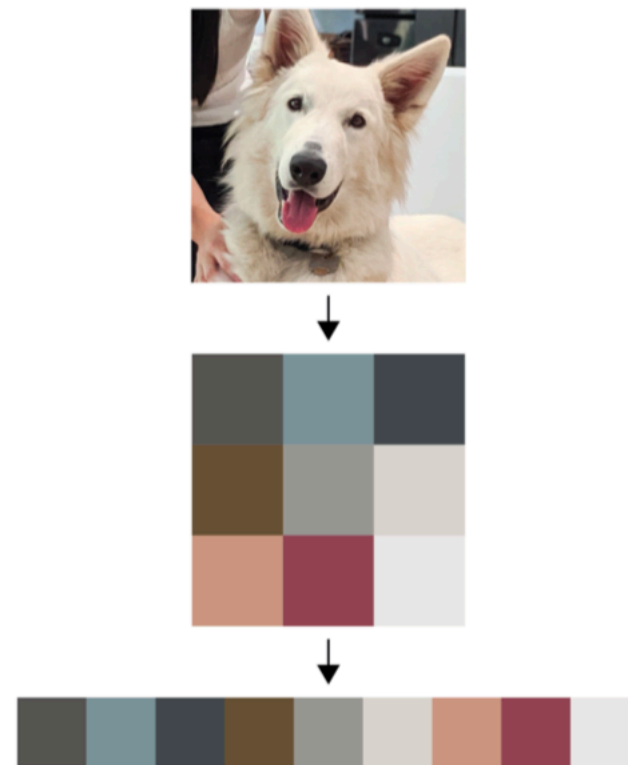
Входные данные значительно большей размерности, чем в NLP задачах ($224^2 \times 3$)

-> не можем использовать внимание с исходными данными

-> уменьшаем Image Resolution ($32^2 \times 3$ или $48^2 \times 3$ или $64^2 \times 3$)

-> переводим картинку из палитры RGB в новую 9-ти битную палитру с помощью k-means кластеризации

-> получаем изображения значительно меньшей размерности (32^2 или 48^2 или 64^2)



