

---

# XGPT: CROSS-MODAL GENERATIVE PRE-TRAINING FOR IMAGE CAPTIONING

---

Qiaolin Xia<sup>♣</sup> Haoyang Huang<sup>◇</sup> Nan Duan<sup>◇</sup> Dongdong Zhang<sup>◇</sup> Lei Ji<sup>◇</sup>  
 Zhifang Sui<sup>♣</sup> Edward Cui<sup>♡</sup> Taroon Bharti<sup>♡</sup> Xin Liu<sup>♡</sup> Ming Zhou<sup>◇</sup>  
<sup>♣</sup>MOE Key Laboratory of Computational Linguistics, Peking University  
<sup>◇</sup>Microsoft Research Asia <sup>♡</sup>Microsoft  
 {xql,szf}@pku.edu.cn {haohua,nanduan}@microsoft.com  
 {dongdong.zhang,leiji,edwac,tbharti,mingzhou}@microsoft.com

## Abstract

While many BERT-based cross-modal pre-trained models produce excellent results on downstream understanding tasks like image-text retrieval and VQA, they cannot be applied to generation tasks directly. In this paper, we propose XGPT, a new method of **Cross-modal Generative Pre-Training for Image Captioning** that is designed to pre-train text-to-image caption generators through three novel generation tasks, **including Image-conditioned Masked Language Modeling (IMLM), Image-conditioned Denoising Autoencoding (IDA), and Text-conditioned Image Feature Generation (TIFG)**. As a result, the pre-trained XGPT can be fine-tuned without any task-specific architecture modifications to create state-of-the-art models for image captioning. Experiments show that XGPT obtains new state-of-the-art results on the benchmark datasets, including COCO Captions and Flickr30k Captions. We also use XGPT to generate new image captions as data augmentation for the image retrieval task and achieve significant improvement on all recall metrics.

## 1 Introduction

Cross-modal pre-training has substantially advanced the state of the art across a variety of Vision-and-Language (VL) tasks. **VL understanding tasks**, such as Image-Text Retrieval [1], Visual Question Answering (VQA) [2], Visual Commonsense Reasoning (VCR) [3], Referring Expression Comprehension [4], require the pre-trained model to learn the representation of visual contents, language semantics, and cross-modal alignments, but don't require generation ability. Recent pre-training methods for understanding tasks [5, 6, 7, 8, 9, 10, 11] have achieved state-of-the-art performance and attracted a lot of attention in both CV and NLP.

Vision-and-language generation tasks (e.g., Image Captioning and Text-to-Image Generation), however, requires the model to not only understand cross-modal representations but also learn generation capabilities. Thus, directly applying a model pre-trained on VL understanding tasks is not feasible. The reason is two-fold. **On one hand, pre-trained models developed for understanding tasks only provides the encoder. To support generation tasks, separate decoders have to be trained**, like the methods proposed by VideoBERT [12] and CBT [13]. **On the other hand, existing VL pre-training objectives are almost all related to the masked region or span prediction**, including VLP [14]. **None of the pre-training tasks is designed for the whole sentence generation**. Compared to studies for understanding tasks, the large-scale pretraining and fine-tuning model for VL generation tasks are still under-developed.

In this paper, we present Cross-modal Generative Pre-Training for Image Captioning (XGPT). The XGPT model uses a cross-modal encoder-decoder architecture and is directly optimized for VL generation tasks. Inspired by UniLM [15], we share parameters between the encoder and decoder to allow more effective cross-task knowledge sharing. To leverage underlying semantic and improve bidirectional generalization

between two modalities, we carefully design three generative pre-training tasks: 1) *Image-conditioned Masked Language Modeling* (IMLM), 2) *Image-conditioned Denoising Autoencoding* (IDA), 3) *Text-conditioned Image Feature Generation* (TIFG), where the encoder takes single-stream or multi-stream data as input, and the decoder is adapted to predict a wide range of generation objectives in an autoregressive manner, such as words, a sentence or image features.

Following VisualBERT [8], we perform two-stage pre-training before fine-tuning, which allows the model to better adapt to the target domain. We propose several model variants and task combinations in a thorough ablation study, which allows us to carefully control a number of factors, including model architectures, training objectives, and whether to keep placeholders for span masks.

In addition to vision-to-language generation, the proposed XGPT can also help understanding tasks like image retrieval, by performing data augmentation using our XGPT as a generator. We verified this idea by retraining a model that has state-of-the-art performance in image retrieval with the augmented data, and achieve significant improvement.

Our contributions can be summarized as follows:

- We introduce XGPT, a new method of Cross-modal Generative Pre-Training for Image Captioning, and design three novel pre-training tasks that are especially effective for image-to-text generation.
- We achieve state-of-the-art (SotA) results on COCO Captions, Flickr30k on all metrics, outperforming existing SotA and concurrent methods by a large margin. We also present extensive experiments and analysis to provide useful insights on the effectiveness of each pre-training task and model variant.
- We employ XGPT to help image retrieval, an understanding task, by performing data augmentation. After retraining, the model that has state-of-the-art performance still achieves significant improvement on all recall metrics.

## 2 Related Work

### 2.1 Pre-training for NLP Tasks

Recently, pre-trained language models (LM) over large language corpus such as ELMo [16], BERT [17], GPT2 [18], and XLNet [19] have shown great advances for NLP tasks. Among numerous works in natural language pre-training, we review three Transformer-based methods that are most relevant to our approach, namely MASS [20], Unicoder [21], and BART [22].

