

BALANCING PRECISION AND RICHNESS IN IMAGE CAPTION SERVICES FOR ENHANCED DESCRIPTIVE ACCURACY

Anonymous authors

Paper under double-blind review

ABSTRACT

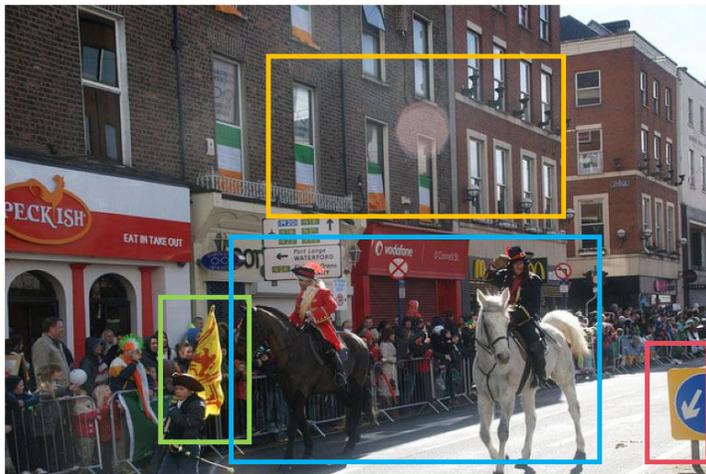
Current image captioning services often learn to generate captions by imitating ground truth references, which are constrained by the limitations of manual annotations. This leads to overlooked details in images, causing captions to lack richness and precise descriptions, critical for enhanced image captioning services. To address this, we propose a CLIP-based image captioning framework designed to balance descriptive precision and richness enhancement. Our approach uses fine-grained pseudo tags for learning and integrates an asymmetric attention multi-modal projector to map and fuse information across modalities effectively. We also introduce an evaluation metric, Tags Coverage, to measure the granularity of generated captions and incorporate it into reinforcement learning to optimize the reward function. This eliminates the need for additional text annotations while addressing unannotated details. Experimental results on the MS-COCO Karpathy’s test set demonstrate the model’s effectiveness, with improvement in CIDEr and Tags Coverage compared to state-of-the-art baselines, highlighting its potential for advancing precision and richness in image captioning services.

1 INTRODUCTION

The field of image captioning has advanced rapidly, emphasizing the integration of cross-modal feature knowledge to support service-driven applications. These services require robust visual perception, natural language generation, and effective multimodal feature alignment to deliver accurate and rich descriptions. Models such as (Vinyals et al., 2015; Ma et al., 2015) utilize CNNs to extract visual features and process them into textual captions. Approaches like Up-Down (Anderson et al., 2018) and Xmodal-Ctx (Kuo & Kira, 2022) enhance precision in descriptions through object detection and context integration. Despite achieving significant metric improvements, existing methods often rely on limited human annotations, leading to captions that lack descriptive richness and fail to balance precision and detail effectively (Rotstein et al., 2024). For real-world service contexts, such as news reporting or visual assistance, solutions must address these limitations by generating captions that incorporate nuanced image details. The focus on balancing precision and richness directly supports scalable service systems for diverse applications in the service computing.

Existing human annotations are often concise, typically describing images in a limited format that simplistically addresses “what is happening in where”. For instance, as illustrated in Fig. 1, annotations usually only hit the main objects in the image, such as “men”, “horse”, “crowd”, and “costume”, yet miss other intricate contextual background details like the “parade”, “houses and windows behind the street”, or “individuals walking in the front with flag”. From the aspect of human perceptions and cognition, rich and diverse contents would be described when human looks at the same image. This does not mean that longer or overly detailed descriptions are better, as too much complexity can lead to confusion. Instead, the captioner should make the generated sentences free from the rigid structure in a more rich and flexible mode. The current features extraction methods and training strategies may lead to perceptual and generative biases. After extracting visual information from the images, captioners often prioritize generating sentences closely imitating the ground truth. However, this approach tends to overlook the inherent details in images, excluding vocabularies that extend beyond the ground truth but accurately capture the image content. Consequently, this results in image captions that lack details, flexibility, and richness.

