# CLIP-based image captioning via unsupervised cycle-consistency in the latent space

**Romain Bielawski**
ANITI
Université de Toulouse, France
rom.bielawski@gmail.com

**Rufin VanRullen**
CerCo
CNRS UMR5549, Toulouse
rufin.vanrullen@cnrs.fr

## Abstract

Image captioning typically involves an image encoder to extract meaningful image features, and a text decoder to generate appropriate sentences. Powerful pretrained models can be used for both image encoding and text decoding; but in this case, a separate multimodal translation stage between image-encoder output features and text-decoder input features must be learned. One exception is when image and text features are already aligned by construction, as in the CLIP model (Contrastive Language and Image Pretraining – a bimodal network pretrained on 400M image-text pairs). Pretrained CLIP-image features can be directly fed to a text-decoder trained to reconstruct captions from their pretrained CLIP-text features. Here we show that this direct captioning method is in fact sub-optimal. Instead, we propose an alternative method to translate CLIP-image features into CLIP-text features in a strictly unsupervised way, using the CycleGAN architecture – originally designed for unpaired image-to-image translation. Our Latent CycleGAN, optimized solely for an unsupervised cycle-consistency objective, generates CLIP-text latent features conditioned on CLIP-image latent features and vice-versa. Using these CLIP-text latent features as input to the text decoder, our method largely outperforms the direct captioning method that uses CLIP-image features – despite the fact that CLIP's large-scale pretraining should have already aligned the two feature spaces. This implies that cycle-consistency on unmatched multimodal data can be efficiently implemented in a bimodal latent space, and that CLIP-based image captioning can be improved without additional supervised training.

## 1 Introduction

Multimodality is gaining popularity due to the recently available online resources that make the creation of huge visio-linguistic datasets possible (Jia et al., 2021). Many models have been created to perform specific bimodal tasks such as Visual Question Answering or Image Captioning (Anderson et al., 2017; Lu et al., 2019; Li et al., 2019; Singh et al., 2019), but some have been designed with a more general objective: producing a multimodal latent vectorial space where images and text can be represented and compared. Among these models, CLIP – an algorithm trained with a multimodal contrastive objective on a large dataset (400M samples) of image-caption pairs – has shown impressive zero-shot learning abilities (Radford et al., 2021). This model has recently been tested on tasks for which it was not initially trained, such as transfer learning and few-shot learning on unimodal and multimodal datasets, or image captioning, establishing new SOTA results on some tasks (Bielawski et al., 2022; Mokady et al., 2021).

In the specific case of image captioning, many studies use pretrained models for image feature encoding as well as for text generation. An end-to-end image-to-caption fine-tuning stage is typically required, however, to align visual and linguistic representations in a supervised way on a matched image-caption dataset (Chen et al., 2021; Fang et al., 2021; Zhou et al., 2019). There is an obvious exception to this rule: when the pretraining of the model already aligned text and image features – as in the case of CLIP. Therefore, here we aim at leveraging this property by implementing a captioning pipeline that does not use matched image-caption data.

We first train a "CLIP-text decoder" to reconstruct captions based on their textual features in CLIP's latent space (a unimodal, linguistic objective); this text decoder is subsequently frozen. Hence, we compare a direct captioning pipeline – feeding the text-decoder with CLIP image features in order to generate a caption – with a pipeline where a CycleGAN (Zhu et al., 2017) inspired translator – trained with only unpaired visual and textual features – is used to convert image features

266

into text features before feeding them to the text-decoder. Even though CLIP's latent space was already pretrained with a brute-force approach to align its visual and linguistic representations on 400M image-caption pairs, we demonstrate that our feature conversion model trained using cycle-consistency in the CLIP latent space significantly improves captioning performance, compared with the direct method.

## 2 Dataset

To train our algorithms, we use the COCO (Lin et al., 2014) train 2014 dataset, composed of images representing complex scenes, along with their descriptions. We simply use the captions and the images independently, as two sets of unpaired uni-modal data from each modality.

For the evaluation, we use the COCO validation 2014 dataset.