MASS [20] adopts the encoder-decoder framework to predict masked fragments given the remaining part of the sentence. We also use the encoder-decoder framework and train our text-only model. Unicoder is a universal language encoder pre-trained based on three pre-training tasks. The new tasks help the model learn mappings among different languages from more perspectives. BART [22] uses a denoising autoencoder for pre-training. Specifically, its pre-training objective is to reconstruct the whole sentence, which is substantially different from the masked language modeling in BERT. Our method is inspired by these works, but since images are not sequential data, we have to tailor our model for cross-modal tasks in particular.

### 2.2 Pre-training for Cross-modal Generation Tasks

Very recently, several attempts have been made to pre-train models for cross-modal generation tasks.

Both VideoBERT [12] and CBT [13] are seeking to conduct pre-training for the video captioning task. But they perform pre-training only for the BERT-based encoder to learn bidirectional joint distributions over sequences of visual and linguistic tokens. So they have to train a separate video-to-text decoder. In contrast, Unified VLP [14], concurrently with our work, uses a shared multi-layer transformer network for both encoding and decoding. Following UniLM [15], they pre-train the model on two masked language modeling (MLM) tasks, like cloze tasks designed for sequence-to-sequence LM. So target prediction is still masked tokens, not the whole sentence. However, we find that by using generative pre-training objectives such as Image-conditioned Masked Language Modeling, Image-conditioned Denoising Autoencoding and Text-conditioned Image Feature Generation, XGPT can outperform Unified VLP significantly on Image Captioning.

### 3 Preliminaries

**Linguistic Representation.** For each token in the input language sequence, its representation is a sum of token embedding and position embedding. We denote the input tokens as  $\mathbf{w} = \{w_1, w_2, \dots, w_M\}$  and the corresponding representations as  $\mathbf{x}^T = \{x_1^T, x_2^T, \dots, x_M^T\}$ .

**Image Representation.** For each input image, we first detect objects using a pre-trained Faster R-CNN model [23]. Here, the top 100 objects with highest confidence scores are selected, each of which has a feature vector computed by mean-pooling the last-layer convolutional feature of its region of interest. To represent the position of each object, we construct a 5-d position vector from its spatial location (normalized top-left and bottom-right coordinates) and the fraction of image area it covered. Next, we concatenate the feature vector and position vector of each object and transform it into another vector by linear projection, to make sure the dimensions  $h$  of linguistic tokens and visual tokens are identical, and we denote as image regions as  $\mathbf{v} = \{v_1, v_2, \dots, v_N\}$  and the corresponding representations as  $\mathbf{x}^I = \{x_1^I, x_2^I, \dots, x_N^I\}$ .

**Image Refining.** Unlike words in text, image regions lack a natural ordering. To better model the relationship among objects in an image, we add an additional image refining layer following AoANet [24] to refine the image representation before feeding them to the encoder. We refer readers to Appendix A.1 for technical details.

## 4 Cross-modal Generative Pre-Training for Image Captioning

### 4.1 Revisit Pre-training Tasks

In this section, we first review the objectives of the classical masked language modeling that is commonly used for understanding tasks.

The objective of masked language and region modeling is to learn joint representations for both vision and language by reconstructing masked tokens  $\bar{\mathbf{w}}$  from a corrupted version  $\hat{\mathbf{w}}$ :

$$\max_{\theta} \mathcal{L}_{MLM} = \log p_{\theta}(\bar{\mathbf{w}}|\hat{\mathbf{w}}) = \sum_{t=1}^T \log p_{\theta}(w_t|\hat{\mathbf{w}})m_t \quad (1)$$

where  $m_t = 1$  if  $w_t$  is masked as corruption, and 0 otherwise. The objective is designed based on *bidirectional* contexts which allows the words in the future to be attended. So only the masked tokens of the sentence or labels of regions are required to be predicted instead of the whole sentence.

The pros and cons of this pre-training objective can be concluded in the following aspects:

- *Downstream tasks:* Many BERT-based cross-modal pre-pretraining models choose masked language and region modeling as their main pre-training task. Because the encoder learned from this task can provide representations of both vision and language based on *bidirectional* context and it is naturally fit for many understanding tasks, including VQA, Image Retrieval, etc.
- *Architecture modification:* For downstream generation tasks, models pre-trained only through mask prediction tasks usually have to train an additional layer for generation, as pointed out by Song et al. [20]. But separately trained decoders could create a pretrain-finetune discrepancy that hurt the generality of the model. This results in a gap between pre-training and fine-tuning on generation tasks.

To address these concerns, we propose XGPT and the new pre-training tasks.

### 4.2 Model Architecture

XGPT has a unified encoder and decoder architecture and can be pre-trained through different generative pre-training tasks. Basically, both encoder and decoder are multi-layer Transformer networks. The encoder reads the source image and sentence and generates a set of representations as introduced in Section 3. Different from other BERT-based encoder-only models, the probability of each target token is estimated by the decoder given the *cross-attention* performed over the final hidden layer of the encoder.