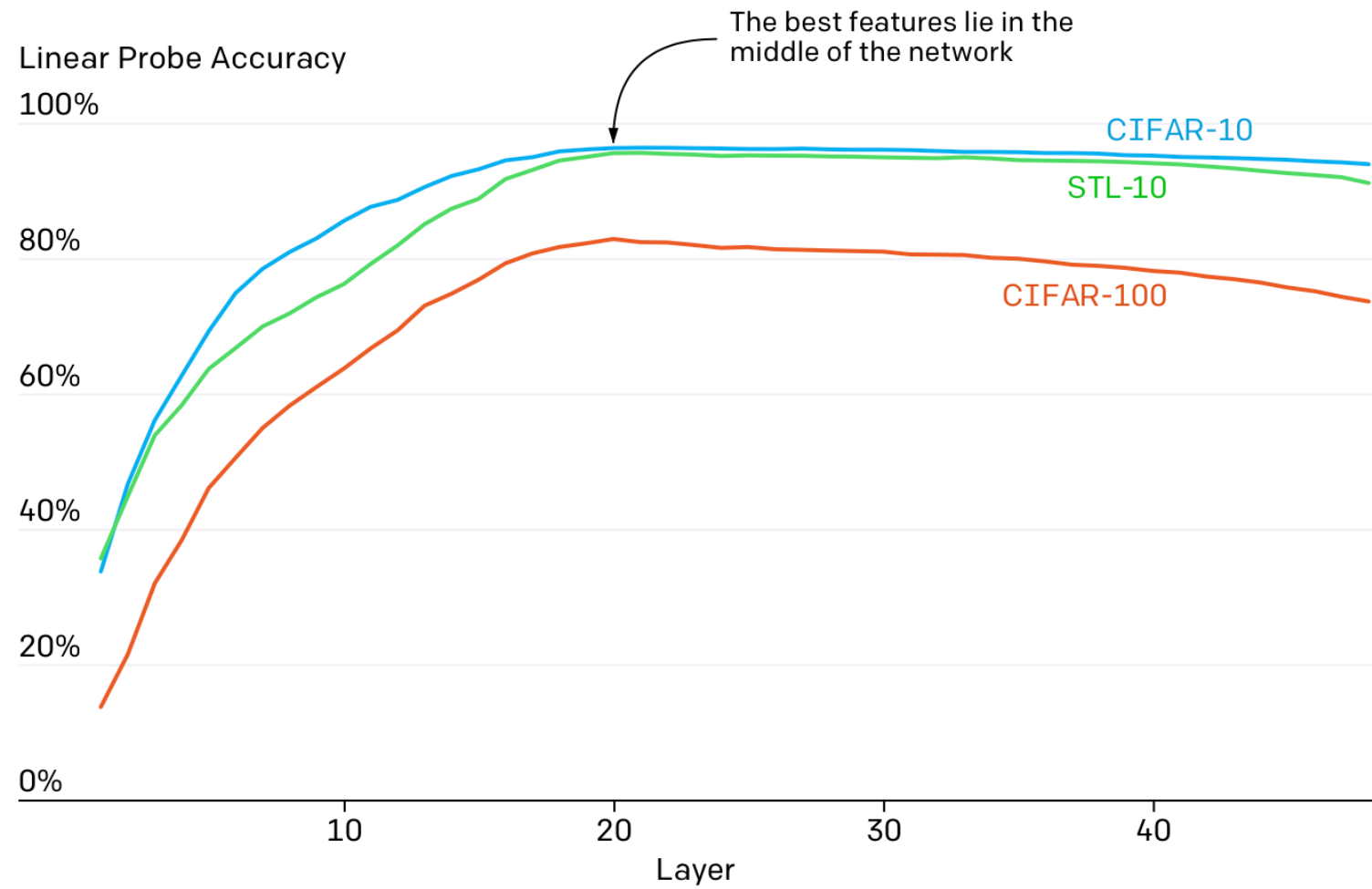
Data and models

- iGPT-XL. $L = 60$, $d = 3072$, parameters: 6.8B
- iGPT-L. $L = 48$, $d = 1536$, parameters: 1.4B
- iGPT-M. $L = 36$, $d = 1024$, parameters: 455M
- iGPT-S. $L = 4$, $d = 512$, parameters: 76M

Datasets:

- CIFAR-10
- CIFAR-100
- STL-10

Results

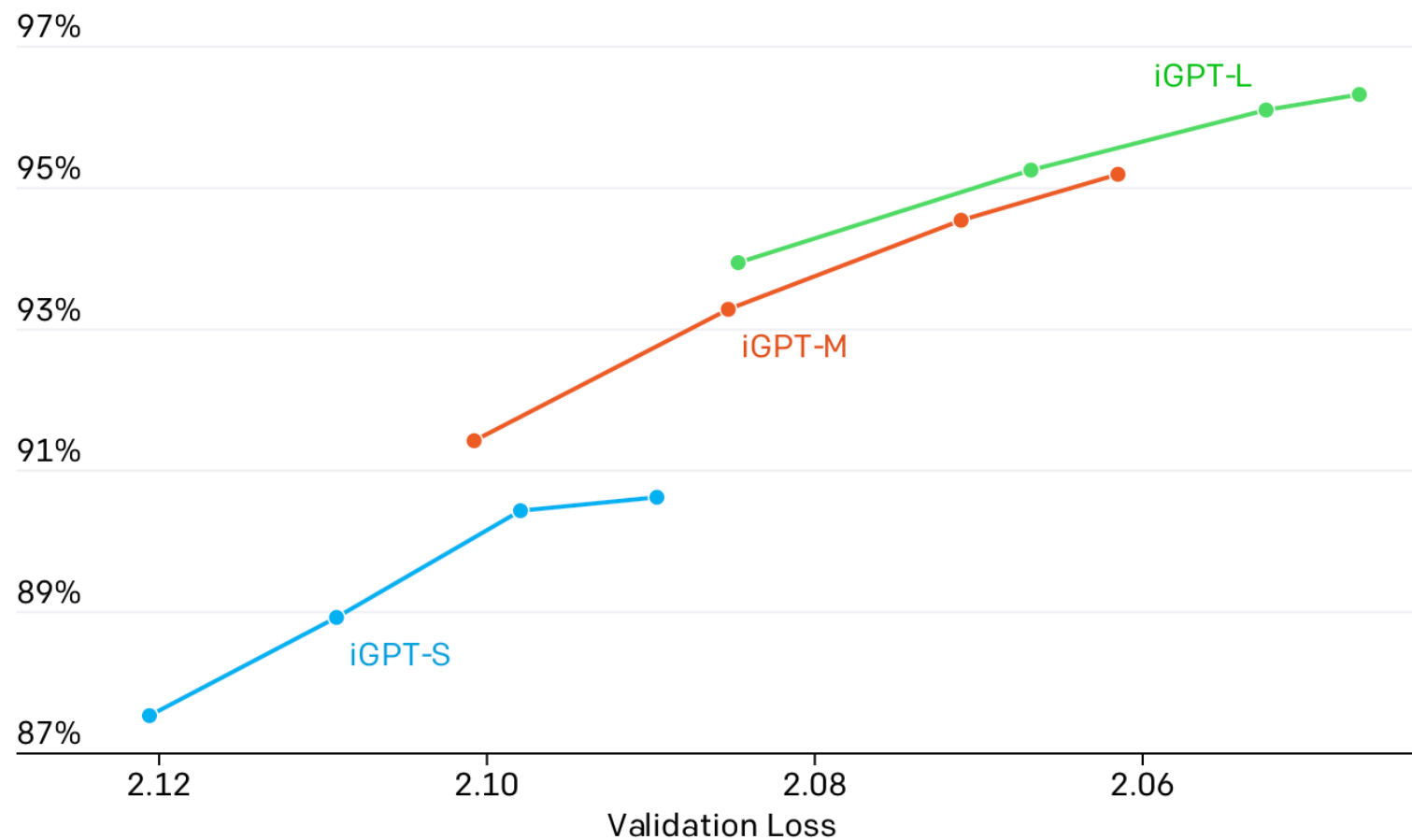


Results

Steps:

- 131k
- 262k
- 524k
- 1000k

CIFAR-10 Linear Probe Accuracy



Results

Качество моделей
предобученных на
ImageNet

EVALUATION	MODEL	ACCURACY	PRE-TRAINED ON IMAGENET	
			W/O LABELS	W/ LABELS
CIFAR-10 Linear Probe	ResNet-152 ⁵⁰	94.0		✓
	SimCLR ¹²	95.3	✓	
	iGPT-L 32x32	96.3	✓	
CIFAR-100 Linear Probe	ResNet-152	78.0		✓
	SimCLR	80.2	✓	
	iGPT-L 32x32	82.8	✓	
STL-10 Linear Probe	AMDIM-L ¹³	94.2	✓	
	iGPT-L 32x32	95.5	✓	
CIFAR-10 Fine-tune	AutoAugment ⁵¹	98.5		
	SimCLR	98.6	✓	
	GPipe ¹⁵	99.0		✓
	iGPT-L	99.0	✓	
CIFAR-100 Fine-tune	iGPT-L	88.5	✓	
	SimCLR	89.0	✓	
	AutoAugment	89.3		
	EfficientNet ⁵²	91.7		✓