## 3 Models

### 3.1 Pretrained models

We use CLIP ViT-B/32, a pretrained Vision-Transformer-based (Dosovitskiy et al., 2020) CLIP checkpoint, as image and text encoder. CLIP's vision encoder will be therafter referred to as just CLIP, and CLIP's text encoder will be called CLIP-T.
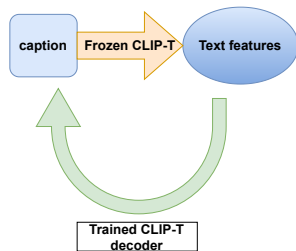


Figure 1: The text decoder is trained to reconstruct COCO train captions from their textual embedding in CLIP's latent space. It learns a mapping from CLIP-T features to prefixes that condition the generation of text with a pretrained (frozen) GPT-2. Note that the text-decoder is trained only with (unimodal) linguistic data.

In order to create our CLIP-T decoder (see Figure 1) we rely on the code provided by Mokady et al. (2021), inspired from Li and Liang (2021). Their decoder was originally trained on the CLIP image features of COCO images, with the objective to reconstruct their corresponding captions (therefore using paired vision-language data to align the

text decoder training with pretrained image features). Instead, our text decoder is trained in a unimodal setting on the CLIP-T textual features of captions from the COCO train set (414K captions), with the objective of regenerating the original text. This decoder uses GPT-2 (Radford et al., 2019) as a frozen language generator, and learns to produce prefixes that condition the generation of text. The parameters of the text decoder are shown in Table 1. Once trained, our text decoder is frozen and used as such in the two captioning pipelines that we compare.
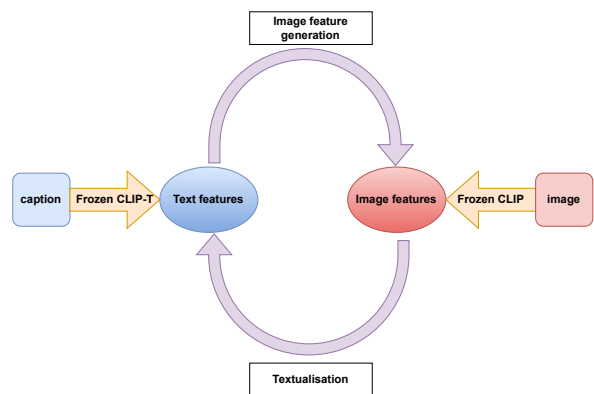
### 3.2 Architecture



Figure 2: The full architecture of the latent CycleGAN. The generators (purple arrows) are trained with un-matched multimodal data from the COCO dataset. One is trained to generate latent image features given a CLIP-T embedding, the other is trained to produce latent text features given the image features, i.e. to "textualize" them. Discriminators are not shown here.
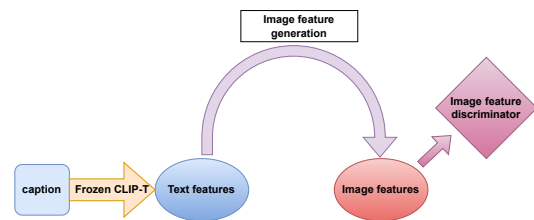


Figure 3: The GAN objective for Image feature generation: the generator must fool the discriminator, which must distinguish between real and fake inputs (here, between real image features and those translated from text features). A similar training objective and discriminator network exists for the other "textualisation" generator (not shown here).

The architecture and training procedure of the Latent CycleGAN are shown in Figure 2 to 4, the parameters of the architecture are displayed in Table 2. It is trained as a CycleGAN on unpaired data

| Parameter | # of epochs | prefix length | CLIP prefix length | mapping type | batch size | fine-tune GPT-2 | # of layers | CLIP version |
|---|---|---|---|---|---|---|---|---|
| Value | 20 | 40 | 40 | Transformer | 64 | False | 8 | ViT-B/32 |

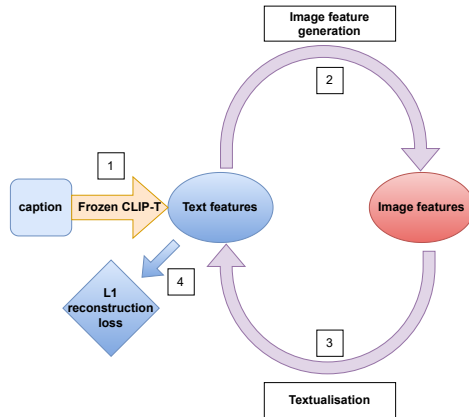Table 1: Parameters used for the training of the text decoder. For details see Mokady et al. (2021).



