

# MultiModal Final Project Report

2000012918 Li Haoyang  
2000012993 Liu Haowen

January 14, 2023

## 1 Introduction to the Paper [3]

Training image captioning models typically requires large datasets of captioned images (image-text pairs), and these are challenging to collect. In order to solve this problem, the paper introduces a new image-captioning method that only uses the CLIP model and additional text data at training time, while no additional images are needed.

It relies on the fact that CLIP is trained to make visual and textual embeddings similar. Therefore, the only thing that needs to learn is how to translate CLIP embeddings back into text using a decoder. This method is named CapDec [3].

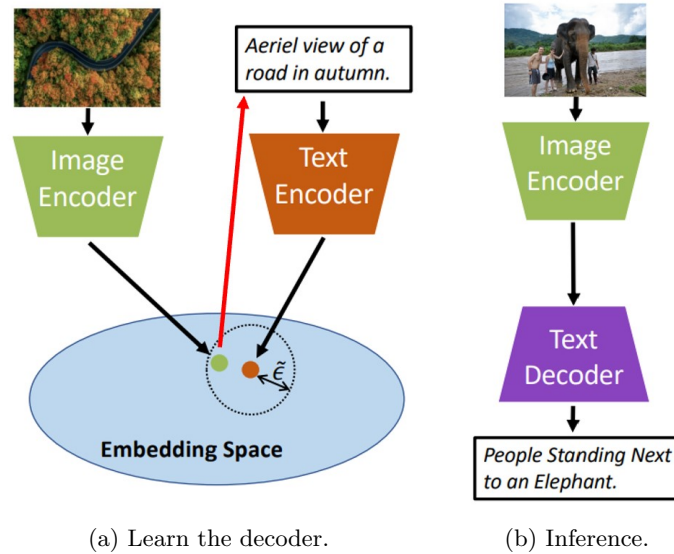


Figure 1: The red line in Fig.1a is the decoding process that we aim to learn. At inference, we use the CLIP image encoder to get the image embedding of the querying image, and feed it to the decoder (Text Decoder in Fig.1b). We can then get the corresponding text, namely, the caption.

In order to train the decoder while using only text, the paper introduces an intuitive way of fixing the gap between the text embeddings and image embeddings via noise injection. This has the effect of creating a ball in embedding space that will map to the same text, and corresponding image embedding is more likely to be inside this ball.

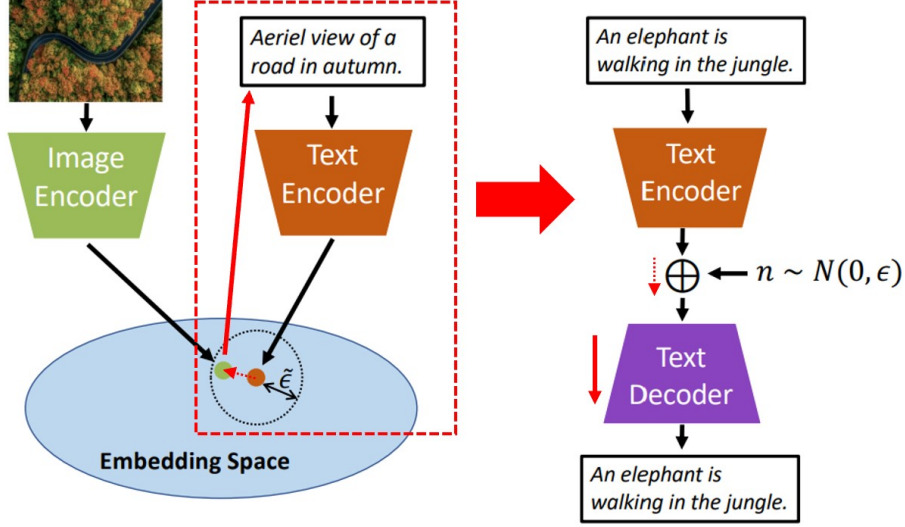


Figure 2: Embedded text is relatively close to its corresponding visual embedding, but with a certain gap. CapDec trains a model that decodes the CLIP embedding of text  $T$  back to text  $T$ , after noise-injection. The encoder remains frozen at training.

The paper demonstrate the effectiveness of the approach by showing SOTA zero-shot image captioning across four benchmarks, including style transfer. The code, data, and models are available at <https://github.com/DavidHuji/CapDec>.

Note that the backbone of the decoder is originate from ClipCap [2], which use embeddings as prefix to guide the language model, GPT-2.

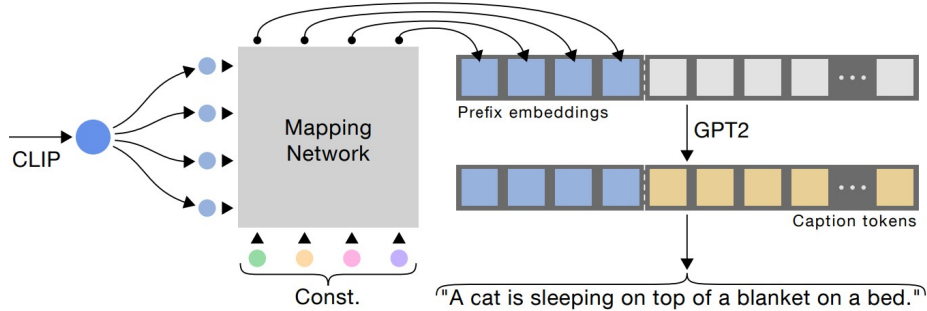


Figure 3: To extract a fixed length prefix, they trains a lightweight transformer-based mapping network from the CLIP embedding space and a learned constant to GPT-2.

## 2 Our Motivations

In the paper, the authors want all text embeddings in a small ball to decode to the same caption, which should also correspond to the visual content mapped to this ball. And they implement this intuition by adding zero-mean Gaussian noise of  $\epsilon$  to the text embedding before decoding it. Specifically, they set  $\epsilon$  to the mean  $\infty$  norm of embedding differences between five captions that correspond to the same image, estimating this based on captions of only 15 MS-COCO images.

Through this way, they calculate that  $\epsilon$  should be  $\sqrt{0.016}$ . And after some experiments of adjusting the  $\epsilon$ , they prove that  $\epsilon = \sqrt{0.016}$  is exactly a proper choice for CapDec.

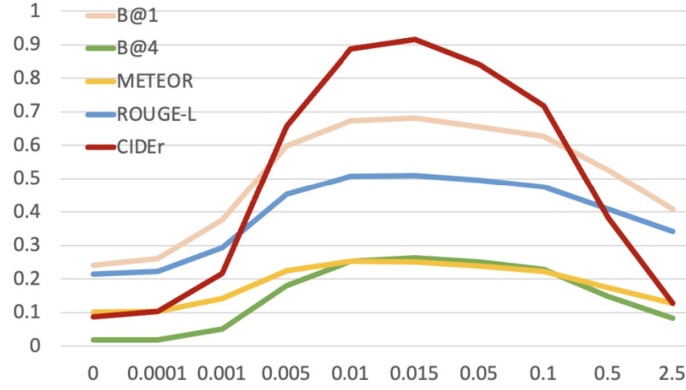


Figure 4: The effect of the noise variance ( $\epsilon^2$ ) on COCO performance. It is shown for analysis purposes only (not for  $\epsilon$  choosing) in their work.