Specially, we use shared parameters for encoder and decoder, and a faster attention strategy in decoder where the weights of self-attention and encoder-decoder attention are shared. We add a signal in the attention network to distinguish whether keys and values are from the output of the encoder for efficient re-use. For

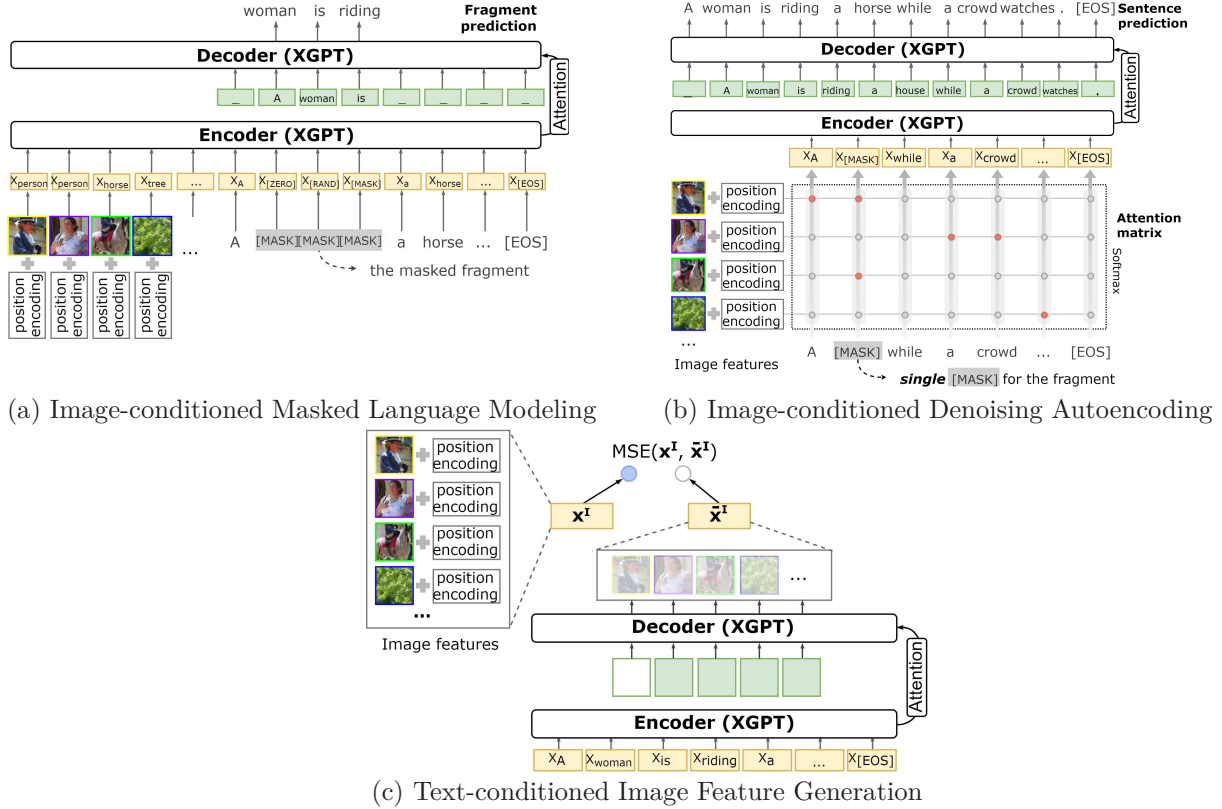


Figure 1: Cross-modal generative pre-training tasks

our base model, we use 12 layers in the encoder and decoder. And we use 6 layers for our tiny model to explore the necessity of the decoder for image captioning in the experimental analysis (Section 6).

### 4.3 Generative Pre-training Tasks

Unlike Unified VLP, we use the image captioning task as a basic generative task in the pre-training stage. It only takes images as inputs (single modality). We also introduce three new cross-modal generative pre-training tasks that can jointly pre-train the encoder and decoder: IMLM, IDA, and TIFG.

**Image Captioning (IC).** A major approach to Image Captioning is encoder-decoder framework. Giving the image regions  $\mathbf{v}$ , the training objective is to generate the caption  $\mathbf{w}$  in an autoregressive manner by minimize the negative log-likelihood

$$\mathcal{L}_{\text{IC}} = -\lambda_{\text{IC}} \sum_{t=1}^T \log p_{\theta}(w_t | \mathbf{w}_{<t}, \mathbf{v}) \quad (2)$$

where  $\mathbf{w}_{<t}$  is the context history produced by the neural generator, and  $\lambda_{\text{IC}}$  is the pre-defined weight of IC loss. This objective requires the model to predict the whole sentence from scratch.

**Image-conditioned Masked Language Modeling (IMLM).** IMLM aims to teach the model to learn the relationship between vision and language by predicting consecutive tokens in the decoder side.

XGPT is trained to reconstruct the n-gram masked words through a sequence to sequence framework. This task is similar in the idea to n-gram MLM in BERT or Masked Seq-to-Seq in MASS. The difference lies in that (1) the task encourages the encoder to learn the cross-modality relationships between the unmasked tokens and image regions, and (2) the decoder has to generate masked tokens of the fragment, and extract useful image-conditioned information from the encoder side.

As shown in Figure 1(a), we concatenate the regions and unmasked tokens as input to the encoder during pre-training. And we let the decoder predict the masked fragment by minimizing the negative log-likelihood loss:

$$\mathcal{L}_{\text{IMLM}} = -\lambda_{\text{IMLM}} \sum_{t=1}^M \log p_{\theta}(w_t | \mathbf{w}_{<t}, \hat{\mathbf{w}}, \mathbf{v}) m_t \quad (3)$$

where  $\hat{\mathbf{w}}$  is the corrupted caption, and  $m_t = 1$  if  $w_t$  is in the masked fragment, and 0 otherwise.  $\lambda_{\text{IMLM}}$  is the pre-defined weight of IMLM loss.