Figure 4: The cycle consistency objective consists in minimizing the L1 loss between a feature vector and its reconstruction after passing successively through both generators (here the translation of text features to and back from image features). The same cycle-consistency objective is also applied with cycles starting from the other (image) modality (not shown here).

| Parameter | # of epochs | batch size | # of layers | Latent space dimensions |
|---|---|---|---|---|
| Value | 20 | 64 | 8 | 512 |

Table 2: Parameters used for the Cycle-consistent architecture.

from the image and text modalities of the COCO train dataset (83K images and 414K captions). The training takes two generators – one of text features, one of image features – and two discriminators – to discriminate between real image (resp. text) features and fake/generated ones (Figure 2).

Just as in any GAN, the objective of each generator is to fool the corresponding discriminator. This is done by generating a fake latent vector in one modality, given a real latent vector from the other (this source vector can thus be considered as the noise that conditions the generation). The discriminator's objective is to guess whether any latent feature vector is real or generated (Figure 3). The generators of a CycleGAN (Zhu et al., 2017) also have specific extra objectives. The cycle consistency objective (Figure 4) minimizes the L1 loss between a feature vector and its reconstruction when passed successively through the two generators (e.g. an image feature vector is passed through the text feature generator, then this vector is passed though the image feature generator: the result of

this operation is a reconstructed image feature vector). The identity objective aims at learning the identity function when the image (resp. the text) generator is fed with an image (resp. text) feature vector.

Each generator is composed of 4 dense layers of dimension 512x512 with Tanh activation; the discriminators are composed of two dense layers, one of dimension 512x256, the other of 256x1, with LeakyReLU activation.

After having trained the Latent Cycle-GAN to convergence, we can then compare the two captioning pipelines illustrated in Figure 5.

The first one uses the fact that in CLIP's latent space, the features extracted from an image are intended to be as close as possible to the features computed by CLIP-T for a matching caption. This similarity was enforced by extensive contrastive training over 400M paired image-captions. Therefore, we may simply feed our CLIP-T decoder (trained on text features) with image features, and generate a corresponding caption.

The second pipeline uses the image-to-text-feature generator (the rest of the Latent CycleGAN was only required during training, i.e. to compute and optimize cycle-consistency). The image-to-text-feature generator is used for what we call here "textualisation", i.e. it generates a text feature vector conditioned on an input image feature vector. After the textualisation of an image vector, the textualised vector is fed to the CLIP-T decoder to generate the caption.

## 4 Task

Given an image from the COCO dataset, the model's task is to reconstruct one of the corresponding captions. Several scores can be used to evaluate the quality of the reconstruction. Here we dispay the BLEU-1 to BLEU-4 (Papineni et al., 2002) (BLEU-n counts matching n-grams in the model output to n-grams in the reference text), the ROUGE_L (Lin, 2004) (measuring the longest common subsequence between the model output and the reference), the CIDEr (Vedantam et al., 2014) (computing the average n-gram cosine similarity between the model output and several descrip-
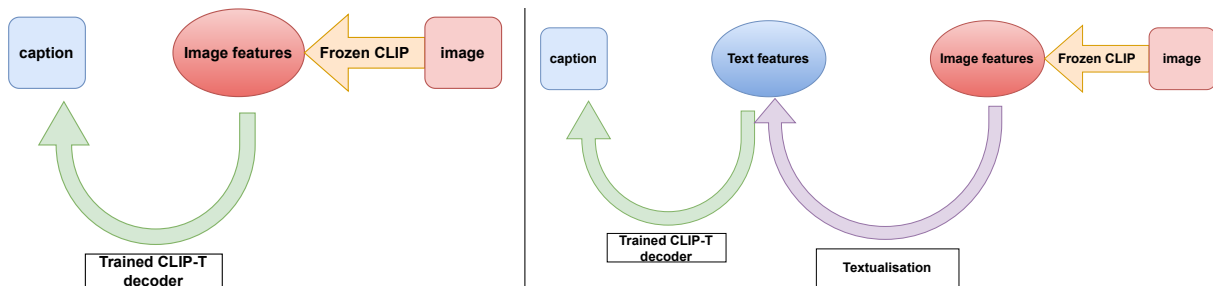
Figure 5: The two pipelines compared here for generating a caption. Our baseline (left) relies on the fact that CLIP was trained to project an image and its caption as close as possible in the latent multimodal space: the text decoder can thus generate a caption when given image features. The second one uses our generator, trained in an unsupervised way with unpaired multimodal data, to textualise the image features before feeding them to the text decoder.

