

CLIP : Connecting Text and Images

Цыганов Артем, БПМИ-181

Мотивация

- Современные методы предсказывают только фиксированные категории
- Если надо усовершенствовать, прибегают к использованию, например, linear probe
- Но это не всегда удобно
- Хотим научиться получать информацию об изображениях на основе их описания
- А затем применять zero-shot transfer

Natural language supervision

- Проще масштабировать
- Связываем векторные представления для изображений с их текстовым описанием

Новый набор данных

Раньше:

- MS-COCO (100 000)
- Visual Genome (100 000)
- YFCC100M (100M)

STANFORD CARS

2012 Honda Accord Coupe (63.3%) Ranked 1 out of 196



✓ a photo of a 2012 honda accord coupe.

✗ a photo of a 2012 honda accord sedan.

✗ a photo of a 2012 acura tl sedan.

✗ a photo of a 2012 acura tsx sedan.

✗ a photo of a 2008 acura tl type-s.

KINETICS-700

country line dancing (99.0%) Ranked 1 out of 700



✓ a photo of country line dancing.

✗ a photo of square dancing.

✗ a photo of swing dancing.

✗ a photo of dancing charleston.

✗ a photo of salsa dancing.

SUN

kennel indoor (98.6%) Ranked 1 out of 723



✓ a photo of a kennel indoor.

✗ a photo of a kennel outdoor.

✗ a photo of a jail cell.

✗ a photo of a jail indoor.

✗ a photo of a veterinarians office.

FLOWERS-102

great masterwort (74.3%) Ranked 1 out of 102



✓ a photo of a great masterwort, a type of flower.

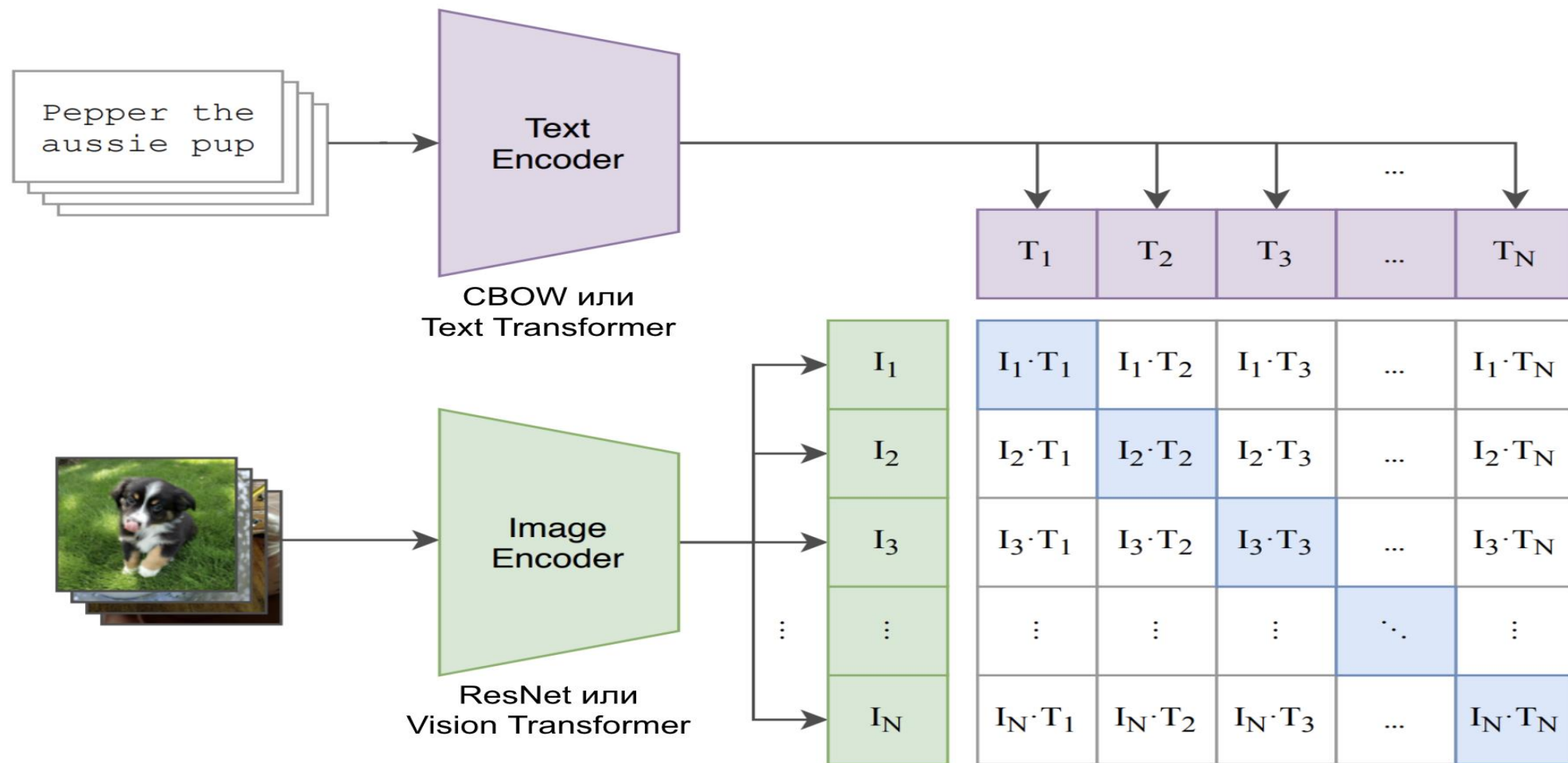
✗ a photo of a bishop of llandaff, a type of flower.