However, their analysis to the noise is not persuasive enough to us. They believe the role of noise is just domain-gap correction. But From Fig.4, we could see that the chosen of the noise has a huge impact on the performance of the CapDec. Why it works so well by just injecting a  $N(0, 0.016)$  noise? Why they don't choose to use a uniform noise instead of a Gaussian noise? If it is possible to learn an adaptable noise if we embed the noise injection step into the decoder? What is the essence of the noise? We believe these problems are still remain researching.

Thus, our research is mainly centered on the noise injection step at training.

### 3 Evaluations

#### 3.1 Settings and Methods

- Python 3.9.15
- Torch 1.13.1
- Transformers 4.10.3
- Clip 1.0
- scikit-image 0.18.1
- pycocotools 2.0.6

We use images from COCO val set as input, and feed them to the pretrained CLIP image encoder. After obtaining the embeddings, we use them as prefixes of our trained decoder, and generate corresponding captions. We use the method from <https://github.com/tylin/coco-caption> to score the generated image-caption pairs.

#### 3.2 Metrics

Same as the original paper, we use 7 evaluation metrics: BLEU.1, BLEU.2, BLEU.3, BLEU.4, METEOR, ROUGE-L and CIDEr.

- BLEU is the metrics applied to machine translation before, and it is evaluated by judging the common parts between two sentences.

- METEOR matches between words and calculates the F value, and also considers the order of words to get the score.
- ROUGE\_L calculates the matching degree between two sentences through the longest common subsequence.
- CIDEr is a combination of BLEU and vector space model. It regards each sentence as a document, and then calculates the TF-IDF vector, which can reflect the weight of some words.

## 4 Reproduction of the Paper’s Work

we train a model locally on COCO dataset with default settings using zero-mean Gaussian noise of  $\epsilon = \sqrt{0.016}$ . There is 10 epochs in total, and for each epoch, it cost nearly 2 hour to train on a TITAN RTX GPU.

Epoch	1	2	3	4	5	6	7	8	9	10
Loss	1.75	1.20	1.05	0.96	0.89	0.84	0.79	0.76	0.73	0.70

Table 1: Average loss on training set of each epoch.

Tab.1 shows the changing of the loss while training the model above. We could see that the loss on the training set isn’t convergent at last, but owing to the high expense of training, we have to limit the total epochs to 10.

Epoch	1	2	3	4	5	6	7	8	9	10
Bleu_1	0.30	<b>0.50</b>	0.48	0.47	0.46	0.46	0.45	0.33	0.31	0.30
Bleu_2	0.14	<b>0.29</b>	0.27	0.27	0.26	0.25	0.25	0.15	0.14	0.14
Bleu_3	0.06	<b>0.17</b>	0.15	0.15	0.14	0.14	0.13	0.07	0.06	0.06
Bleu_4	0.03	<b>0.09</b>	0.08	0.08	0.08	0.07	0.07	0.03	0.03	0.03
METEOR	0.12	<b>0.16</b>	0.17	0.17	0.17	0.17	0.16	0.12	0.12	0.12
ROUGE_L	0.25	<b>0.36</b>	0.35	0.35	0.35	0.34	0.34	0.27	0.25	0.25
CIDEr	0.15	<b>0.37</b>	0.35	0.32	0.31	0.30	0.28	0.16	0.15	0.17

Table 2: Metrics on validation set of each epoch.

From Tab.2, we could find out that the results at evaluation are not quite ideal. There is an apparent gap between our reproduction and the original work in the paper. In our results, though Tab.1 demonstrates that the loss is decreasing while training, it seems that  $Epoch = 2$  is the best while evaluating. There is no doubt that our model comes into overfitting. And We think that the rest of the training process is just overfitting the text encoder-decoder reconstruction, which has no benefits to the eventual image captioning process.

We scan the code again, and find out that it offers an option before noise injection, that is, normalize the CLIP text embeddings. Although this trick seems unreasonable because the gap between the CLIP image embeddings and the CLIP text embeddings is calculated in the original CLIP embedding space, it maybe can fix the overfitting problem when training and worth trying. We use this trick and train the model again using zero-mean Gaussian noise of  $\epsilon = \sqrt{0.016}$ .

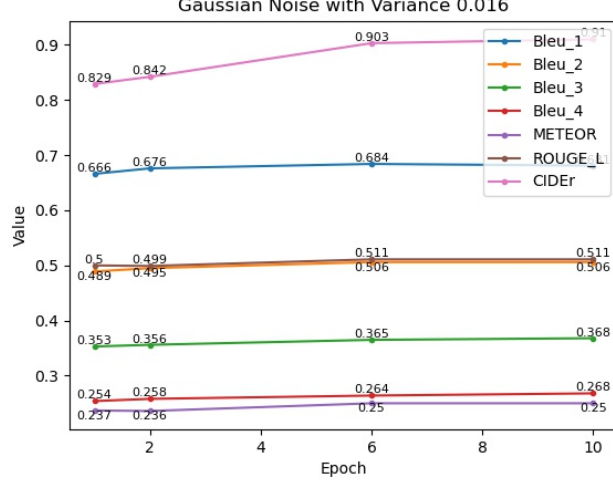


Figure 5: Use the trick that normalizes the CLIP text embeddings before inject zero-mean Gaussian noise of  $\epsilon = \sqrt{0.016}$ . The overfitting problem is fixed.

From Fig.5, we could see that this time, when  $Epoch = 10$ , the model reaches its highest performance. And all its metrics are familiar to the paper’s best model’s metrics.

Besides, there are pretrained weights on its official Github. We download the weights directly and test its performance on different metrics. Here are the results of our model using the normalization trick ( $Epoch = 10$ ) and the official pretrained model.

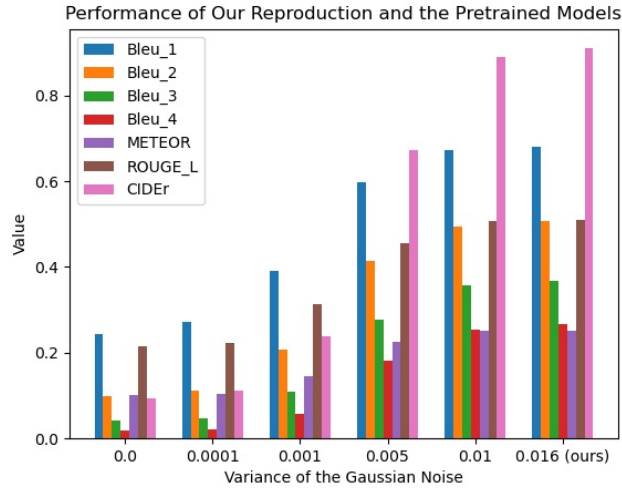


Figure 6: The official website doesn’t offer the pretrained weights of  $\epsilon^2 = 0.016$ . We choose some of the pretrained weights and evaluate. Among them, the pretrained weights of  $\epsilon^2 = 0.01$  has the similar performance with our reproduction, which is consistent with the previous Fig.4 in the paper.