Results

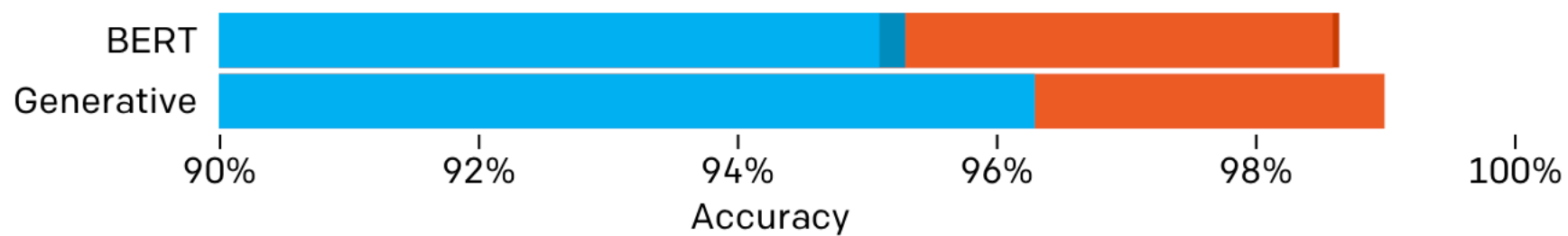
Качество на
ImageNet

METHOD	INPUT RESOLUTION	FEATURES	PARAMETERS	ACCURACY
Rotation ⁵³	original	8192	86M	55.4
iGPT-L	32x32	1536	1362M	60.3
BigBiGAN ³⁷	original	16384	86M	61.3
iGPT-L	48x48	1536	1362M	65.2
AMDIM ¹³	original	8192	626M	68.1
MoCo ²⁴	original	8192	375M	68.6
iGPT-XL	64x64	3072	6801M	68.7
SimCLR ¹²	original	2048	24M	69.3
CPC v2 ²⁵	original	4096	303M	71.5
iGPT-XL	64x64	3072 x 5	6801M	72.0
SimCLR	original	8192	375M	76.5