✗ a photo of a pincushion flower, a type of flower.

✗ a photo of a globe flower, a type of flower.

✗ a photo of a prince of wales feathers, a type of flower.

Архитектура



Предобучение

- Основной критерий – временные и вычислительные затраты
- Первоначальный подход работал слишком долго
- Упрощенная задача – как весь текст в целом подходит в качестве описания изображения

Предобучение

- Батч размера N
- Encoders
- Согласование размеров
- Матрица cosine similarity
- Cross-entropy loss

Cosine similarity between text and image features



Предобучение

```
# image_encoder - ResNet or Vision Transformer
# text_encoder - CBOW or Text Transformer
# I[n, h, w, c] - minibatch of aligned images
# T[n, l] - minibatch of aligned texts
# W_i[d_i, d_e] - learned proj of image to embed
# W_t[d_t, d_e] - learned proj of text to embed
# t - learned temperature parameter

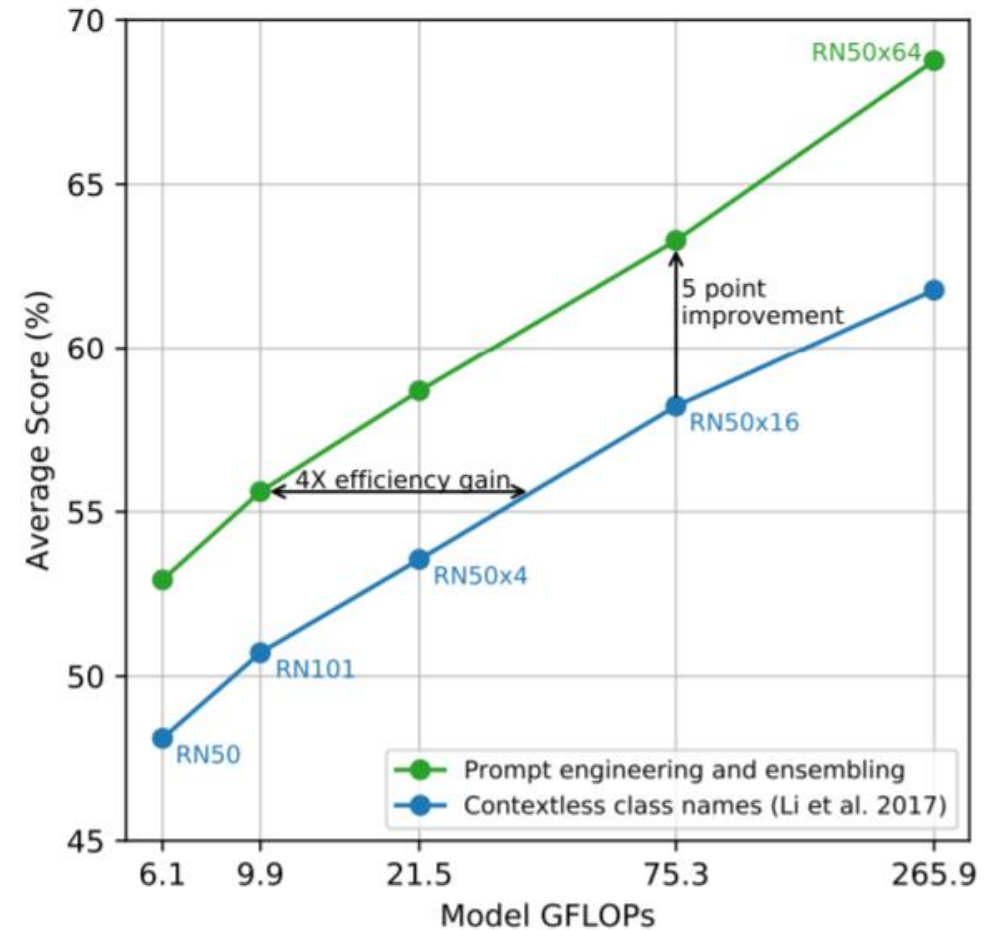
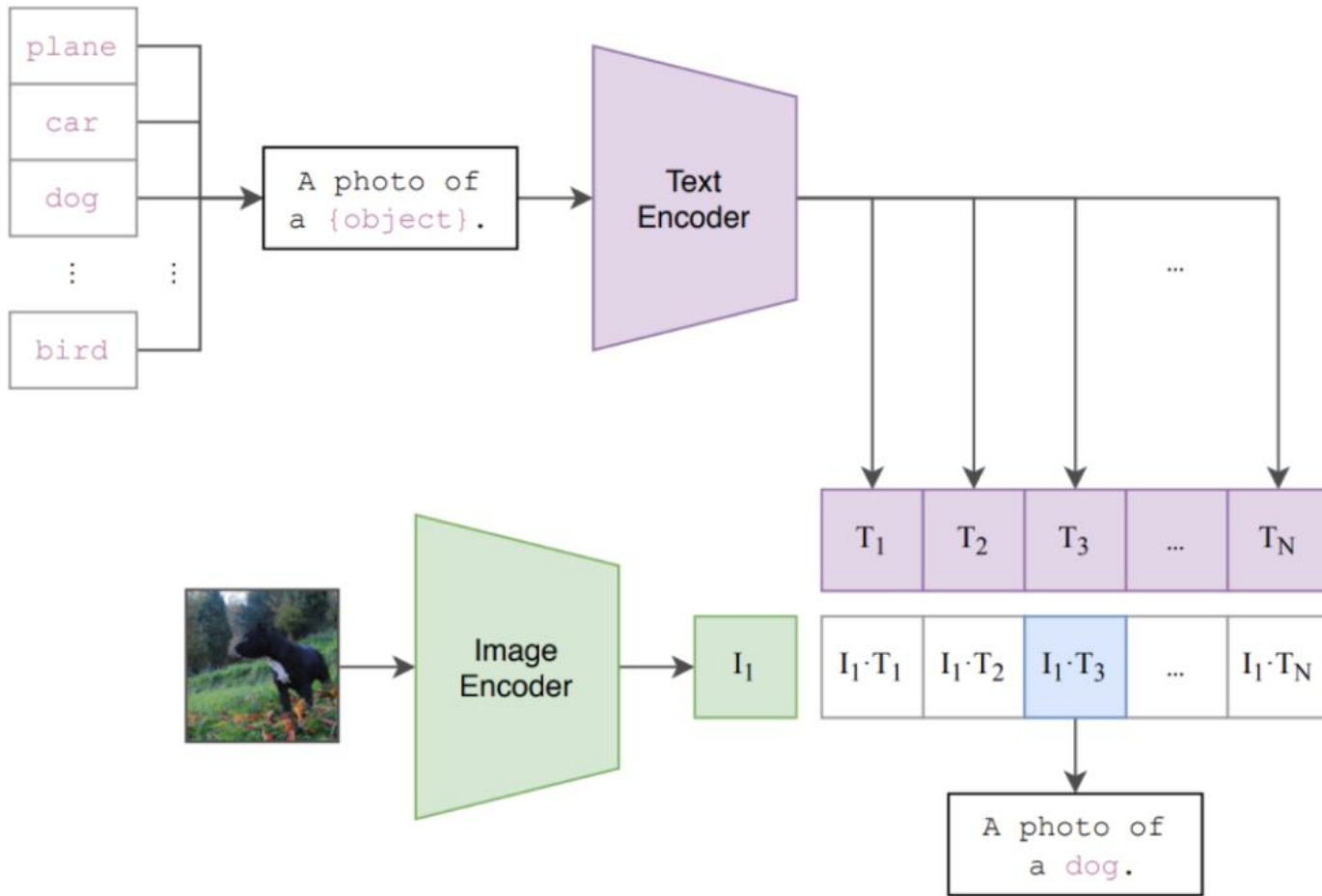
# extract feature representations of each modality
I_f = image_encoder(I) #[n, d_i]
T_f = text_encoder(T) #[n, d_t]
```

```
# joint multimodal embedding [n, d_e]
I_e = l2_normalize(np.dot(I_f, W_i), axis=1)