To testify our reproduction of the original paper, we also draw a graph whose format is quite similar to Fig.4.

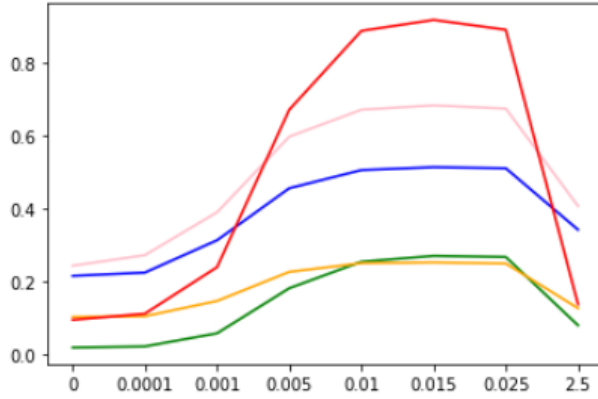


Figure 7: The re-present sketch of Fig.4, showing that our reproduction is successful.

## 5 Our Proposals and Implementations

**Note** that all of our experiments are based on the following settings:

- We use COCO training set to train and COCO validation set to evaluate, the version of COCO is 2014. And we parse the COCO dataset according to [karpathy-splits](#), after downloading the COCO dataset from its official website using the command "wget".
- We train the model on TITAN RTX GPU on a server, and set the epochs of training to 10. It cost 20-40 hours for each model to train, which depends on the specific proposal and implementation. And it cost nearly 3 hours to evaluate for each model.
- In order to guarantee the fairness at evaluation, all the models are trained without pretrained weights.
- Most of our models employ the trick mentioned above to avoid overfitting. If not, we will use the "w/o norm" tag to label these models in the following report.

### 5.1 Using Uniform Noise Instead of Gaussian Noise

When implementing the noise injection, we fetch the embeddings in a small ball with  $\epsilon$  radius around the text embeddings via uniform sampling, instead of Gaussian sampling. It means that every vector in the ball has the same probability to be chosen.

```
def get_uniform_ball_noise(input_shape, radius=0.1):
    uniform_noise_ball = torch.randn(input_shape, device=device) # normal distribution
    uniform_noise_sphere = torch.nn.functional.normalize(uniform_noise_ball, dim=1)
    u = torch.rand(input_shape[0], device=device) # unified distribution
    u = u ** (1. / input_shape[1])
    uniform_noise_ball = (uniform_noise_sphere.T * u * radius).T
    return uniform_noise_ball
```

Figure 8: The code of adding a uniform noise on the original text embeddings.

We adjust the radius of the uniform ball by changing the  $\epsilon$ . Here, we test the cases that  $\epsilon = \sqrt{0.001}$ ,  $\epsilon = \sqrt{0.016}$ , and  $\epsilon = \sqrt{0.1}$ , using Gaussian noise with  $\epsilon = \sqrt{0.016}$  as baseline. The results are as follows.

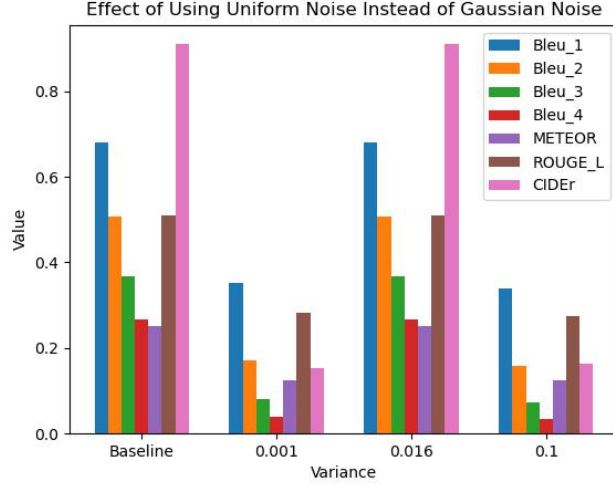
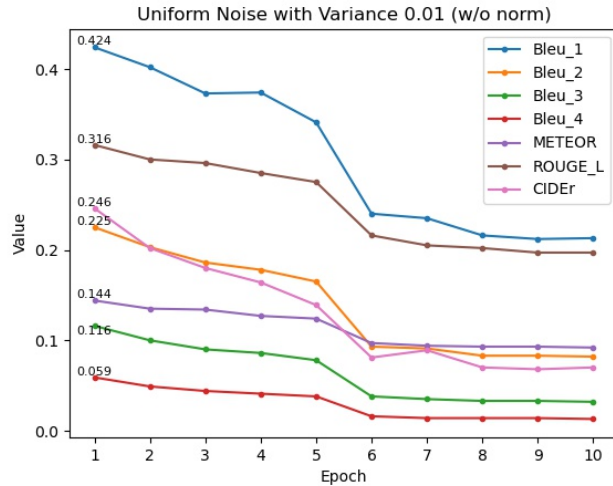
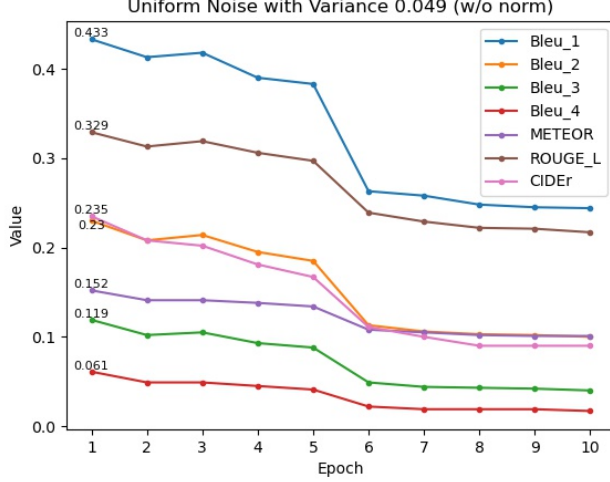


Figure 9: We notice that the uniform noise with  $\epsilon = \sqrt{0.016}$  reaches the same performance with the Gaussian noise with  $\epsilon = \sqrt{0.016}$ , while the uniform noise with  $\epsilon = \sqrt{0.001}$  and  $\epsilon = \sqrt{0.1}$  have poor performance.

We also train the "w/o norm" models using  $\epsilon = \sqrt{0.01}$  and  $\epsilon = \sqrt{0.049}$ . The results are shown in Fig.10a and Fig.10b. It turns out that for these models, the loss is declining with epochs at training but the best performance occurs on  $Epoch = 1$ , meaning that they all meet the overfitting problem.



(a) Inject  $U(0, \sqrt{0.01})$  noise (w/o norm).



(b) Inject  $U(0, \sqrt{0.049})$  noise (w/o norm).