**Image-conditioned Denoising Autoencoding (IDA).** In IDA we take distance between feature spaces of two modalities into account, and use an attention matrix to model the underlying text-image alignments. Besides, IDA forces the model to reconstruct the whole sentence without giving it the length of the masked fragment, as illustrated in Figure 1(b). Specifically, we use

- *Single [MASK] token for fragments.* Inspired by text filling task in [22], we first sample n-gram fragments to be mask and, then, replace each with a single [MASK] token. This is more challenging because the model have to predict not only the missing tokens but also the length of the original sentence. We also compare this method with conventional method which keeps the placeholder for all masked tokens in Section 6.
- *The attention-driven image-text matching.* To model the text-image alignments, we first compute an attention matrix for each token-region pair  $(w_i, v_j)$ :  $A_{ij} = W[x_i^T, x_j^I, x_i^T \odot x_j^I]$  where  $W \in \mathbf{R}^{3 \times h}$  is a trainable weight and  $\odot$  is elementwise multiplication. Then, we represent each word as the weighted sum of all region representations based on the attention matrix:  $x_i^T = \sum_{j=1}^N \text{softmax}(A_{ij}) x_j^I$ . Finally, XGPT takes new  $x_i^T$  as input and tries to predict the original word sequence  $\mathbf{w}$ . The loss function is defined as

$$\mathcal{L}_{\text{IDA}} = -\lambda_{\text{IDA}} \sum_{t=1}^M \log p_{\theta}(w_t | \mathbf{w}_{<t}, \hat{\mathbf{w}}, \mathbf{v}) \quad (4)$$

where  $\hat{\mathbf{w}}$  is the corrupted caption, and  $\lambda_{\text{IDA}}$  is the pre-defined weight of IDA loss.

**Text-conditioned Image Feature Generation (TIFG).** Text-to-image (T2I) generation can be regarded as the inverse problem of image captioning, a Image-to-Text (I2T) problem. It is natural and reasonable to unify the model to leverage the underlying semantic in both domains.

In contrast to Uniter [11], which involved image feature generation by learning to regress the transformer output of each masked region to its visual features, TIFG aims to regress the decoder output of all image regions conditioned on text descriptions rather than only the masked regions. As shown in Figure 1(c), we employ the encoder-decoder pipeline to convert linguistic representations into  $\bar{x}_i^I$  of the same length and dimension as image representations  $x_i^I$ . Then we train with mean squared error to supervise the XGPT to generate semantically consistent image features. Mathematically, this loss can be expressed as:

$$\mathcal{L}_{\text{TIFG}} = \lambda_{\text{TIFG}} \frac{1}{N} \sum_{i=1}^N \|x_i^I - \bar{x}_i^I\|_2^2 \quad (5)$$

where  $\lambda_{\text{TIFG}}$  is the weight of TIFG loss.

**Multi-task pre-training.** Following Unicoder [21], we calculate task-specific loss in turns and update the model for each task in every pre-training iteration. We include these tasks with individual weight of loss to study how each objective works for pre-training.

## 5 Experiments and Results

### 5.1 Training Stages

**Out-of-domain pre-training stage.** We first conduct pre-training on Conceptual Captions (CC) dataset [25] which contains about 3M image-caption pairs scraped from alt-text enabled web images. The automatic collection leaves some noise (e.g., not relevant and too short) in the dataset but brings a massive scale. So we use it only as our out-of-domain dataset for the first pre-training stage. We set the weight to 1 for all pre-training tasks on CC.

**In-domain pre-training stage.** Before fine-tuning XGPT on the final image captioning task, we find it beneficial to further pre-train the model using the data from downstream tasks with the proposed pre-training objectives. This step allows the model to adapt to the target domain. So we reduce the weights of the cross-modal tasks (i.e.,  $\lambda_{\text{IMLM}}$ ,  $\lambda_{\text{IDA}}$ ,  $\lambda_{\text{TIFG}}$ ) and keep the image caption task unchanged.

**Fine-tuning stage.** In this step, the model only takes image features and position information as input, and the decoder is trained to predict the whole sentence in an autoregressive manner. We also applied other inference approaches like beam search as well. More details follow Section 5.3.

## 5.2 Evaluation datasets

The datasets for downstream tasks include COCO Captions [26] and Flickr30k [27]. In these datasets, each image is labeled with 5 captions. We follow Karpathy’s split<sup>1</sup>, which gives 113.2k/5k/5k and 29.8k/1k/1k images for train/val/test splits respectively. We use standard metrics for Image Captioning, including BLEU@4, METEOR, CIDEr, SPICE, to evaluate the propose method and compare with other methods.

## 5.3 Implementation Details

In all experiments, the backbone Transformer of XGPT follows Vaswani et al. [28], and we modify it to the BERT-based encoder-decoder architecture with 768 hidden units, 8 heads, GLEU activations used as GPT [29]. The dropout rate is 0.1.