| Scoring method | Image features as input | Textualised image features as input |
|---|---|---|
| BLEU-1 | 0.281 | 0.407 |
| BLEU-2 | 0.120 | 0.231 |
| BLEU-3 | 0.047 | 0.121 |
| BLEU-4 | 0.020 | 0.062 |
| ROUGE_L | 0.239 | 0.341 |
| METEOR | 0.098 | 0.161 |
| CIDEr | 0.064 | 0.247 |

Table 3: Scores for image captioning on the COCO validation set, for the two pipelines displayed in Figure 5. Higher scores indicate better captioning. The captioning pipeline with image features as input to the text decoder underperforms, compared to the one with features textualised using the text feature generator.

tions of reference and several n) and the METEOR (Denkowski and Lavie, 2014) (a variation of BLEU that aligns the reference and the output differently by incorporating semantic knowledge) scores.

## 5 Results

Results for the captioning task are displayed in Table 3. Despite the fact that CLIP's latent space was specifically designed and trained so that the encoding of an image and its description are as similar as possible, the strategy of directly using image latent features as input to the CLIP-T decoder does not perform well (for all scoring methods).

By simply training a Latent CycleGAN on unmatched COCO images and description (i.e. training in an unsupervised way on 82K images and their 500K descriptions, compared to the 400M image-text pairs of CLIP's initial training set) the improvement in score can go up to more than a factor 3.

Some uncurated examples of images and output captions from the COCO validation set can be seen

in Appendix A.

## 6 Discussion and conclusion

CLIP's bimodal alignment can allow image captioning at a SOTA level, but this requires a fine-tuning with paired image-caption data (Mokady et al., 2021). Since image and text are projected in the same latent space, it is also possible to use a direct captioning method with a trained CLIP-T decoder, without requiring any bimodal training; however, as we show here, this method is sub-optimal. We show how, using only unpaired images and captions, it is possible to significantly improve performance, while still taking advantage of CLIP's latent space multimodal alignment. Nonetheless, the results of the unpaired translation method implemented here remain far from the SOTA reached with supervised image captioning. Moreover, in our experiment, each caption implicitly matched an image from the training dataset, even though the matching was not given to the model. In future work, one might try enlarging the training domain of each modality, and incorporating data from separate, potentially larger unimodal datasets.

Our work suggests that the geometries of the representation of the vision and language modalities differ in CLIP's latent space. That is, CLIP's training has not properly brought together the two modalities. If it had, the image features would be directly usable by the text decoder, and our unsupervised "textualisation" training would not help the caption generation. This means that CLIP can represent vision and language in the same space, but vectors extracted from one domain are not fully multimodal in the sense that they are not indistinguishable from vectors from the other domain – in other words, modality-specific information appears

to interfere with full multimodality. Our Latent CycleGAN helps bridge the gap between the two latent representations, by enabling a unimodal text decoder to better understand image features, once they have been "textualised" by the text feature generator.

The recently proposed DALL-E2 (Ramesh et al., 2022) model, which uses a diffusion process to generate images from a caption, appears to have been based on a similar realization. Their diffusion image generator was trained to reconstruct an image given its CLIP image feature vector; however, for text-to-image generation, they did not directly feed the CLIP-T embedding into the diffusion generator, but first "translated" it into a suitable image-feature latent vector, exactly as we propose here.

## Acknowledgments

## References

Peter Anderson, Xiaodong He, Chris Buehler, Damien Teney, Mark Johnson, Stephen Gould, and Lei Zhang. 2017. Bottom-up and top-down attention for image captioning and VQA. *CoRR*, abs/1707.07998.

Romain Bielawski, Benjamin Devillers, Tim Van De Cruys, and Rufin Vanrullen. 2022. When does CLIP generalize better than unimodal models? when judging human-centric concepts. In *Proceedings of the 7th Workshop on Representation Learning for NLP*, pages 29–38, Dublin, Ireland. Association for Computational Linguistics.

Jun Chen, Han Guo, Kai Yi, Boyang Li, and Mohamed Elhoseiny. 2021. Visualgpt: Data-efficient image captioning by balancing visual input and linguistic knowledge from pretraining. *CoRR*, abs/2102.10407.

Michael Denkowski and Alon Lavie. 2014. Meteor universal: Language specific translation evaluation for any target language. In *Proceedings of the EACL 2014 Workshop on Statistical Machine Translation*.

Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. 2020. An image is worth 16x16 words: Transformers for image recognition at scale. *CoRR*, abs/2010.11929.

