



Faculty of Computer Science

Научно-исследовательский
семинар

2022

CoCa: Contrastive Captioners are Image-Text Foundation Models

Review

Akulov Dmitry



Суть работы

- Может решать большой спектр задач
- Легко делать Fine-tuning и zero-shot learning

**Encoder
Classifier**

**Dual-Encoder
Contrastive**

**Encoder-Decoder
Captioner**

Публикация

- Залито на arXiv 04.05.2022
- Опубликовано в журнале Transactions on Machine Learning Research
- Цитирования: 63
 - Продолжения нет
 - Одно интересное упоминание в качестве ImageNet SOTA

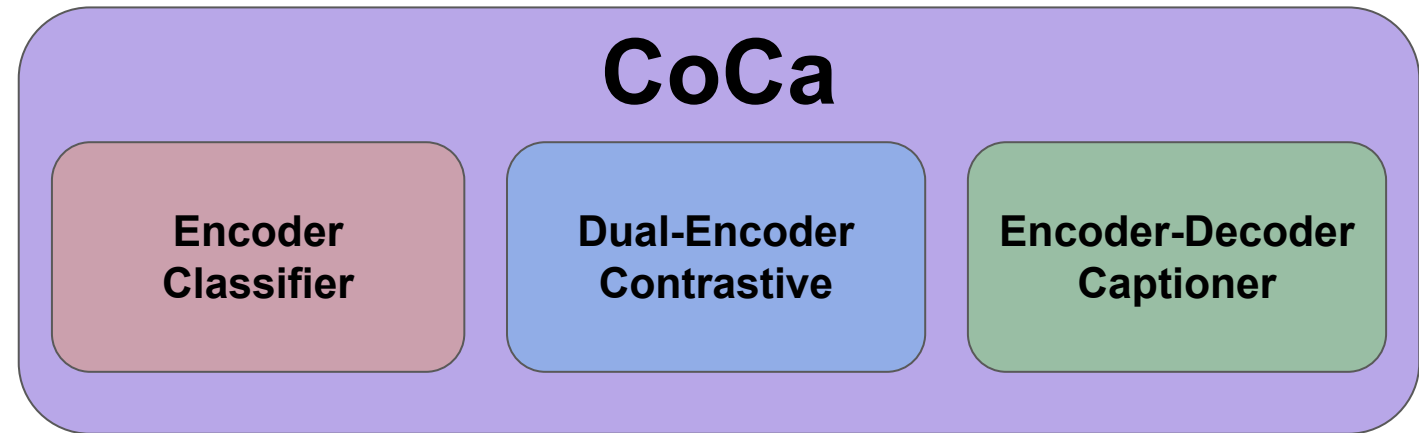


Суть работы

- Может решать большой спектр задач
- Легко делать Fine-tuning и zero-shot learning

Публикация

- Залито на arXiv 04.05.2022
- Опубликовано в журнале Transactions on Machine Learning Research
- Цитирования: 63
 - Продолжения нет
 - Одно интересное упоминание в качестве ImageNet SOTA





Авторы

С 2017 года
60 статей
8446 цитирований

CV

NLP

Audio

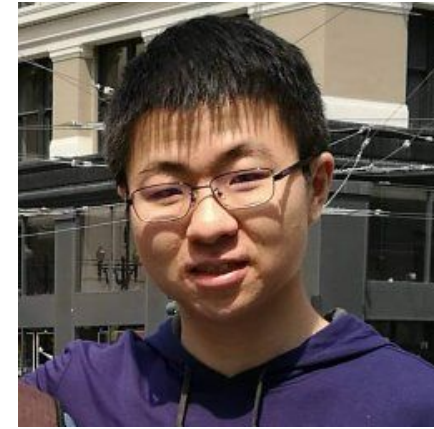


Jiahui Yu

Совместная статья (08.2021):

Encoder-Decoder
Captioner

Simvlm: Simple visual language model pretraining with weak supervision



Zirui Wang

С 2019 года
13 статей
775 цитирований

NLP

Opt

CV



Аналог

Florence: A New Foundation Model for Computer Vision (22.11.2021)



Аналог

Florence: A New Foundation Model for Computer Vision (22.11.2021)

Datasets

CoCa

- | | | |
|---|----|---|
| ImageNet
Kinetics | 1. | visual recognition |
| MSCOCO (ZH)
Flickr30K (ZH)
MSR-VTT (ZH) | 2. | crossmodal alignment |
| VQA v2
SNLI-VE
NLVR2
MSCOCO | 3. | image captioning
and
multimodal understanding |

Datasets

ImageNet
Kinetics
CIFAR10
Cross-Domain
Few-Shot learning
benchmark (FS)

MSCOCO (ZH+FT)
Flickr30K (ZH+FT)