We train XGPT with mixed-precision training and FP16, which makes use of GPUs more efficiently. The Adam [30] with  $\beta_1 = 0.9$ ,  $\beta_2 = 0.98$  is used for optimization. The learning rate is varying from  $1e - 4$  to  $2e - 5$  for out-of-domain pre-training with invert square root decay [28]. The weight of the individual loss is set to 1. During the in-domain pre-training and fine-tuning stage, we take an average of the top-4 pre-trained weights and reduce the initial learning rate to  $1e - 5$ . We include  $\lambda_{\text{IMLM}}$ ,  $\lambda_{\text{IDA}}$ ,  $\lambda_{\text{TIFG}}$  with 0.3 and  $\lambda_{\text{IC}} = 1$  for in-domain pre-training, and turn off other tasks during fine-tuning (image captioning only). The two-stage pre-training takes about 8 days to converge on 8x V100 GPUs with a total batch size of 512 by gradient accumulation. We fine-tuned the model 30 epochs with four GPUs.

For caption inference, we use greedy search on the validation set and beam search with beam size 2 on the test set. We describe more training details in Appendix A.2.

## 5.4 Experimental Settings

We compare XGPT with state-of-the-art methods on image captioning in two settings:

- **Text Pre-training** is our baseline setting where XGPT is trained from scratch with two pre-trained objectives: Masked Language Model in BERT [17] and Masked Seq-to-Seq in MASS [20], and directly fine-tuned on the image-caption task.
- **Image-Text Pre-training.** Weights are also initialized from Text Pre-training. To fuse the information from Image-Text during the pre-training stage, XGPT continue pre-training on three proposed tasks along with the image captioning task, and fine-tuned only on the image captioning task in the end.

## 5.5 Comparisons against SotAs

Results comparing our methods and SotA methods on the test set are in 1. We include state-of-the-art works (mostly without pre-training) and a recent work [14] that also use pre-trained models, and our methods (the lowest part). All methods use a single model for the image captioning for a fair comparison. Our full model (Base) significantly outperforms SotA methods on all four metrics on both COCO Captions and Flickr30k.

**Compare with baseline.** Compared to our baseline model that only uses text pre-training, cross-modal pre-training tasks improves the performance on all metrics, which validates the importance of **Image-Text Pre-training** for generation tasks. Unified VLP also employs two masked language

<sup>1</sup>cs.stanford.edu/people/karpathy/deepimagesent/caption\_datasets.zip



Model	Flick30k				COCO			
	C	B@4	M	S	C	B@4	M	S
Approaches that <i>do not</i> use any pre-trained model								
BUTD [23]	56.6	27.3	21.7	16.0	113.5	36.2	27.0	20.3
NBT (with BBox) [31]	57.5	27.1	21.7	15.6	107.2	34.7	27.1	20.1
GCN-LSTM (spa) [32]	-	-	-	-	115.6	36.5	27.8	20.9
GCN-LSTM (sem) [32]	-	-	-	-	116.3	36.8	27.9	20.9
GVD [33]	62.3	27.3	22.5	16.5	-	-	-	-
AoANet [24]	-	-	-	-	119.8	37.2	28.4	21.3
Approaches that <i>do</i> use pre-trained models or pre-trained language models								
Unified VLP* [14]	56.8	27.6	20.9	15.3	114.3	35.5	28.2	21.0
Unified VLP [14]	67.4	30.1	23.0	17.0	116.9	36.5	28.4	21.2
<b>XGPT*</b>	53.0	25.7	20.3	14.7	113.0	34.4	27.8	20.8
<b>XGPT</b>	<b>70.9</b>	<b>31.8</b>	<b>23.6</b>	<b>17.6</b>	<b>120.1</b>	<b>37.2</b>	<b>28.6</b>	<b>21.8</b>

Table 1: Comparison with the previous state-of-the-art methods. **Bold** indicates best value overall. Unified VLP\* and XGPT\* perform **Text Pre-training**. The former is initialized from UniLM, while the latter is pre-trained from scratch with less text data, which is detailed in Section 5.4. Both Unified VLP and XGPT perform **Image-Text Pre-training** (see Section 5.4) where the weights are initialized from Text Pre-training and pre-trained on different tasks, respectively.

Stage	Pre-training Tasks	COCO			
		C	B@4	M	S
Out-of-domain (CC)	IC	116.4	35.9	28.2	21.1
	IC + IMLM	117.7	36.2	28.2	21.2
	IC + IDA	118.1	36.4	28.3	21.3
	IC + TIFG	117.3	36.0	28.2	21.2
	IC + IMLM + IDA + TIFG	<b>118.3</b>	<b>36.5</b>	<b>28.4</b>	<b>21.3</b>
Out-of-domain (CC) + In-domain (COCO)	IC + IMLM	119.1	36.7	28.5	21.5
	IC + IDA	119.2	36.6	28.5	21.6
	IC + TIFG	118.2	36.4	28.4	21.3
	IC + IMLM + IDA + TIFG	<b>120.1</b>	<b>37.2</b>	<b>28.6</b>	<b>21.8</b>

Table 2: Ablation analysis of pre-training tasks on COCO Captions.

Model	C	B@4	M	S
Tiny Enc	112.1	34.1	27.9	20.6
Tiny EncDec	110.8	33.7	27.6	20.5
Tiny EncDecShare	<b>112.7</b>	<b>34.6</b>	<b>27.9</b>	<b>20.7</b>
Base EncDecShare	<b>112.9</b>	<b>35.0</b>	<b>27.9</b>	<b>20.8</b>

Table 3: Evaluation results on COCO Captions using different model structures. We use 6-layer Transformers for Tiny models, and 12-layer for Base models and directly train the model on image captioning task without any pre-training.

pre-training tasks and provides a baseline initialized from a language pre-trained model (UniLM)[15]. The improvement is significant on both benchmark datasets, and is particular sound on Flickr30k, including absolute 17.9% gain on CIDEr, 6.1% on BLEU@4, 3.3% on METEOR, and 2.9% on SPICE. Comparing the gain from **Image-Text Pre-training** tasks, XGPT achieves higher improvement than Unified VLP, demonstrating the effectiveness of the three designs in XGPT.