Results

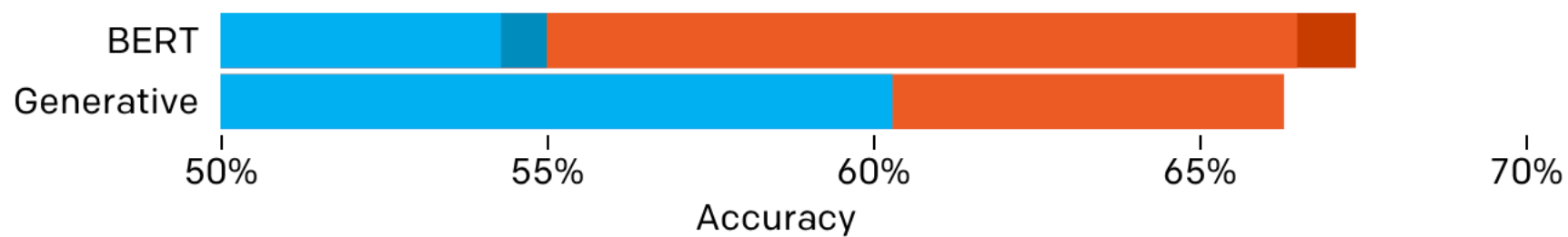
CIFAR-10

● Linear Probe ● Fine-tune



ImageNet

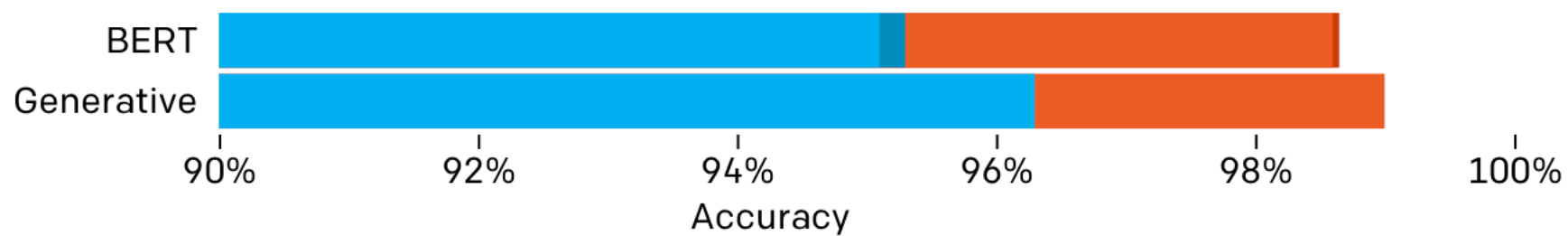
● Linear Probe ● Fine-tune



Results

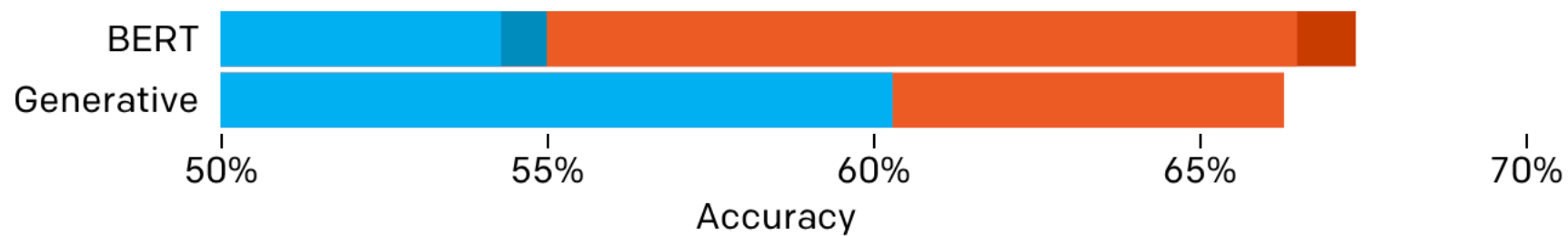
CIFAR-10

● Linear Probe ● Fine-tune



ImageNet

● Linear Probe ● Fine-tune



Results

MODEL	40 LABELS	250 LABELS	4000 LABELS
Improved GAN ⁵⁵	—	—	81.4 ± 2.3
Mean Teacher ⁵⁶	—	67.7 ± 2.3	90.8 ± 0.2
MixMatch ⁵⁷	52.5 ± 11.5	89.0 ± 0.9	93.6 ± 0.1
iGPT-L	73.2 ± 1.5	87.6 ± 0.6	94.3 ± 0.1
UDA ⁵⁸	71.0 ± 5.9	91.2 ± 1.1	95.1 ± 0.2
FixMatch ⁵⁹ RA	86.2 ± 3.4	94.9 ± 0.7	95.7 ± 0.1
FixMatch CTA	88.6 ± 3.4	94.9 ± 0.3	95.7 ± 0.2