Figure 10: We inject  $U(0, \sqrt{0.01})$  and  $U(0, \sqrt{0.049})$  noise without normalize the CLIP text embeddings before. The overfitting problem is even more severe than the  $N(0, \sqrt{0.016})$  in Tab.2

## 5.2 Use Trainable Modality Offset

In the paper, the authors try to add the modality offset based on the noise injection. In other words, by calculating the shift between the mean of text embeddings and the mean of image embeddings in COCO, we can then add that shift to the image embeddings to “correct” for the gap, and apply the CapDec to the resulting embedding.

However, in the paper, it turns out that adding the shift (modality offset) before noise injection, namely, injecting a  $N(shift, \epsilon^2)$  noise rather than a  $N(0, \epsilon^2)$  noise, can’t enhance the performance in general. Still, the simple  $N(0, 0.016)$  noise is the best for COCO.

This result seems to be strange. Given the the law of CLT (Central Limit Theorem) in probability statistics, which states that if we take sufficiently large random samples from a statistic then the distribution of the sample means will be approximately normally distributed, the gap between the image embeddings and text embeddings should be a statistic that follows a normal distribution  $N(\mu, \epsilon^2)$ . There is no reason that the mean of the distribution  $\mu$  should be 0. Therefore, we think that maybe by learning a  $\mu$ , rather than setting a constant shift or just setting it to 0, could we improve the performance.



```

class ClipCaptionModel_embedded_with_noise(nn.Module):

    def get_dummy_token(self, batch_size: int, device: torch.device) -> torch.Tensor:
        return torch.zeros(batch_size, self.prefix_length, dtype=torch.int64, device=device)

    def forward(self, tokens: torch.Tensor, prefix: torch.Tensor, mask: Optional[torch.Tensor] = None,
                labels: Optional[torch.Tensor] = None):
        prefix = prefix + (torch.randn(prefix.shape, device=device) * 0.4) + self.modality_offset
        prefix = torch.nn.functional.normalize(prefix, dim=1)
        embedding_text = self.gpt.transformer.wte(tokens)
        prefix_projections = self.clip_project(prefix).view(-1, self.prefix_length, self.gpt_embedding_size)
        embedding_cat = torch.cat((prefix_projections, embedding_text), dim=1)
        if labels is not None:
            dummy_token = self.get_dummy_token(tokens.shape[0], tokens.device)
            labels = torch.cat((dummy_token, tokens), dim=1)
        out = self.gpt(inputs_embeds=embedding_cat, labels=labels, attention_mask=mask)
        return out

    def __init__(self, prefix_length: int, clip_length: Optional[int] = None, prefix_size: int = 512,
                num_layers: int = 8, mapping_type: MappingType = MappingType.MLP):
        super(ClipCaptionModel_embedded_with_noise, self).__init__()
        self.modality_offset = nn.Parameter(torch.zeros([1, prefix_size], requires_grad=True))
        self.prefix_length = prefix_length
        self.gpt = GPT2LMHeadModel.from_pretrained('gpt2')
        self.gpt_embedding_size = self.gpt.transformer.wte.weight.shape[1]
        if mapping_type == MappingType.MLP:
            self.clip_project = MLP((prefix_size, (self.gpt_embedding_size * prefix_length) // 2,
                                                self.gpt_embedding_size * prefix_length))
        else:
            self.clip_project = TransformerMapper(prefix_size, self.gpt_embedding_size, prefix_length,
                                                clip_length, num_layers)

```

Figure 11: The code of the new class which basically derives from ClipCaptionModel class. In this class, We embed the noise injection step in its forwarding, and setting a new learnable parameter "self.modality\_offset", which represents the mean  $\mu$  of the normal noise. To make the training process more controllable, the std of the normal noise is set to  $\epsilon = \sqrt{0.016}$ , a constant, rather than a learnable parameter as well.

Here are the results.

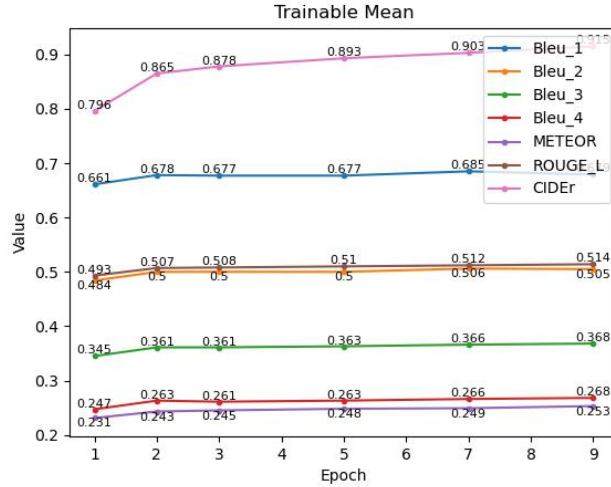


Figure 12: When  $Epoch = 10$ , it reaches the highest performance.

From Fig.12, we could see its performance is slightly better than the baseline refers to Fig.5. And we examine the L2-norm of the trainable mean after each epoch, finding that the mean is still close to 0 after a long period of time training the model, as shown in Fig.13. Therefore, we could draw a conclusion that the best performance of  $N(shift, 0.016)$  happens when  $shift \approx 0$ .

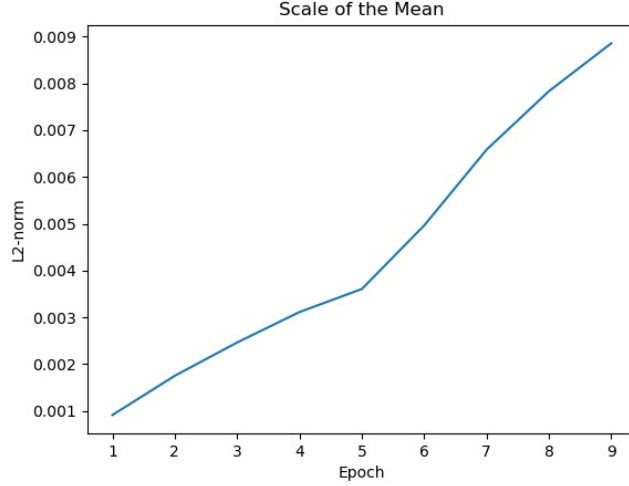


Figure 13: The L2-norm of the trainable Gaussian noise’s mean stays low.

We also try the ”w/o norm” model, which means we don’t normalize the prefix at the first step in the forwarding.

The results are shown in Fig.14, it turns out that the best performance occurs when  $Epoch = 1$ . Still, if we don’t use the normalizing trick, the overfitting problem happens. But comparing it to the original simple  $N(0, 0.016)$  (w/o norm) model in Tab.2, it has a better performance overall.