054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107



Human annotation

- A couple of men riding horses down a street with tall buildings.
- Men riding on horses in street next to buildings.
- A crowd is watching horses go down the street.
- A man dressed in red riding a horse through town.
- people in costume riding down the road on horses.

Hit		Miss		
horse	street	parade	window	flag
crowd	costume	signage	marching	...

Figure 1: An example showcasing the limitations of human annotations - some vocabularies they successfully hit and what they miss reflecting in the corresponding color boxes in the image.

In order to describe images in more detail, we propose a framework for fine-grained pseudo tags, without the additional requirement of text-annotations and task extensions for the overlooked details. Inspired by the Xmodal-CTX (Kuo & Kira, 2022), which introduced a cross-modal retrieval module utilizing CLIP (Radford et al., 2021) to retrieve relevant text descriptions, providing complementary information to objects, we similarly construct relevant tags. Leveraging the data and training advantages of CLIP, we utilize it to extract visual features conditioned to generate the captions of the given image. A set of objects features, and grids features at different details levels extracted by CLIP-I are used. These features complement the visual information for both the main subject and background details with the advantages of preceding models. We construct rich detailed pseudo tags and employ asymmetric attention to adequately capture attention from the visual to the textual side. So that we can attain cross-modal understanding and fusion, which would support the task of generating captions with details.

However, directly utilizing the augmented image information is not sufficient to ensure the richness of the generated image captions. BLIP-2 (Li et al., 2023) and FUSECAP (Rotstein et al., 2024) employ frozen vision encoders in conjunction with powerful large language models to produce highly contextualized and descriptive captions. Nevertheless, the lack of explicit model guidance during the caption generation process makes it challenging to balance precision and richness within a unified end-to-end framework. To this end, we introduce a novel metric, Tags Coverage, to measure the level of details in captions. The constructed pseudo-tags and the metric are applied in NSC (Luo, 2020) for optimizing the reinforcement learning reward function (Rennie et al., 2017). So that, we enhance the richness and details in our generated captions without compromising much on their ac-

108 curacy. The main contributions of this paper are as follows, with a focus on balancing precision and
 109 richness in image captioning:

110 (1) *Enhanced Feature Representation for Image Captioning*: We leverage frozen pre-trained CLIP
 111 to construct pseudo tags, enabling the model to capture intricate image details. To support service-
 112 driven applications, we propose a multi-modal projector module that fuses feature information
 113 across modalities using asymmetric attention, enhancing descriptive precision and richness.

114 (2) *Richness-Oriented Evaluation Metric*: A new Tags Coverage (TC) metric is introduced to quan-
 115 tify caption granularity. TC is integrated into reinforcement learning, aligning training objectives by
 116 balancing accuracy and detail richness in generated captions.

117 (3) *Performance Gains*: Compared to baselines, our model achieves improvement in CIDEr, with a
 118 notable enhancement in TC over human annotations. These gains demonstrate our model’s capabil-
 119 ity to deliver precise and richly detailed image captioning services.