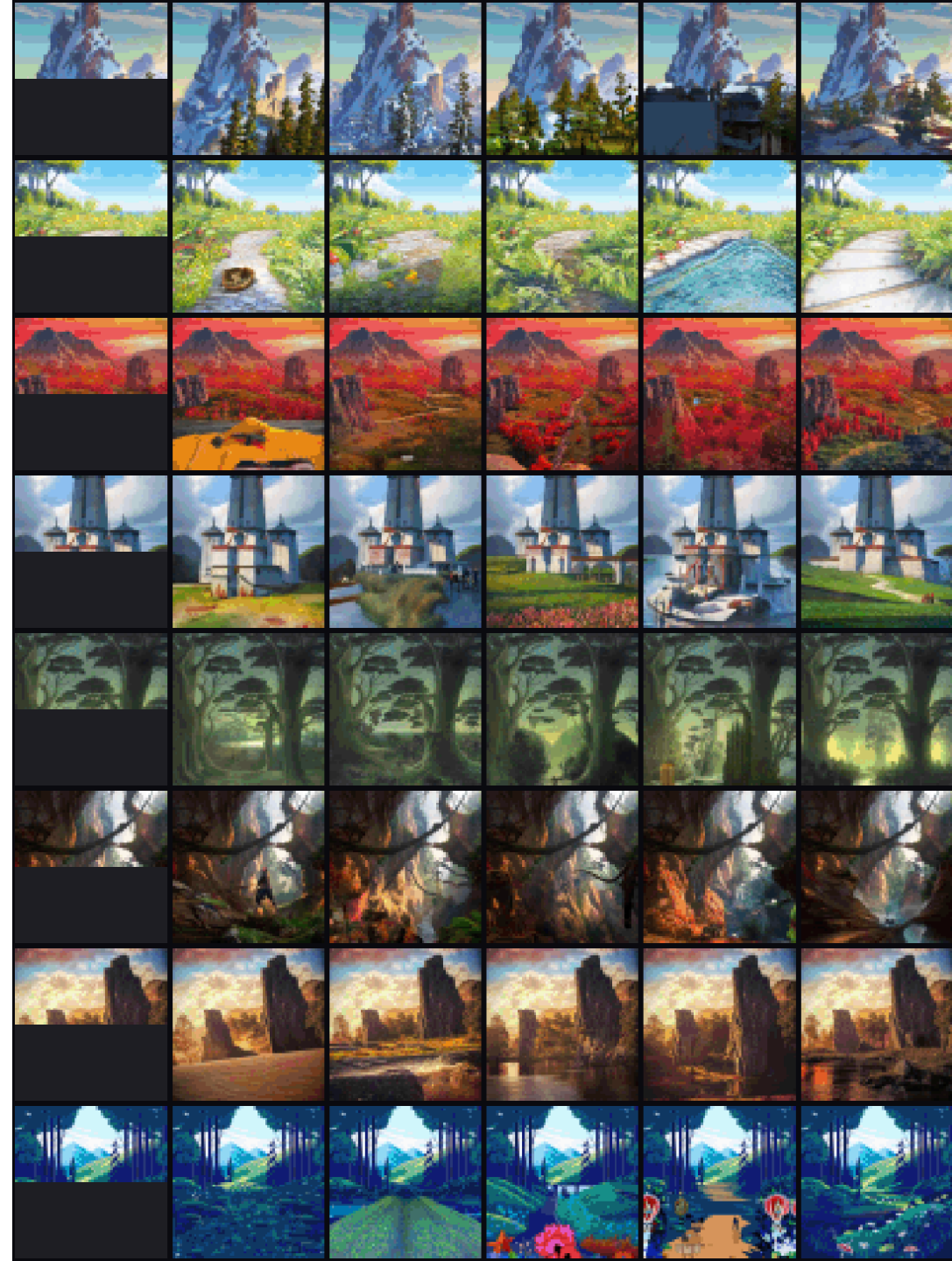