Zhiyuan Fang, Jianfeng Wang, Xiaowei Hu, Lin Liang, Zhe Gan, Lijuan Wang, Yezhou Yang, and Zicheng Liu. 2021. Injecting semantic concepts into end-to-end image captioning. *CoRR*, abs/2112.05230.

Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc V. Le, Yun-Hsuan Sung, Zhen Li, and Tom Duerig. 2021. Scaling up visual and vision-language representation learning with noisy text supervision. *CoRR*, abs/2102.05918.

Liunian Harold Li, Mark Yatskar, Da Yin, Cho-Jui Hsieh, and Kai-Wei Chang. 2019. Visualbert: A simple and performant baseline for vision and language. *CoRR*, abs/1908.03557.

Xiang Lisa Li and Percy Liang. 2021. Prefix-tuning: Optimizing continuous prompts for generation. *CoRR*, abs/2101.00190.

Chin-Yew Lin. 2004. ROUGE: A package for automatic evaluation of summaries. In *Text Summarization Branches Out*, pages 74–81, Barcelona, Spain. Association for Computational Linguistics.

Tsung-Yi Lin, Michael Maire, Serge J. Belongie, Lubomir D. Bourdev, Ross B. Girshick, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. 2014. Microsoft COCO: common objects in context. *CoRR*, abs/1405.0312.

Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. 2019. Vilbert: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks. *CoRR*, abs/1908.02265.

Ron Mokady, Amir Hertz, and Amit H. Bermano. 2021. Clipcap: CLIP prefix for image captioning. *CoRR*, abs/2111.09734.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation.

Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. 2021. Learning transferable visual models from natural language supervision. *CoRR*, abs/2103.00020.

Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners.

Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. 2022. Hierarchical text-conditional image generation with clip latents.

Amanpreet Singh, Vivek Natarajan, Meet Shah, Yu Jiang, Xinlei Chen, Dhruv Batra, Devi Parikh, and Marcus Rohrbach. 2019. Towards VQA models that can read. *CoRR*, abs/1904.08920.

Ramakrishna Vedantam, C. Lawrence Zitnick, and Devi Parikh. 2014. Cider: Consensus-based image description evaluation. *CoRR*, abs/1411.5726.

Luowei Zhou, Hamid Palangi, Lei Zhang, Houdong Hu, Jason J. Corso, and Jianfeng Gao. 2019. Unified vision-language pre-training for image captioning and VQA. *CoRR*, abs/1909.11059.

Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A. Efros. 2017. Unpaired image-to-image translation using cycle-consistent adversarial networks. *CoRR*, abs/1703.10593.

# A Appendix: Uncurated captioning examples



| Ground truth | Pair of commodes side by side in unfinished bathroom area. |
| --- | --- |
| | A torn apart bathroom with some toilets inside of it. |
| | A demolished bathroom with two toilets and a window |
| | The floor and wall of the bathroom are coming apart. |
| | A toilet and bidet sit in a bathroom that is under construction. |
| Direct method | Damaged CCTV image of restaurant staff posing |
| | as uncanny and uncanny people. |
| Textualised input | A view of a rough and dingy bathroom with many objects in it. |

| Ground truth | A group of people are standing on the sandy beach. |
|---|---|
| | Several people on the beach with their surf boards. |
| | Three men and three women posing on a beach in front of surf boards. |
| | A group of young people standing next to each other on a beach. |
| | A group of people pose for a picture near surfboards. |
| Direct method | A photograph of a young Irish kitty with a sun-dappled beach, |
| | and her friends at the bottom of the ocean. |
| Textualised input | The group of people posing and holding surfboards and a surf board. |

| Ground truth | A cheesy pizza with red peppers is in a box. |
|---|---|
| | A meal from japan or china on a tray. |
| | A cheesy casserole covered with toppings is depicted. |
| | A pizza with cheese and vegetables in a box. |
| | A large square shaped pizza covered in melted cheese and veggies. |
| Direct method | A quick chili sauce knife cut in the background, |
| | and green beans, muffin, and muffin |
| Textualised input | A bunch of cheese, ready to go and baked in a cheesy tortilla |

| Ground truth | A man stands beside his black and red motorcycle near a park. |
| --- | --- |
| | A man in black jacket next to a red motorcycle. |
| | An older man is standing beside a red motorcycle. |
| | A man standing by a motor cycle on a street. |
| | A man riding on the side of a red motorcycle. |
| Direct method | A Frank Miller Fun Road shot taken from the time I was born in 2006. |
| Textualised input | This person is showing on the road with some |
| | fresh motorcycle parts on the horizon. |