122 2 RELATED WORK

124 2.1 ACCURACY AND FINE-GRAINED IN IMAGE CAPTIONING

125 Traditional metrics like BLEU (Papineni et al., 2002), ROUGE (Lin, 2004), METEOR (Banerjee
 126 & Lavie, 2005), CIDEr (Vedantam et al., 2015), and SPICE (Anderson et al., 2016) effectively as-
 127 sess the accuracy of image captioning by comparing generated captions with reference texts. Early
 128 image captioning models trained via time-wise cross-entropy performed well on these metrics but
 129 faced challenges like exposure bias (Ranzato et al., 2015) and misalignment between training and
 130 evaluation. Reinforcement learning approaches, including self-critical sequence training (Rennie
 131 et al., 2017), have mitigated these issues by using non-differentiable rewards such as CIDEr and
 132 BLEU, stabilizing reward variance and improving accuracy (Luo, 2020; Bujimalla et al., 2020). De-
 133 spite these advances, models often fail to provide rich and flexible captions, overlooking important
 134 details in images. To address this, approaches like using image-text retrieval scores as rewards (Luo
 135 et al., 2018) and integrating CLIP-S (Hessel et al., 2021) have been explored to encourage distinct
 136 and descriptive captions. However, the lack of reference captions in these methods limits their gran-
 137 ularity. In this study, we propose a service-oriented solution that introduces a novel metric, Tags
 138 Coverage (TC), to evaluate caption granularity. By leveraging CLIP (Radford et al., 2021), we re-
 139 trieve fine-grained tags from images and use them as benchmarks to assess and enhance the richness
 140 of captions. This metric is incorporated into the reinforcement learning framework, aligning with
 141 service computing goals to balance descriptive precision and richness. Unlike previous methods,
 142 our approach eliminates reliance on extensive manual annotations or task extensions, ensuring scal-
 143 ability and adaptability for service applications like news reporting or assistive technologies. This
 144 balance between accuracy and granularity advances image captioning as a robust service, offering
 145 enhanced descriptive capabilities.

146 2.2 MODAL INTERACTION IN IMAGE CAPTIONING

147 Image captioning services rely on effective modal interaction to generate accurate and rich de-
 148 scriptions, a critical aspect of service computing. Earlier models, such as attention-based encoder-
 149 decoders (Gu et al., 2018; Wang et al., 2017; Lu et al., 2017), focused on feature extraction within
 150 visual and language domains. The introduction of object detectors (Anderson et al., 2018) signifi-
 151 cantly enhanced the identification of salient regions, improving caption accuracy. Advances in atten-
 152 tion mechanisms (Vaswani et al., 2017; Huang et al., 2019) further enabled the use of Transformers
 153 for capturing global image contexts. \mathcal{M}^2 Transformer (Cornia et al., 2020) utilized meshed mem-
 154 ory storage to enhance encoder-decoder interactions, while the Multi-modal Transformer (Yu et al.,
 155 2019) unified intra- and inter-modal attention blocks. BLIP-2 (Li et al., 2023) and FUSECAP (Rot-
 156 stein et al., 2024) employ frozen vision encoders in conjunction with powerful large language mod-
 157 els to produce highly contextualized and descriptive captions. Building on recent advancements, our
 158 approach leverages frozen pre-trained large-scale models to extract detailed pseudo-labels, includ-
 159 ing object, grid, and tag features. The integration of asymmetric attention enables more effective
 160 mapping and fusion of these features, enriching the core image content with finer semantic details.
 161 Moreover, we use explicit model guidance during the caption generation process, which effectively
 balances precision and richness within a unified framework.

162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215

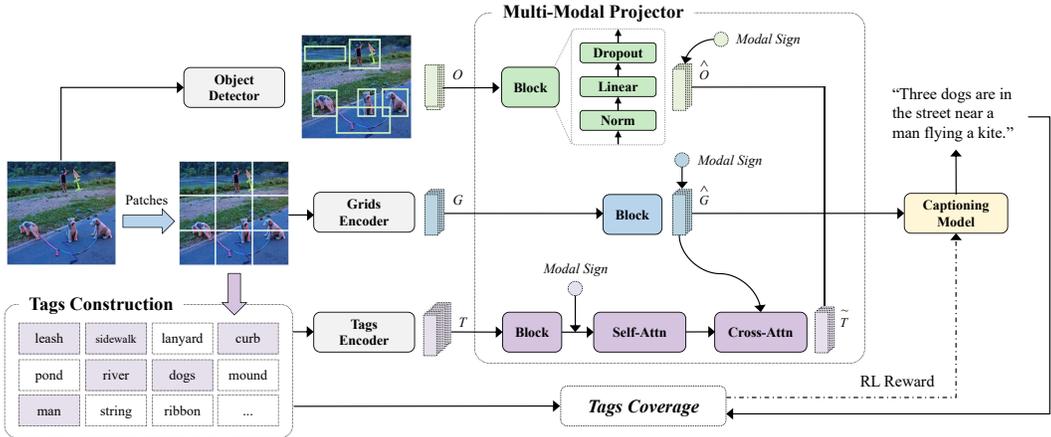


Figure 2: Overview of our proposed model architecture.