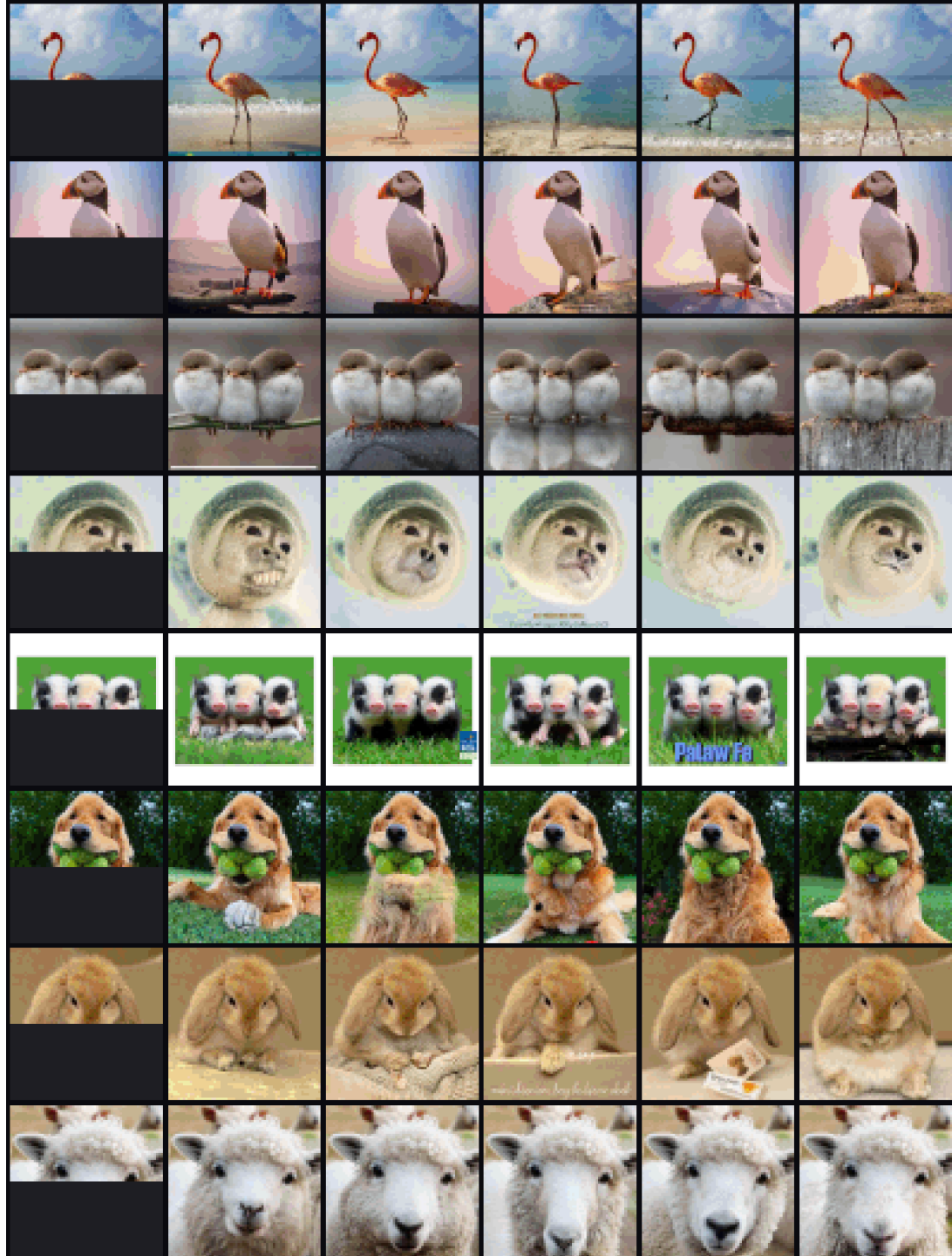
Results