T_e = l2_normalize(np.dot(T_f, W_t), axis=1)

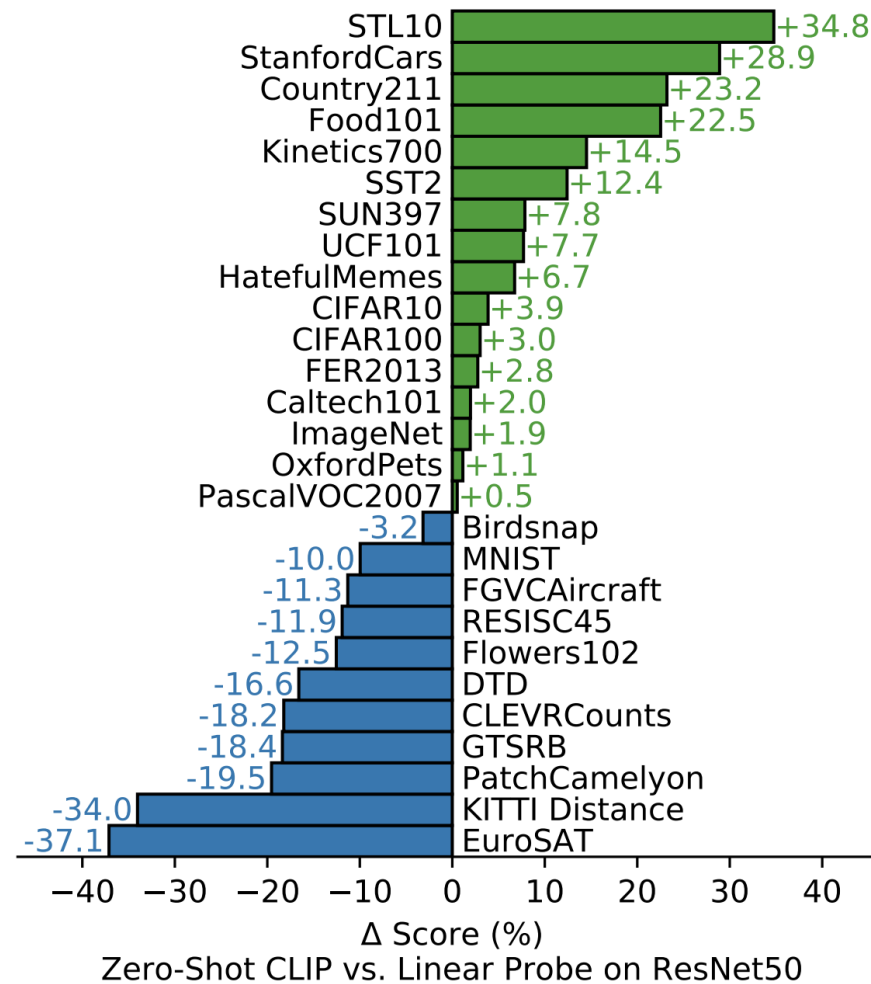
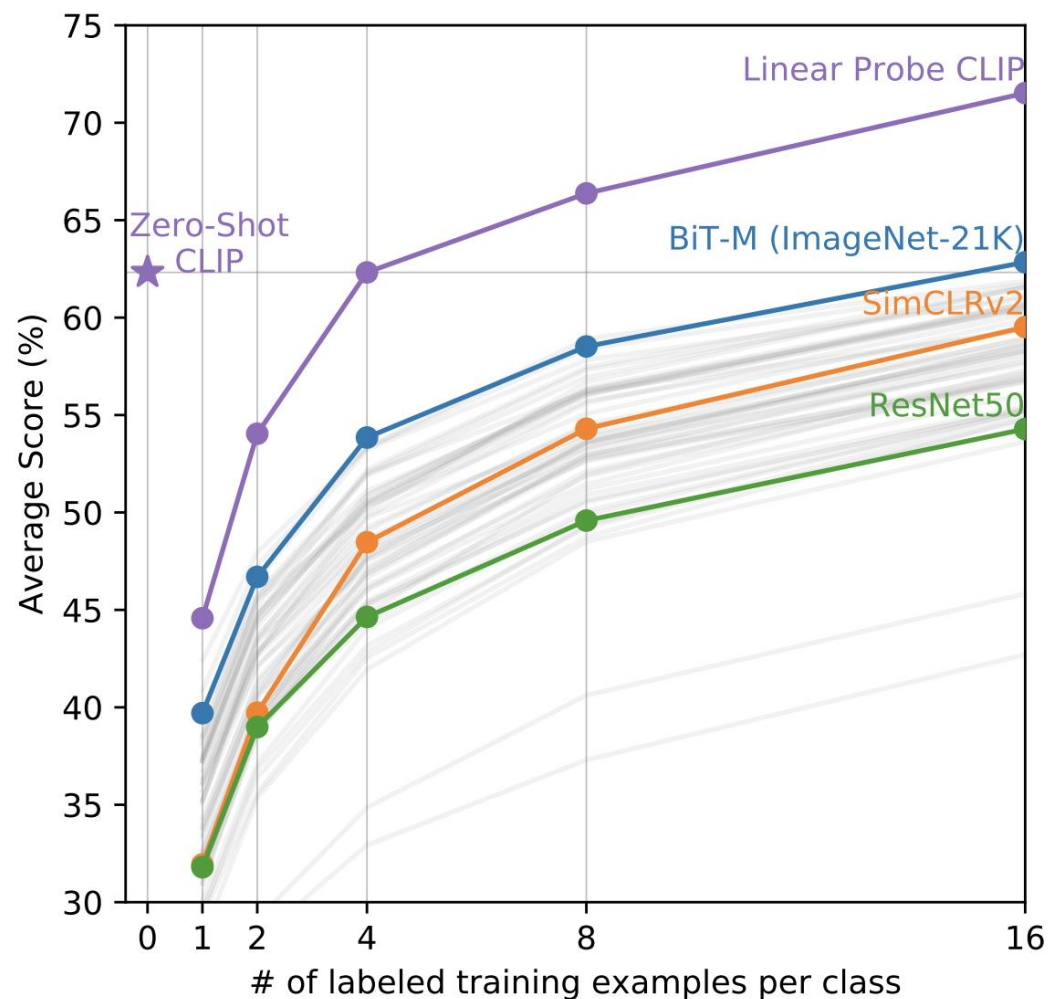
# scaled pairwise cosine similarities [n, n]
logits = np.dot(I_e, T_e.T) * np.exp(t)

# symmetric loss function
labels = np.arange(n)
loss_i = cross_entropy_loss(logits, labels, axis=0)
loss_t = cross_entropy_loss(logits, labels, axis=1)
loss = (loss_i + loss_t)/2
```