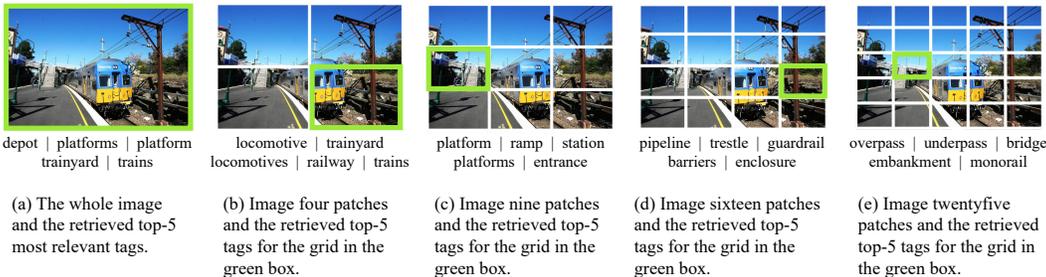


Figure 3: Construct top-5 most relevant tags from different grids for (a) the whole image, (b) image four patches, (c) image nine patches, (d) image sixteen patches, (e) image twenty-five patches.

3 METHOD

We begin with an overview of the model architecture, followed by a introduction of the pseudo-tag construction module and the multi-modal projector. Next, we present the optimized reinforcement learning strategy based on NSC (Luo, 2020), where the proposed Tags Coverage metric is employed as the reward to enhance caption quality.

3.1 MODEL ARCHITECTURE

As shown in Fig. 2, it consists of multi-modal extractors, a multi-modal projector, and a captioning model. We leverage the advantages of existing cross-modal large-scale models to perform features extractors for objects, grids, and tags within images. Specifically, following the methodology proposed in Xmodal-CTX (Kuo & Kira, 2022), we use the same objects features extracted by the object detector pre-train on Visual Genome (Krishna et al., 2017). Employing frozen CLIP-I and CLIP-T from different branches of CLIP (Radford et al., 2021), we extract features for image grids and construct fine-grained pseudo tags. Subsequently, we design a multi-modal projector aiming to fuse the extracted features, thoroughly utilizing self-attention mechanisms and cross-attention mechanisms for sampling tags features. Finally, a captioning model \mathcal{M}^2 Transformer (Cornia et al., 2020) is connected and the fusion of cross-modal features is fed in. The tags constructed before are used for calculating the Tags Coverage metric, which, in turn, will be employed for model fine-tuning, encouraging the generation of image captions containing a more comprehensive set of tags from the tags repository.

3.2 TAGS CONSTRUCTION

The pseudo tags are proposed to enable the model to learn from, beyond human annotations but reflecting images realistically. And it encourages the model to learn without the necessity of addi-

216 tional text-annotations and other vision-language tasks during training. We similarly transform the
 217 tags construction problem into a cross-modal retrieval task, performing specific retrieval for each
 218 sub-region of images from the tags repository. Our goal is to focus on the details of images from
 219 different levels. Through this approach, we can focus on different local regions of the image, thereby
 220 perceiving more detailed local information.

221 The first step is to construct a tags repository for each image. While existing different datasets
 222 contain different annotations for images, to ensure the consistency of the dataset and leverage the
 223 enriched knowledge obtained by large-scale models, we create the tags repository from the ground
 224 truth annotations within the COCO dataset (Lin et al., 2014) (It is also used in our experiments).
 225 After tokenization, deduplication, and cleaning processes, the tags repository is formed using various
 226 vocabulary tokens. Employing the text branch CLIP-T, we encode the tags from the constructed tags
 227 repository as tags repository features T^* .