VQA v2

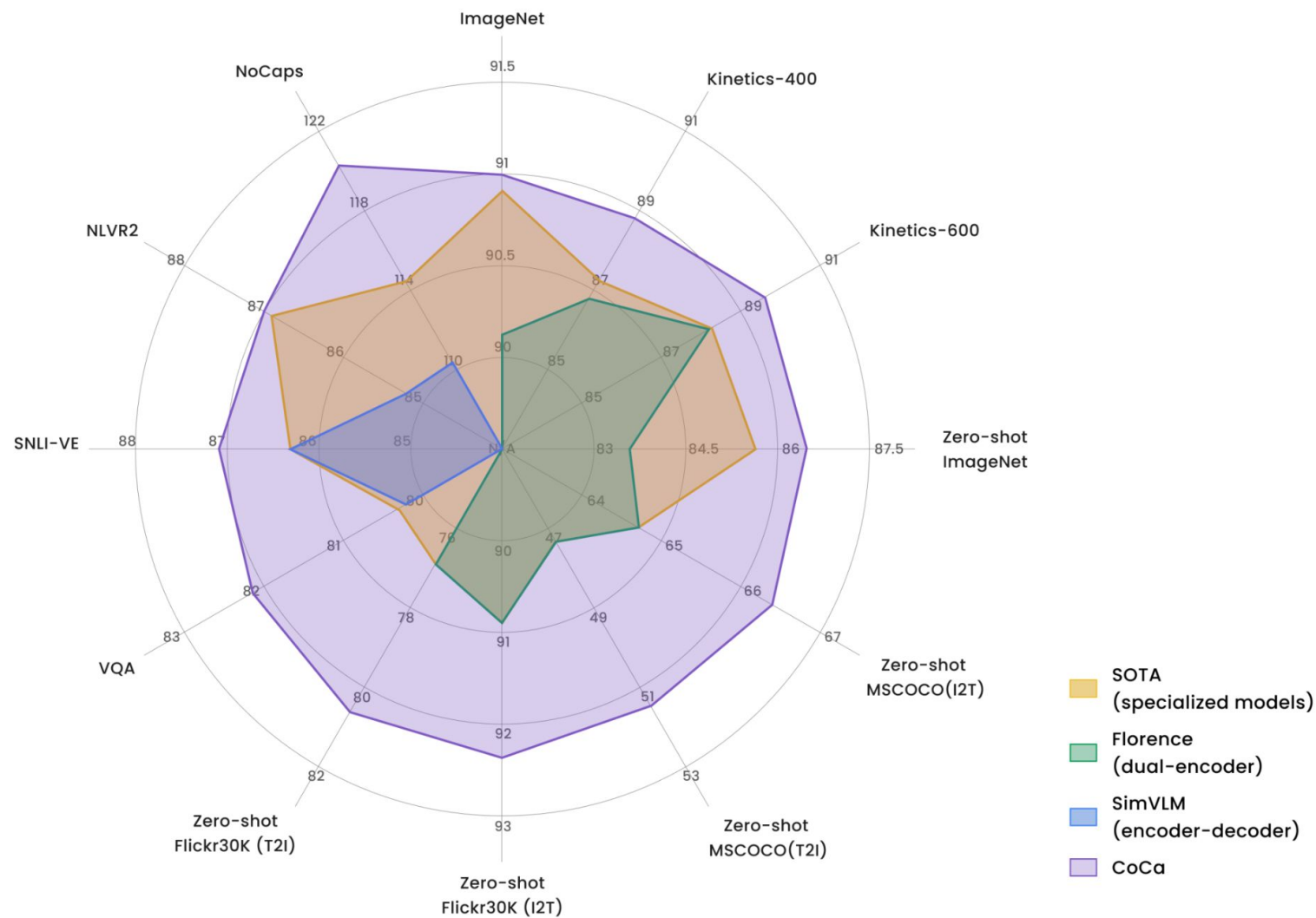
...

Florence

- | | |
|----|---|
| 1. | Classification |
| 2. | Retrieval |
| 3. | VL Representation Learning |
| 4. | Object Detection
Text-to-Video
Video Action Recognition |



Результат





Плюсы

- Хорошее техническое улучшение
- Сами подходы, которые были объединены, хорошо описаны



Минусы

- Были описаны задачи, которые уже покрываются теми моделями, что были объединены



Идеи по улучшению от авторов

- Может оказаться чувствительной к corrupted(~поврежденным) изображениям.
- Исследовать ширину используемости модели



Идеи по улучшению от меня

- Суть работы - объединение подходов. Как много еще можно наобъединять?
 - Что было бы полезно добавить именно к их модели?
 - Почему выбрали именно такие? Можно ли что-то выкинуть?



ИСТОЧНИКИ

- CoCa: Contrastive Captioners are Image-Text Foundation Models
<https://arxiv.org/abs/2205.01917v2>
- SimVLM: Simple Visual Language Model Pretraining With Weak Supervision
<https://arxiv.org/abs/2108.10904>
- Florence: A New Foundation Model for Computer Vision <https://arxiv.org/abs/2111.11432>