	C	B@4	M	S
multi [MASK]	117.8	36.1	28.2	21.3
single [MASK]	<b>118.1</b>	<b>36.4</b>	<b>28.3</b>	<b>21.3</b>

Table 4: Comparison of two masking methods on COCO Captions.

	R@1	R@5	R@10
ViLBERT [5]	58.2	84.9	91.5
ViLBERT + augmentation	<b>60.4</b>	<b>86.4</b>	<b>91.9</b>

Table 5: Results of image retrieval task on Flickr30k.

## 6 Analysis

**Is Decoder Necessary?** To find the best model structure for image captioning, we also designed three model variants. **Enc** is a single multi-layer Transformer that takes either image region features for the unidirectional model or a pair of text and image packed together, uses different self-attention masks to control the access to the context, which is similar to UniLM and Unified VLP. **EncDec** is a Transformer encoder-decoder architecture in which all the weights are initialized randomly and not shared between the encoder and decoder. **EncDecShare** is like EncDec, but the parameters between encoder and decoder are shared. We also use a signal in the attention network to control whether keys and values are from the encoder. This greatly reduces the memory footprint of the model.

In all settings, models are trained on the image captioning task without any text or image pre-training.

Table 3 reports results of these settings on two model sizes: Base (layers=12), Tiny (layers=6). With a similar number of parameters, the shared setup with 6-layer (Tiny EncDecShare) performs better than the encoder and decoder parameters are not shared. Tiny Enc model which simply reuses the encoder for decoding can outperform Tiny EncDec, although it performs less well than EncDecShare models. We also notice that the model with shared 12-layer encoder and decoder parameters (Base EncDecShare) performs best. We use this as the optimal architecture that achieves the new state-of-the-art results in Table 1.

**Effectiveness of Proposed Tasks.** We analyze the effectiveness of different pre-training tasks through ablation studies over COCO Captions and Flickr30K. The results are shown in Table 2. Firstly, we establish our baseline: Row 1 in Table 2 shows the results of the model pre-trained on out-of-domain datasets through the image captioning task only.

As for the out-of-domain pre-training stage, there are significant improvements across all three tasks (comparing Row 2,3,4 with Row 1 baseline). Among the three, we observe IDA which helps the model to learn text-image alignments achieves the biggest improvement, while TIFG the smallest. This is probably because of the discrepancy of the decoder which is originally designed to predict captions and the task objective which is to predict all image region features. When combining all three tasks, we find that they are complementary to each other and see the highest gain of approximately +1.9 on CIDEr over Row 1.

Comparing with one-stage pre-training, we find that each task combination pre-trained after the second stage (Row 6-8) gains approximately +2 on CIDEr. This indicates that two-stage pre-training with in-domain data enables the model to adapt to the downstream data better than only using out-of-domain pre-training. Combining all three tasks leads to the highest score on all metrics. We use this as the optimal pre-training setting for further experiments.

**How to Efficiently Mask?** We conduct further experiments to compare two masking strategies for IDA: (1) Multi [MASK]: replace tokens in the sampled fragment with exactly the same number of [MASK] tokens, (2) Single [MASK]: replace the fragment with a single [MASK] token. Results shown in Table 4 suggest that the effectiveness of IDA highly dependent on the masking strategy. Single [MASK] is obviously better than Multi [MASK]. Multi [MASK] provides the decoder with full position information of masked token, which reduced the difficulty for the decoder to predict correct words, thus, gives a less well performance (-0.3 on CIDEr). Therefore, we use the single [MASK] as the optimal pre-training setting.




	<b>Human-generated captions</b>
	A person trying desperately <u>not to be photographed</u> by putting their sweater over ...
	A person wearing red pants <u>hides their head</u> under a <u>black jacket</u> in front of a desk ...
	A person with red pants with cover over her head <u>sitting</u> in front of multiple computers.
	<u>The woman</u> tries to <u>hide from work</u> under a black sweatshirt, but her red corduroy ... <u>a child</u> is <u>hiding</u> under a <u>sweater</u> in a chair.
	<b>XGPT-generated captions</b>
	<u>A woman</u> wearing red pants is <u>sitting</u> at a desk in front of a computer.
	A person in a <u>blue sweatshirt</u> is <u>sleeping</u> in a chair.
	A woman in a <u>blue sweatshirt</u> is <u>sleeping</u> in a chair in front of a computer.

Table 6: An example of generated captions for the given image. Underlined text shows the difference between captions. We can see that in the original training data, underlined text is usually people’s guess with personal emotions, e.g., *hide from work*. While the generated captions provide more modifier variants (e.g., *blue*) and verb variants (e.g., *sleeping*) according to what can be seen in the picture.


	<b>Human-generated captions</b>
	A young boy wearing a black shirt, BROWN pants and a black watch has his hand on a ...
	A young man wearing a black shirt takes a folding chair from a large stack.
	A person helping to set up chairs for a big event.
	A young man in a black shirt stacks chairs. A boy is setting up folding chairs.
	<b>XGPT-generated captions</b>
	A man in a black t shirt and <u>black shorts</u> is putting up a white chair.
	A man in a black t shirt and <u>black t shirt</u> works on a folding chair.

