

Semi-Supervised learning image-text retrieval

Pasit Tiwawongrut
Asian Institute of Technology
Klong Luang Pathumthani 12120, Thailand
Pasit.Tiwawongrut@ait.asia

Dr. Chaklam Silpasuwanchai
Asian Institute of Technology
Klong Luang Pathumthani 12120, Thailand
chaklam@ait.asia

Abstract

The ABSTRACT is to be in fully-justified italicized text, at the top of the left-hand column, below the author and affiliation information. Use the word “Abstract” as the title, in 12-point Times, boldface type, centered relative to the column, initially capitalized. The abstract is to be in 10-point, single-spaced type. Leave two blank lines after the Abstract, then begin the main text. Look at previous ICCV abstracts to get a feel for style and length.

1. Introduction

Vision-language retrieval is a crucial task to create a search related image and text depends on each modalities. With the growth in vision language pre-trained models, these models become foundation models for many downstream task [2]. CLIP [11] propose pre-trained vision-language model results in remarkable performance in many vision language tasks by training with web-scaled image-text pairs. However when adapted to speicalized domain such as medical images and remote sensing, the model still struggle to get alignment of these specialized image-text [10].

2. Related work

2.1. Vision-Language model

In the past few years, many works have shown the ability to utilize textual information with the image task by training with image text pair *e.g.* CLIP [11], UNITER [5], Blip [9, 8], BEiT [14] and CoCa [15]. We can roughly divide the vision language model architecture into two categories. First, vision and language encoder *e.g.* CLIP, CoCa, ALIGN, and mPlug. These model focus on maximize alignment of two encoders for vision and language encoding. By training with a large amount of the image-text pair dataset, the ALIGN model could make up for the noisy image description and surpass the model, which was trained with the benchmark dataset in the zero shot image classifica-

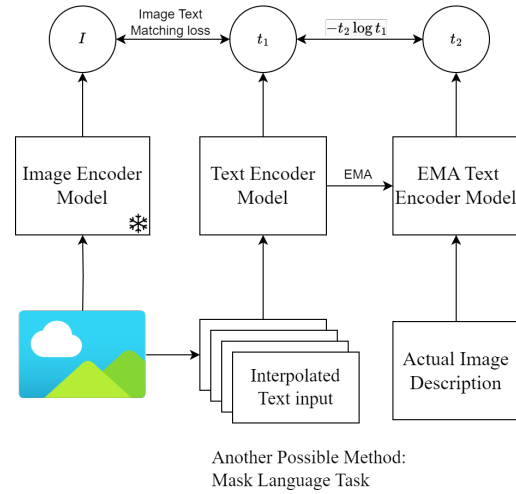


Figure 1. Overview of proposed method of applying moving average teacher to produce robust text encoder in pre-trained vision language model.

tion task. Recently **Contrastive Captioner** (CoCa) [15] proposed a vision-language encoder-decoder model which was trained with image-text contrastive loss and captioning loss. Cross attention layers were added to join image-text modality. Second methods are single encoder jointly trained with both modalities *e.g.* Uniter [5], BEiT-3 [14], and VLMO [1]. These methods concatenated both image and text embedding and utilize multi-head self-attention to joined vision and language modalities. In this research, we choose to experiment with the vision and language encoders method same as CLIP due to separable encoders for distillation.

2.2. Knowledge Distillation and Self-Distillation

Knowledge Distillation was firstly proposed by [7] to compress the model size while maintaining the model performance as much as possible. The method contained a smaller student model and a single or multiple larger teacher model. The knowledge was transferred by optimizing the student model output to match the teacher's out-

put. [6] investigated knowledge distillation using a student model size the same as the teacher model, showing improvement in the student model. Such a method is called self-distillation. The self-distillation has widely adopted in semi-supervised image classification tasks, such as Mean Teacher [13], EMAN [3] and FixMatch [12]. DINO [4] proposed self-distillation pre-training without using any label, which resulted in performance improvement. In this paper, we extended the self-distillation by creating representation which was image-text combined representation, and we trained the student model to match teacher softmax outputs.

3. Methodology

In this section we provided our self-distillation method and experiment setup details.

3.1. Self-Distillation

3.2. Evaluation

References

- [1] Hangbo Bao, Wenhui Wang, Li Dong, Qiang Liu, Owais Khan Mohammed, Kriti Aggarwal, Subhojit Som, Songhao Piao, and Furu Wei. Vlmo: Unified vision-language pre-training with mixture-of-modality-experts. *Advances in Neural Information Processing Systems*, 35:32897–32912, 2022.
- [2] Rishi Bommasani, Drew A Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, Michael S Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, et al. On the opportunities and risks of foundation models. *arXiv preprint arXiv:2108.07258*, 2021.
- [3] Zhaowei Cai, Avinash Ravichandran, Subhansu Maji, Charles Fowlkes, Zhuowen Tu, and Stefano Soatto. Exponential moving average normalization for self-supervised and semi-supervised learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 194–203, 2021.
- [4] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 9650–9660, 2021.
- [5] Yen-Chun Chen, Linjie Li, Licheng Yu, Ahmed El Kholy, Faisal Ahmed, Zhe Gan, Yu Cheng, and Jingjing Liu. Uniter: Universal image-text representation learning. In *Computer Vision – ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXX*, page 104–120, Berlin, Heidelberg, 2020. Springer-Verlag.
- [6] Tommaso Furlanello, Zachary Lipton, Michael Tschannen, Laurent Itti, and Anima Anandkumar. Born again neural networks. In *International Conference on Machine Learning*, pages 1607–1616. PMLR, 2018.
- [7] Geoffrey Hinton, Jeff Dean, and Oriol Vinyals. Distilling the knowledge in a neural network. pages 1–9, 03 2014.
- [8] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. *arXiv preprint arXiv:2301.12597*, 2023.
- [9] Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In *International Conference on Machine Learning*, pages 12888–12900. PMLR, 2022.
- [10] Sangwoo Mo, Minkyu Kim, Kyungmin Lee, and Jinwoo Shin. S-clip: Semi-supervised vision-language learning using few specialist captions. *Advances in Neural Information Processing Systems*, 36, 2024.
- [11] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021.
- [12] Kihyuk Sohn, David Berthelot, Nicholas Carlini, Zizhao Zhang, Han Zhang, Colin A Raffel, Ekin Dogus Cubuk, Alexey Kurakin, and Chun-Liang Li. Fixmatch: Simplifying semi-supervised learning with consistency and confidence. *Advances in neural information processing systems*, 33:596–608, 2020.
- [13] Antti Tarvainen and Harri Valpola. Mean teachers are better role models: Weight-averaged consistency targets improve semi-supervised deep learning results. *Advances in neural information processing systems*, 30, 2017.
- [14] Wenhui Wang, Hangbo Bao, Li Dong, Johan Bjorck, Zhiliang Peng, Qiang Liu, Kriti Aggarwal, Owais Khan Mohammed, Saksham Singhal, Subhojit Som, et al. Image as a foreign language: Beit pretraining for vision and vision-language tasks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 19175–19186, 2023.
- [15] Jiahui Yu, Zirui Wang, Vijay Vasudevan, Legg Yeung, Mojtaba Seyedhosseini, and Yonghui Wu. Coca: Contrastive captioners are image-text foundation models. *Transactions on Machine Learning Research*, 2022.