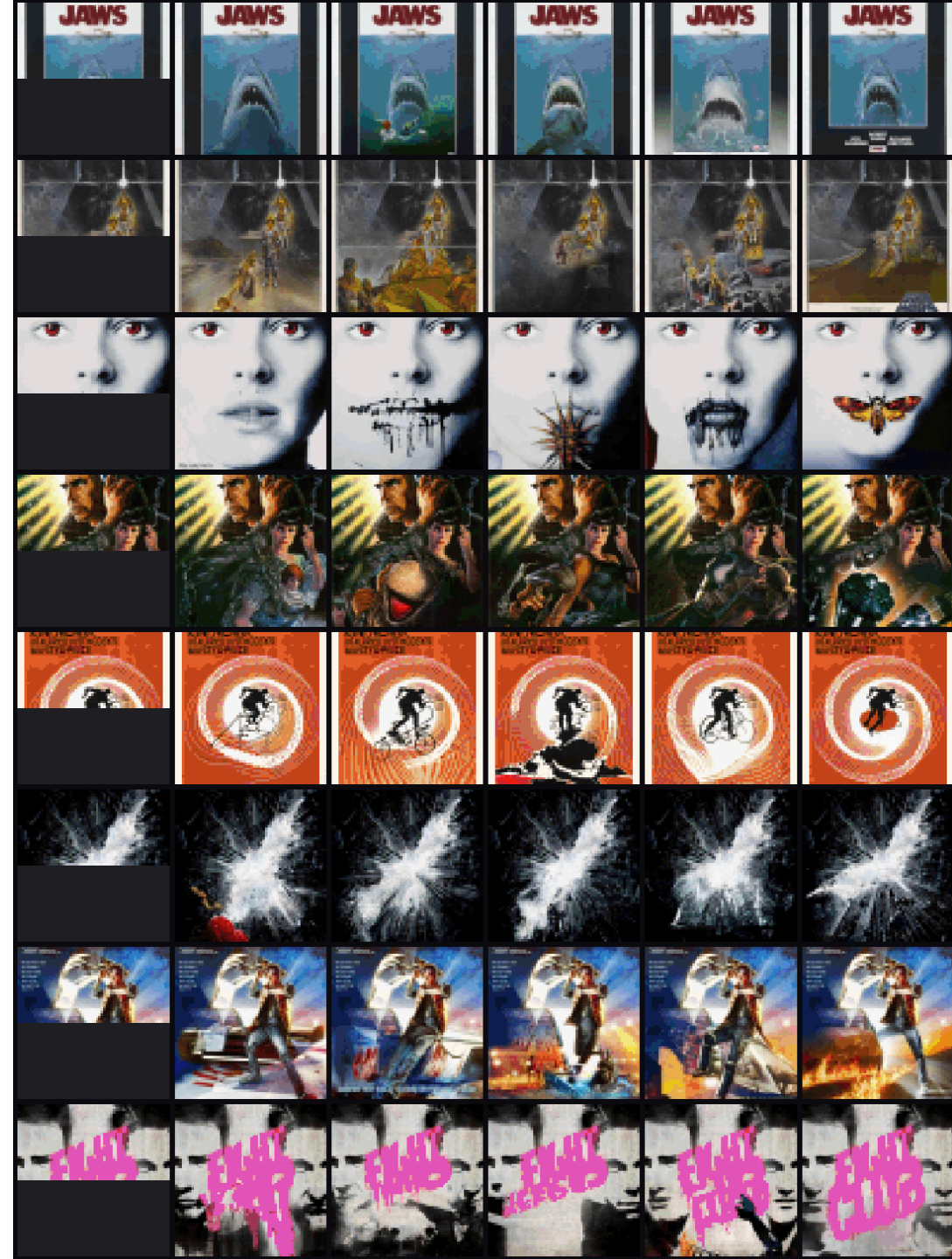
MODEL	40 LABELS	250 LABELS	4000 LABELS
Improved GAN ⁵⁵	—	—	81.4 ± 2.3
Mean Teacher ⁵⁶	—	67.7 ± 2.3	90.8 ± 0.2
MixMatch ⁵⁷	52.5 ± 11.5	89.0 ± 0.9	93.6 ± 0.1
iGPT-L	73.2 ± 1.5	87.6 ± 0.6	94.3 ± 0.1
UDA ⁵⁸	71.0 ± 5.9	91.2 ± 1.1	95.1 ± 0.2
FixMatch ⁵⁹ RA	86.2 ± 3.4	94.9 ± 0.7	95.7 ± 0.1
FixMatch CTA	88.6 ± 3.4	94.9 ± 0.3	95.7 ± 0.2

Conclusion

- Качество моделей в различных задачах не уступает сверточным сетям
 - Приходится использовать данные с более низким разрешением, чем сверточные сети
 - Большое количество параметров, долгое обучение
- > На данный момент достаточно непрактично для использования в реальных задачах







Links

- [Статья](#)
- [Блог пост](#)