Table 7: A negative example of the generation results. The first predicted the wrong color of the pants (brown→black). And the second generated caption duplicate the same phrase (black t shirt).

## 6.1 Data Augmentation for Image Retrieval.

In addition to generation tasks, our XGPT can also help vision-and-language understanding tasks, such as image retrieval, by performing data augmentation as an image description generator. Image retrieval is a task of identifying an image from a pool given a caption describing its content. We generate 62k more captions for all 29k images (about 2.1 captions for each) in the Flickr30k training set, which originally contains 145k captions. We continue to fine-tune the open-source state-of-the-art model<sup>2</sup> introduced in [5] on the combination of the augmentation and the original training data.

Comparing Row 1 (trained only with the original training data) against Row 2 (fine-tuned on augmented data) in Table 5, the improvement is significant (2.2% on R@1, 1.5% on R@5, and 0.4% on R@10). The higher relative gain on R@1 also indicates that the generator can produce high-quality image captions which can help the model better understand images.

A negative example of the generation results is provided in Table 7. The first sentence contains wrong information (brown→black); the second has a duplicated phrase. Both can be considered as noise.

Table 6 shows a positive example of XGPT-generated captions. We can see that the generated captions are grammatically and semantically correct, and also can increase the diversity of the data.

## 7 Conclusion

In this paper, we present XGPT, Cross-modal Generative Pre-Training for Image Captioning. Three main pre-training tasks are proposed and the ablation study shows that the effectiveness of each task is different. The combination of all tasks achieves stronger performance on all evaluation metrics suggested that they are complementary to each other. After in-domain and out-of-domain pre-training, XGPT outperforms state-of-the-art models by a significant margin. For future works, we are curious about extending XGPT to cross-modal understanding tasks, such as VQA and VCR.

<sup>2</sup>[https://github.com/jiasenlu/vilbert\\_beta](https://github.com/jiasenlu/vilbert_beta)

## References

- [1] Andrej Karpathy and Li Fei-Fei. Deep visual-semantic alignments for generating image descriptions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3128–3137, 2015.
- [2] Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C Lawrence Zitnick, and Devi Parikh. Vqa: Visual question answering. In *Proceedings of the IEEE international conference on computer vision*, pages 2425–2433, 2015.
- [3] Rowan Zellers, Yonatan Bisk, Ali Farhadi, and Yejin Choi. From recognition to cognition: Visual commonsense reasoning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6720–6731, 2019.
- [4] Sahar Kazemzadeh, Vicente Ordonez, Mark Matten, and Tamara Berg. Referitgame: Referring to objects in photographs of natural scenes. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, pages 787–798, 2014.
- [5] Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. Vilbert: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks. In *Advances in Neural Information Processing Systems*, pages 13–23, 2019.
- [6] Chris Alberti, Jeffrey Ling, Michael Collins, and David Reitter. Fusion of detected objects in text for visual question answering. *EMNLP-IJCNLP*, 2019.
- [7] Hao Tan and Mohit Bansal. Lxmert: Learning cross-modality encoder representations from transformers. In *EMNLP-IJCNLP*, pages 5103–5114, 2019.
- [8] Liunian Harold Li, Mark Yatskar, Da Yin, Cho-Jui Hsieh, and Kai-Wei Chang. Visualbert: A simple and performant baseline for vision and language. *arXiv preprint arXiv:1908.03557*, 2019.
- [9] Gen Li, Nan Duan, Yuejian Fang, Daxin Jiang, and Ming Zhou. Unicoder-vl: A universal encoder for vision and language by cross-modal pre-training. *AAAI*, 2020.
- [10] Weijie Su, Xizhou Zhu, Yue Cao, Bin Li, Lewei Lu, Furu Wei, and Jifeng Dai. Vl-bert: Pre-training of generic visual-linguistic representations. *arXiv preprint arXiv:1908.08530*, 2019.
- [11] Yen-Chun Chen, Linjie Li, Licheng Yu, Ahmed El Kholy, Faisal Ahmed, Zhe Gan, Yu Cheng, and Jingjing Liu. Uniter: Learning universal image-text representations. *arXiv preprint arXiv:1909.11740*, 2019.
- [12] Chen Sun, Austin Myers, Carl Vondrick, Kevin Murphy, and Cordelia Schmid. Videobert: A joint model for video and language representation learning. *Proceedings of the IEEE International Conference on Computer Vision*, 2019.
- [13] Chen Sun, Fabien Baradel, Kevin Murphy, and Cordelia Schmid. Contrastive bidirectional transformer for temporal representation learning. *ArXiv*, abs/1906.05743, 2019.
- [14] Luowei Zhou, Hamid Palangi, Lei Zhang, Houdong Hu, Jason J Corso, and Jianfeng Gao. Unified vision-language pre-training for image captioning and vqa. *AAAI*, 2020.
- [15] Li Dong, Nan Yang, Wenhui Wang, Furu Wei, Xiaodong Liu, Yu Wang, Jianfeng Gao, Ming Zhou, and Hsiao-Wuen Hon. Unified language model pre-training for natural language understanding and generation. *arXiv preprint arXiv:1905.03197*, 2019.
- [16] Matthew E Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. Deep contextualized word representations. In *NAACL, 2018. 2*, 2018.
- [17] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics*, pages 4171–4186, 2019.
- [18] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language models are unsupervised multitask learners. *OpenAI Blog*, 1(8), 2019.
- [19] Zhilin Yang, Zihang Dai, Yiming Yang, Jaime Carbonell, Ruslan Salakhutdinov, and Quoc V Le. Xlnet: Generalized autoregressive pretraining for language understanding. *arXiv preprint arXiv:1906.08237*, 2019.
- [20] Kaitao Song, Xu Tan, Tao Qin, Jianfeng Lu, and Tie-Yan Liu. Mass: Masked sequence to sequence pre-training for language generation. *arXiv preprint arXiv:1905.02450*, 2019.