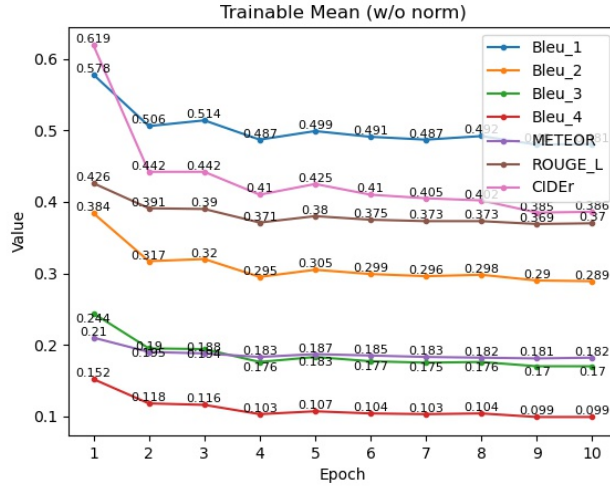


Figure 14: The overfitting problem happens. When  $Epoch = 1$ , it reaches the highest performance.

### 5.3 Auxiliary Training with Partial Images

If we are working on a different scenario: few image-text pairs with a plenty of text-only training data, which is actually more common in our daily life. Then, we may use image-text pairs for CapDec auxiliary training.

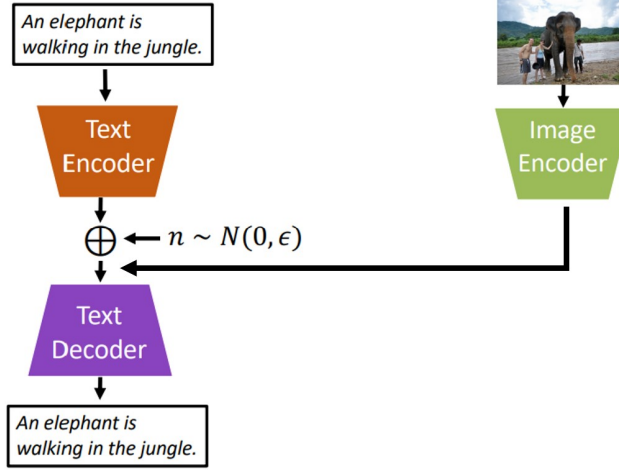


Figure 15: Auxiliary training using partial images from COCO training set.

Just as Fig.15 depicted, if we have some image-text pairs, we can then use the CLIP image encoder to obtain the image embeddings. And this time, since we don't need to bridge the gap between the embedding space, we can feed the image embeddings directly to the decoder without noise injection step. On the contrary, for those text-only data, we still need to add a noise based on the text embeddings to infer the corresponding image embeddings in the ball before feed them to the decoder.

```

noise = torch.randn(x.shape, device=device)
if not_text_only:
    noise_mask = torch.ones_like(noise, dtype=torch.bool, device=device)
    noise[is_image] = 0.0
    noise = noise_mask * noise
x = x + (noise * std)

```

(a) Partially Adding Noise.

```

if not_text_only:
    global text_image_rate
    text_mask = torch.ones_like(self.prefixes, dtype=torch.bool)
    for i in range(self.prefixes.shape[0]):
        if i % text_image_rate == 0:
            text_mask[i] = False
    image_mask = (~text_mask)
    self.prefixes = text_mask * all_data["clip_embedding_text_dave"] + image_mask * all_data["clip_embedding"]

```

(b) Partially Adding Images.

Figure 16: A part of the code of our implementation. We add images to the original text-only training set, rewrite the Dataset, DataLoader and the noise injection part in our code. We use a global "text\_image\_rate" to control the scale of the images proportion in the training set, and also ensure that there will be only few image embedding in every specific amount of embeddings when training (not sampling randomly, meaning that the distribution of the text-image pairs and text-only data is fixed). Practically, after a particular number of text embeddings, we will feed an image embeddings w/o noise injection to train the decoder.

Here, we set the text-to-image ratio to 4:1, 1:1 and evaluate their performance. And we also train a model using no text but all images. Note all of these models are trained "w/o norm", using  $N(0, 0.016)$  as noise injection. Here are the results.

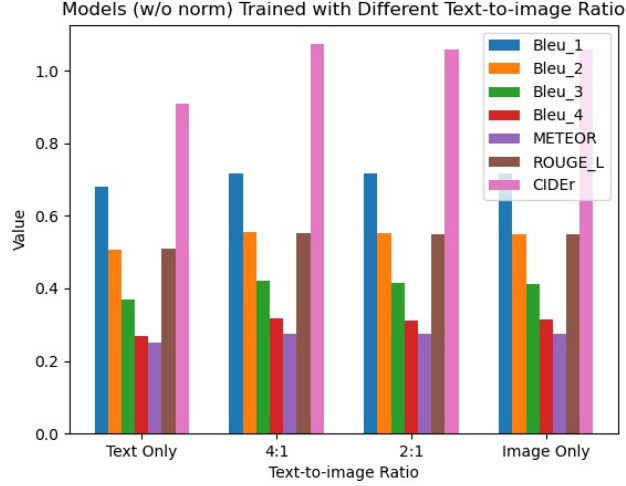


Figure 17: Note that the "Text Only" model here is the model shown in Fig.6, which is our reproduction of the paper’s best work, equipped with the normalizing trick mentioned above. Here we choose it as a baseline.

From Fig.17, we could see that this method works well. And it is astonishing to find out that the performance gaps among the models with "4:1", "2:1" and "Image Only" text-to-image ratio are really tiny. This result inspires us that by adding a small ratio of images-text pairs to the training set, we could enhance the model’s performance. And we believe that in our daily life, the real tasks must consist of some image-text pairs. By using this auxiliary training method based on the text-only CapDec training, we can gain a better performance.

On top of that, we abstract some checkpoints and evaluates their performance to see the trend when training the model (w/o norm) with 4:1 text-to-image ratio. Here are the results.

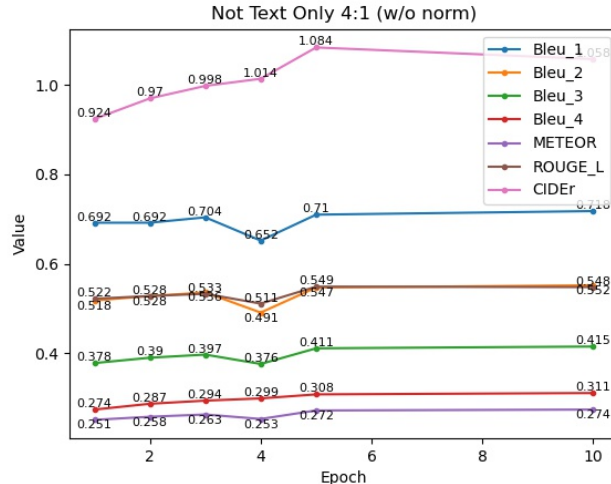


Figure 18: Note that we only sample some checkpoints for evaluation, because it is quite time-consuming.

From Fig.18, we could find out that this time, although we don’t use the normalizing trick, the overfitting problem doesn’t appear. We believe it is because the partial image-text pairs that we add

when training. It brings disturbance to the pure text encoder-decoder reconstruction process, making it hard to overfit.

We also check its training loss on the text encode-decode reconstruction process and compare it with the "Text Only" model (w/o norm) in Tab.1. Note that the text-only model (w/o norm) here is not the "Text Only" model in Fig.17, which uses the normalizing trick.