228 Then, we can retrieve the corresponding tags from different features of sub-regions within im-
 229 ages, based on the cross-modal joint embedding of CLIP. CLIP trained on vast datasets and aligns
 230 images and texts through contrastive learning to ensure textual and visual consistency. We ex-
 231 plicitly consider dividing the image I into one, four, nine, sixteen, and twenty-five patches as
 232 varying region sizes (finer partitions), which would impact the level of details within the im-
 233 age. Using CLIP-I image branch, we encode the patches of the image I into grid features
 234 $G = \{G^{\#1}, G^{\#4}, G^{\#9}, G^{\#16}, G^{\#25}\}$, where $G^{\#i} = \{g_j^{\#i} | j \in \{1, 2, \dots, i\}\}$, and $g_j^{\#i}$ denotes
 235 the j -th grid feature from i sub-patch(es) in the Image I where each $g \in \mathbb{R}^{d_g}$. Utilizing the grids
 236 features $G^{\#i}$ as queries and tags repository T^* as retrieval keywords, we retrieve the top- k tags ac-
 237 cording to the highest cosine similarity scores, use CLIP-T text branch and obtain the tags features
 238 $T = \{T^{\#1}, T^{\#4}, T^{\#9}, T^{\#16}, T^{\#25}\}$, where $T^{\#i} = \{t_{j,k}^{\#i} | j \in \{1, 2, \dots, i\}, k \in \{1, 2, \dots, \text{top-}k\}\}$,
 239 and $t_{j,k}^{\#i}$ denotes the k -th tag feature in top- k tags features corresponding to the j -th grid in i patch(s)
 240 where each $t \in \mathbb{R}^{d_t}$. Some examples of the top-5 results are shown in Fig. 3.

241
 242 After processing mentioned above, we obtain the corresponding tags for each image patch. These
 243 tags will be utilized in two aspects: one will be fed into the model for joint training within the
 244 multi-modal projector, and the other will be used in our proposed metric TC for calculating.

245 3.3 MULTI-MODAL PROJECTOR

246
 247 The multi-modal projector aims to map inputs from different modalities into a shared representation
 248 space to achieve information fusion for cross-modalities. In our work, we utilize frozen pre-trained
 249 models CLIP (Radford et al., 2021) to extract features from image I , including grids features G and
 250 tags features T . Let $O = \{o_1, o_2, \dots, o_n\}$, where each $o \in \mathbb{R}^{d_o}$, represents a set of number n objects
 251 features o_i detected by a frozen pre-trained object detector. Within the multi-modal projector, we
 252 use a Block including Norms, Fcs, and Drops to map features from different modalities to adapt to
 253 downstream sequence generation, which can be modeled as:

$$\begin{aligned} \hat{O} &= [\text{Drop}(\text{FC}(\text{Norm}(O))), m_O], \\ \hat{G} &= [\text{Drop}(\text{FC}(\text{Norm}(G))), m_G], \\ \hat{T} &= [\text{Drop}(\text{FC}(\text{Norm}(T))), m_T], \end{aligned} \quad (1)$$

254
 255 where Norm denotes the normalization layer, FC denotes the fully connected layer, Drop refers
 256 to the dropout layer, $[\cdot, \cdot]$ is concatenation and m denotes a sign from different modals of objects,
 257 grids, and tags. In our works, distinct FC layers are employed to encode various features, addressing
 258 different modalities and granular levels.