- [21] Haoyang Huang, Yaobo Liang, Nan Duan, Ming Gong, Linjun Shou, Daxin Jiang, and Ming Zhou. Unicoder: A universal language encoder by pre-training with multiple cross-lingual tasks. *arXiv preprint arXiv:1909.00964*, 2019.
- [22] Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Ves Stoyanov, and Luke Zettlemoyer. Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. *arXiv preprint arXiv:1910.13461*, 2019.
- [23] Peter Anderson, Xiaodong He, Chris Buehler, Damien Teney, Mark Johnson, Stephen Gould, and Lei Zhang. Bottom-up and top-down attention for image captioning and visual question answering. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6077–6086, 2018.
- [24] Lun Huang, Wenmin Wang, Jie Chen, and Xiao-Yong Wei. Attention on attention for image captioning. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 4634–4643, 2019.
- [25] Piyush Sharma, Nan Ding, Sebastian Goodman, and Radu Soricut. Conceptual captions: A cleaned, hypernymed, image alt-text dataset for automatic image captioning. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2556–2565, 2018.
- [26] Xinlei Chen, Hao Fang, Tsung-Yi Lin, Ramakrishna Vedantam, Saurabh Gupta, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco captions: Data collection and evaluation server. *arXiv preprint arXiv:1504.00325*, 2015.
- [27] Peter Young, Alice Lai, Micah Hodosh, and Julia Hockenmaier. From image descriptions to visual denotations: New similarity metrics for semantic inference over event descriptions. *Transactions of the Association for Computational Linguistics*, 2:67–78, 2014.
- [28] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008, 2017.
- [29] Dan Hendrycks and Kevin Gimpel. Bridging nonlinearities and stochastic regularizers with gaussian error linear units. 2016.
- [30] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *international conference on learning representations*, 2015.
- [31] Jiasen Lu, Jianwei Yang, Dhruv Batra, and Devi Parikh. Neural baby talk. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7219–7228, 2018.
- [32] Ting Yao, Yingwei Pan, Yehao Li, and Tao Mei. Exploring visual relationship for image captioning. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 684–699, 2018.
- [33] Luowei Zhou, Yannis Kalantidis, Xinlei Chen, Jason J Corso, and Marcus Rohrbach. Grounded video description. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6578–6587, 2019.
- [34] Philipp Koehn, Hieu Hoang, Alexandra Birch, Chris Callison-Burch, Marcello Federico, Nicola Bertoldi, Brooke Cowan, Wade Shen, Christine Moran, Richard Zens, et al. Moses: Open source toolkit for statistical machine translation. In *Proceedings of the 45th annual meeting of the association for computational linguistics companion volume proceedings of the demo and poster sessions*, pages 177–180, 2007.

## A Appendices

### A.1 AoA refining layer

Instead of directly feeding image region features to the encoder, we build a refining network to refine their representation using an AoA module following [24]. The AoA module adopts the multi-head attention function [28] where  $\mathbf{Q} = W_Q \mathbf{x}$ ,  $\mathbf{K} = W_k \mathbf{x}$ , and  $\mathbf{V} = W_v \mathbf{x}$  are three individual linear projections of the region features  $\mathbf{x}$ . The AoA layer is formulated as

$$\begin{aligned} AoA(f_{att}, \mathbf{Q}, \mathbf{K}, \mathbf{V}) &= \sigma(W_1 \mathbf{Q} + W_2 f_{att}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) + b_1) \\ &\odot (W_3 \mathbf{Q} + W_4 f_{att}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) + b_2) \end{aligned}$$

where  $W_1, W_2, W_3, W_4 \in \mathbb{R}^{h \times h}$ ,  $b_1, b_2 \in \mathbb{R}^h$  are learnable weights,  $f_{att}$  is a conventional attention module which operates on some queries, keys and values and generates some weighted average vectors.

The refined region features  $\mathbf{x}'$  are calculated by

$$\mathbf{x}' = LayerNorm(\mathbf{x} + AoA(f_{att}, W_q \mathbf{x}, W_k \mathbf{x}, W_v \mathbf{x}))$$

### A.2 More Implementation Details

For tokenization, we follow the line of [34] to build vocabulary from English Wikipedia, and use byte-pair encoding (BPE) to process image captions. And we trim the max sequence length to 60 and represent each input image as 100 object regions.

Following [17], the masked tokens in the encoder will be a [MASK] token 80% of the time, a random token 10% of the time and an unchanged token 10% of the time. For IMLM, we set the fragment length as roughly 50% of the total number of tokens in the sentence. And the span lengths are drawn from a Poisson distribution ( $\lambda = 3$ ) for IDA. Each span is replaced with a single [MASK].