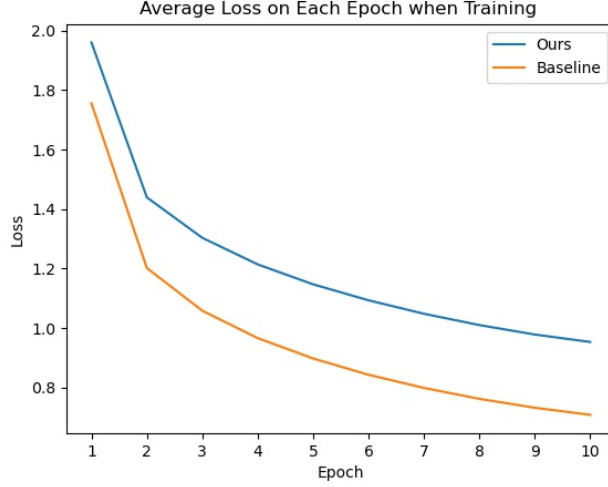


Figure 19: Comparison of the loss on the text encode-decode reconstruction process, here "Ours" means the model (w/o norm) trained with 4:1 text-to-image ratio, and "Baseline" means the text-only model (w/o norm) in Tab.1. We could see that though the "Baseline" has a lower training loss at each epoch, "Ours" has an apparently better performance. This result further verifies that, to some extent, our auxiliary training method can avoid the overfitting problem brought by text-only training.

We also train a model using the normalizing trick before noise injection and set the text-to-image ratio to 4:1 (4:1 "w/o norm" as baseline) in order to see if it can further boost the performance. Here are the results.

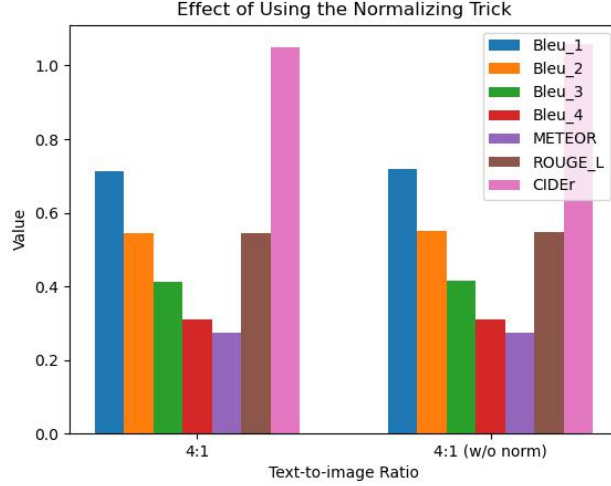


Figure 20: There is no apparent change after adding the normalizing trick before noise injection. This result further proves that our auxiliary training method has already fixed the overfitting problem thoroughly, and there is no need to do the normalization.

## 5.4 Adversarial Training

We are wondering if we can use some other methods to replace the noise injection step and enhance its performance. After a brainstorm, we come up with adversarial training.

Its behaviour is to some extent similar to the noise injection. Or in other words, generating adversarial samples is a typical way of injecting purposeful noise based on the original samples. Though the basic goal of adversarial training is to enhance the robustness of the model, and the authors claim that the noise injection is not a kind of data augmentation. We believe it is still worth trying.

Here, we implement adversarial training using the FGSM [1] method. Let’s first have a look at what FGSM is.

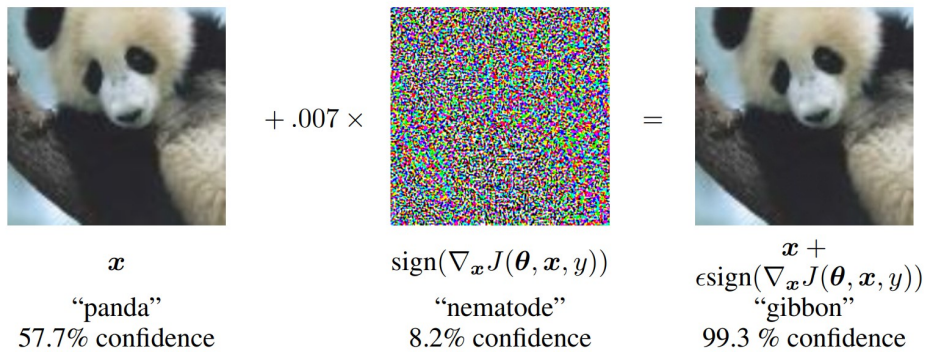


Figure 21: A demonstration of fast adversarial example generation applied to GoogLeNet (Szegedy et al., 2014a) on ImageNet. By adding an imperceptibly small vector whose elements are equal to the sign of the elements of the gradient of the cost function with respect to the input, we can change GoogLeNet’s classification of the image.

Just like the example shown in Fig.21, standard supervised training does not specify that the chosen function be resistant to adversarial examples. By training on a mixture of adversarial and

clean examples, a neural network could be regularized somewhat.

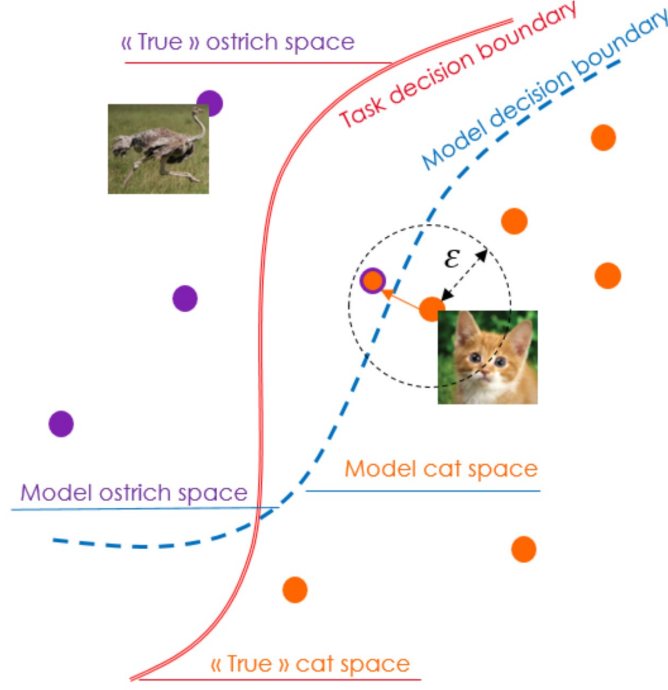


Figure 22: The essence of adversarial examples.

As shown in Fig.22, using adversarial examples to train the model actually means expanding the model decision boundary, letting it get closer to the true decision boundary. And we think this intuition could be used in our task. Adversarial training makes the surrounding embeddings mapping to a same category if they are closer in the semantic latent space, and the goal of CapDec is also to let the corresponding CLIP text embeddings and image embeddings map to the same captions. Thus, by replacing noise injection step with FGSM, which perturbs on the CLIP text embeddings, we may gain better results.