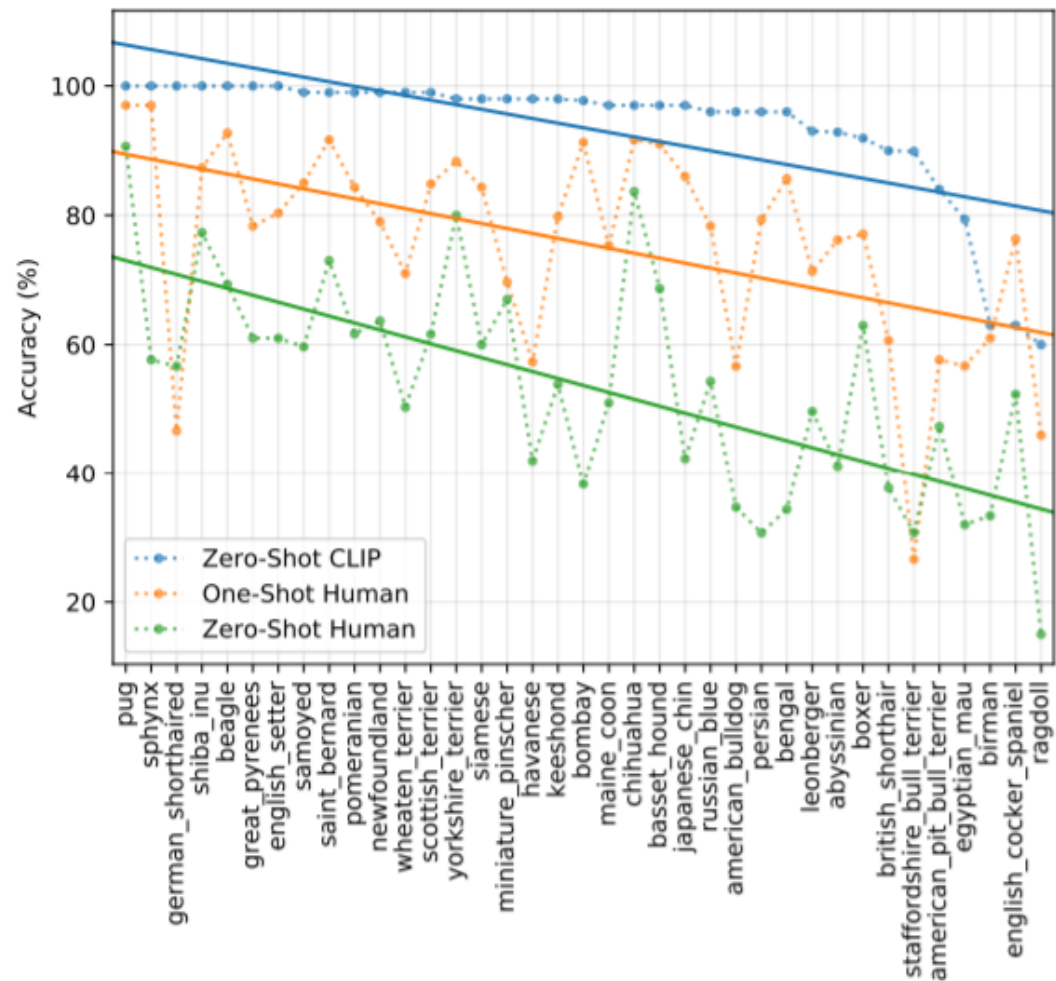

Transfer learning



Результаты



Недостатки



Granny Smith	85.6%
iPod	0.4%
library	0.0%
pizza	0.0%
toaster	0.0%
dough	0.1%



Granny Smith	0.1%
iPod	99.7%
library	0.0%
pizza	0.0%
toaster	0.0%
dough	0.0%

ИТОГИ

Плюсы

- Высокая производительность
- Широкое применение
- Обобщаемость и гибкость

Минусы

- Не использует уже полученные знания при повторном обучении
- Типографические атаки
- zero-shot не применим для узконаправленных задач

ИСТОЧНИКИ

- <https://cdn.openai.com/papers/Learning Transferable Visual Models From Natural Language Supervision.pdf>
- <https://habr.com/ru/post/539312/>
- <https://distill.pub/2021/multimodal-neurons/#person-neurons>
- <https://habr.com/ru/post/540312/>
- <https://openai.com/blog/clip/>
- <https://habr.com/en/post/537334/>