In the paper, training with an adversarial objective function based on the fast gradient sign method was an effective way to implement [1]:

$$\tilde{J}(\theta, \mathbf{x}, y) = \alpha J(\theta, \mathbf{x}, y) + (1 - \alpha) J(\theta, \mathbf{x} + \epsilon \text{sign}(\nabla_{\mathbf{x}} J(\theta, \mathbf{x}, y))) \quad (1)$$

In our task, since we only cares about the corresponding image embeddings, rather than the original text embeddings, we set the objective function to:

$$\tilde{J}(\theta, \mathbf{x}, y) = J(\theta, \mathbf{x} + \epsilon \text{sign}(\nabla_{\mathbf{x}} J(\theta, \mathbf{x}, y))) \quad (2)$$

We believe by using this method, the embeddings in the surrounding ball with radius  $\epsilon$  of the CLIP text embeddings will be decoded to the same captions.



```

outputs = model(tokens, prefix, mask)
logits = outputs.logits[:, dataset.prefix_length - 1: -1]
loss = nnf.cross_entropy(logits.reshape(-1, logits.shape[-1]), tokens.flatten(), ignore_index=0)
loss.backward()

# FGSM
if args.adv:
    prefix_grad = prefix.grad.data
    prefix = fgsm_attack(prefix, args.noise_variance, prefix_grad)
    optimizer.zero_grad()
    outputs = model(tokens, prefix, mask)
    logits = outputs.logits[:, dataset.prefix_length - 1: -1]
    loss_adv = nnf.cross_entropy(logits.reshape(-1, logits.shape[-1]), tokens.flatten(), ignore_index=0)
    loss_adv.backward()

```

(a) Code of adversarial training.

```

def fgsm_attack(data, variance, data_grad):
    sign_data_grad = data_grad.sign()
    std = math.sqrt(variance)
    perturbed_data = data + torch.randn(1, device=device) * std * sign_data_grad
    return perturbed_data

```

(b) Code of generating adversarial examples.

Figure 23: A part of the code of our implementation. We use "loss\_adv" to employ the eventual backpropagation that updates the parameters. The original "loss", however, is only used to capture the gradient of the prefixes (CLIP text embeddings). It takes nearly 4 hours to train a single epoch, which is a great cost. Therefore, we only train two model (with norm and w/o norm) using fixed disturbance scale  $\epsilon = \sqrt{0.016}$ , corresponding to the std of the Gaussian noise when using noise injection approach.

We use the model in Fig.5 as baseline and implement the adversarial training on it. Here are the results.

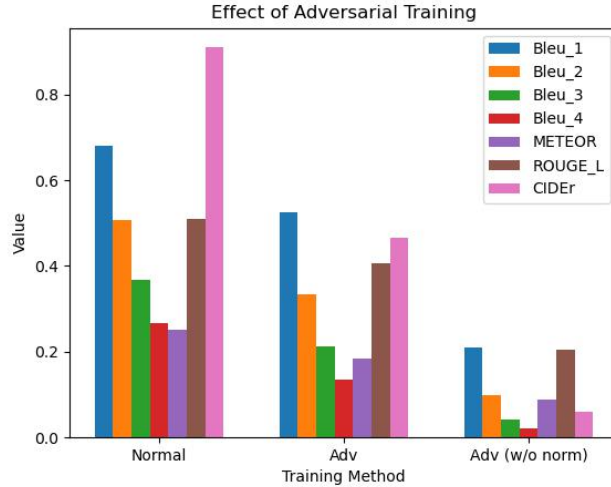


Figure 24: Note "Normal" here means the normal method of training with noise injection. "Adv" and "Adv (w/o norm)" are methods that use adversarial samples to train the model, and they don't need the noise injection step.

As shown in Fig.24, the effect of adversarial training, however, is very disappointing. We guess that the adversarial training is aimed at the small gap between the model division space and the real space, but the image and text embeddings in this image captioning task have an overall gap in the CLIP latent space, and the adversarial training cannot deal with it. Besides, image captioning task is



more complicated than a single object classification task, in which we can't solve the task by simply learn some decision boundaries.

Therefore, it is reasonable that the adversarial training can't be an effective approach compared to the noise injection.

## 6 Conclusion

Methods	BLEU_1	Bleu_2	Bleu_3	Bleu_4	METEOR	ROUGE_L	CIDEr
<b>Baseline:</b> $N(0, 0.016)$	0.684	0.506	0.365	0.264	0.250	0.511	0.903
$N(0, 0.016)$ w/o norm	0.478	0.273	0.152	0.083	0.165	0.352	0.352
$U(0, \sqrt{0.01})$ w/o norm	0.402	0.203	0.100	0.049	0.135	0.300	0.202
$U(0, \sqrt{0.049})$ w/o norm	0.413	0.208	0.102	0.049	0.141	0.313	0.208
$U(0, \sqrt{0.001})$	0.353	0.170	0.080	0.039	0.125	0.282	0.154
$U(0, \sqrt{0.016})$	0.681	0.506	0.368	0.268	0.250	0.511	0.910
$U(0, \sqrt{0.1})$	0.391	0.204	0.106	0.055	0.138	0.308	0.218
Trainable Mean w/o norm	0.506	0.317	0.195	0.118	0.190	0.391	0.442
Trainable Mean	0.679	0.505	0.368	0.268	0.253	0.514	0.915
Not Text Only (4:1) w/o norm	0.718	0.552	0.415	0.311	0.274	0.548	1.058
Not Text Only (4:1)	0.712	0.546	0.411	0.311	0.274	0.546	1.049
Not Text Only (1:1) w/o norm	0.719	0.549	0.411	0.309	0.271	0.546	1.044
All Images w/o norm	0.716	0.549	0.413	0.314	0.274	0.548	1.058
Adv w/o norm	0.163	0.098	0.041	0.020	0.088	0.205	0.059
Adv	0.525	0.335	0.213	0.136	0.183	0.407	0.465

Table 3: Metrics of different models when reaching their highest performance.

Based on the above data, we have the following conclusion.

First of all, we find out that at practice, correcting the modality offset by simply injecting a noise is not enough. It may cause severe overfitting because of the inconsistency between the text-only training process and the image-only evaluating process. Though the loss of the text encode-decode reconstruction process may seem to be low, the performance of the image captioning process can be awful eventually. To fix this overfitting problem, we need to normalize the CLIP text embeddings before the noise injection step. By doing so, the loss will get bigger, but the performance of image captioning will get better.

Moreover, we find out that the uniform noise works similar to the Gaussian noise, but is more complicated to realize in code level. This may explain why the original paper adopts the Gaussian noise rather than the uniform noise. And we presume that the authors of the paper may also have tried to compare Gaussian noise with other existed noise in mathematical field, and it is likely that similar conclusion has been drawn.

Additionally, We find out that the effect of setting a trainable noise is a bit better than the original method, which is understandable. But there is no particularly obvious improvement. Therefore, Probably due to its complexity, the original paper adopts a fixed Gaussian noise rather than a trainable noise.

Most importantly, the models trained with partial image-text pairs have significantly improvement, indicating that by adding a specific amount of image-text pairs, the text-only training model can better learn how to convert between images and text. At the same time, we find that the ratio of images has little effect on the results, indicating that only part of the images are needed to guarantee its performance. This verifies the results of the original paper to a certain extent, that is, there is a gap between the picture and the text, and we realize the image captioning by covering this gap. We think this result is of great practical significance, because the data in reality is often a large amount of text with a small number of pictures, and the datasets that completely consist of image-text pairs or text-only data are both rare. If we can take full advantages of these data, we may be able to get more

good result. Besides, this auxiliary training method can also avoid the overfitting problem brought by text-only training (for more details, please check Fig.18 and Fig.19), which means there is no need to use the normalizing trick before the noise injection step.

Last but not least, we intuitively implement adversarial training on CLIP embeddings level and undo the noise injection, but the effect of it is very unsatisfactory. We guess that the adversarial training is aimed at the small gap between the model division space and the real space, but the image and text embeddings in the image captioning task have an overall gap in the latent space, and therefore the adversarial training cannot cope with it.

## 7 Group Work

- **Li Haoyang**, in charge of implementing new proposals, training models and visualizing results (section 1, 2, 4 and 5).
- **Liu Haowen**, in charge of environment configuration, paper reproduction and evaluating models (section 3, 4, 6 and 8).

## 8 Appendix



(a) **Generated caption:** a man standing on a beach holding a surfboard.



(b) **Generated caption:** a surfborder winds up at the top of the gulf coast.

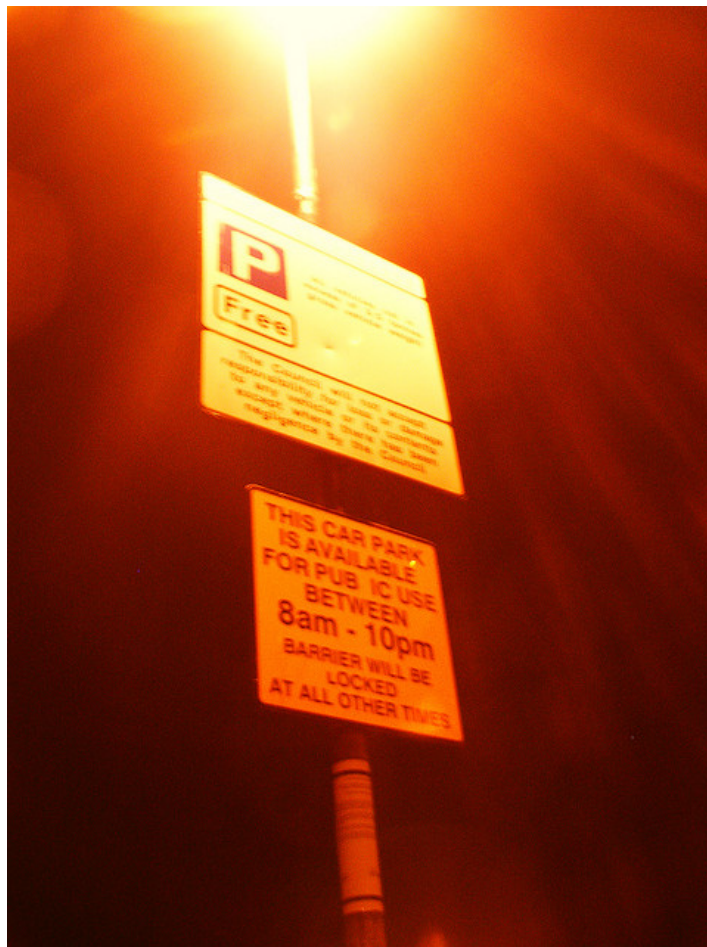


(c) **Generated caption:** a group of surfers hold their umbrellas as they celebrate the ocean waves.

Figure 25: Some **correct** captions generated from our auxiliary training model mentioned in Fig.17. It works well on semantic synthesis correlative to some typical topics like "surfers".

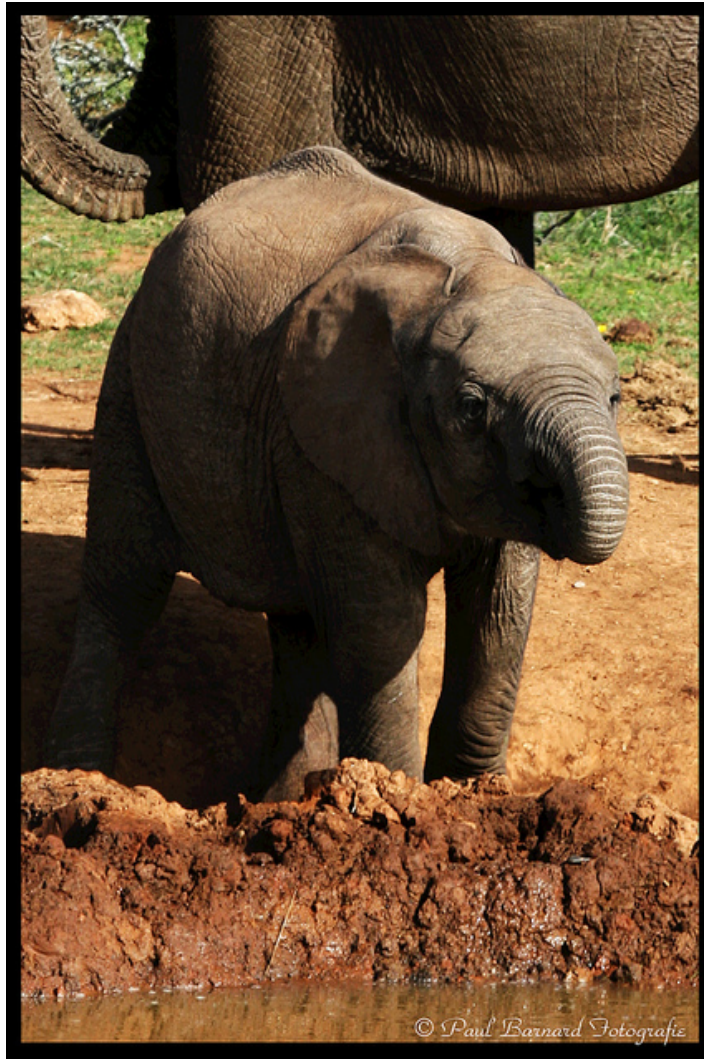


(a) **Generated caption:** a young boy stares into the distance at the words "california".



(b) **Generated caption:** a photograph of a parking lot sign stating "no parking".





(c) **Generated caption:** a baby giraffe is proud to be an adult in the zoo enclosure.

Figure 26: Some **incorrect** captions generated from our auxiliary training model mentioned in Fig.17. Fig.26a doesn't have any "california". And the sign in Fig.26b is written "parking free", rather than "no parking", which is definitely opposite. In Fig.26c, it depicts a baby elephant, but not a baby giraffe. These incorrect results, we believe, are caused by the modality gap between CLIP text embeddings and image embeddings during training and inferring.

## References

- [1] Ian J. Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples, 2014.
- [2] Ron Mokady, Amir Hertz, and Amit H. Bermano. Clipcap: Clip prefix for image captioning, 2021.
- [3] David Nukrai, Ron Mokady, and Amir Globerson. Text-only training for image captioning using noise-injected clip, 2022.