259 To thoroughly explore inter-modal information and enhance cross-modal integration for the correla-
 260 tion between grids and tags, we introduce asymmetric attention with summed shortcut connections
 261 and design a self-attention layer (SA) and a cross-attention layer (CA). Specifically, the tags features
 262 \hat{T} , processed through a Block, are input into the self-attention layer to obtain \hat{T}_{SA} . Subsequently,
 263 grid features \hat{G} are injected into the tag features \hat{T}_{SA} via the cross-attention layer to produce \hat{T}_{CA} .
 264 This result is then fed into the Feed-Forward Network (FFN), where it is added to the output of the
 265 cross-attention \hat{T}_{CA} , generating the visual-perceptive tag representation \tilde{T} :

$$\hat{T}_{\text{SA}} = \text{Norm}(\text{SA}(\hat{T}) + \hat{T}), \quad (2)$$

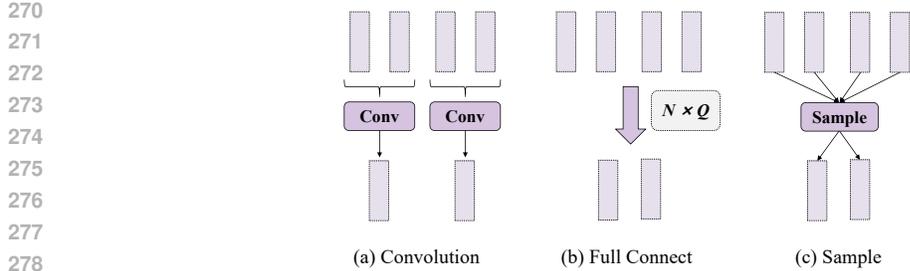


Figure 4: Three types of refinements in cross attention.

$$\hat{T}_{CA} = \text{Norm}(\text{CA}(\hat{T}_{SA}, \hat{G}) + \hat{T}_{SA}), \quad (3)$$

$$\tilde{T} = \text{Norm}(\text{FNN}(\hat{T}_{CA}) + \hat{T}_{CA}), \quad (4)$$

where FFN consists of two linear layers, a ReLU layer, and a Dropout layer.

It is noteworthy that when computing cross-attention, we do not use the entire \hat{T}_{SA} as the query Q . We explicitly guide the model to learn fine-grained tags from different levels, but it does not imply that every tag is beneficial for the model. Meanwhile, to reduce the computational cost, we perform refine the tags features after self-attention \hat{T}_{SA}^* by three ways in Fig. 4. We envision utilizing a fully connected layer (FC) for a linear transformation of all tag features, employing sampling (SP) to extract features from certain tags features, and leveraging convolution (CN) to capture local information of tags features. We employ the scaled dot-product attention, which operates on three sets of vectors. Based on the similarity distribution between the query and key, we calculate the weighted sum of the value:

$$\begin{aligned} Q_{FC} &= W_q \text{FC}(\hat{T}_{SA}), \\ Q_{SP} &= W_q \text{SP}(\hat{T}_{SA}), \\ Q_{CN} &= W_q \text{CN}(\hat{T}_{SA}), \end{aligned} \quad (5)$$

$$\text{CA}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d}}\right)V, \quad (6)$$

where W_q is learnable weights, $Q \in \{Q_{FC}, Q_{SP}, Q_{CN}\}$ is a matrix of tags features after refinements \hat{T}_{SA}^* , K and V both contain \hat{G} keys and values, and d is a scaling factor. Subsequently, the mapped objects features \hat{O} , grids features \hat{G} , and tags features \tilde{T} are concatenated along the sequence dimension to form $z = [\hat{O}, \hat{G}, \tilde{T}]$, which is fed into a captioning model for sequence-to-sequence learning.

3.4 TAGS COVERAGE METRIC

The prevailing metric for assessing the accuracy of generated image captions is CIDEr (Vedantam et al., 2015), which calculates the cosine similarity between the generated and reference sentences based on N-gram. However, higher CIDEr scores may result in generated sentences that are relatively short and lack details. To assess the granularity of generated sentences, one straightforward approach is to examine whether the output captions contain more detailed information. Inspired by (Zhao et al., 2019), we propose a novel metric Tags Coverage (TC), employing the fine-grained tags construction in Section 3.2 as the details benchmark. This metric serves as a precision-like score in Eq. equation 7, quantifying the proportion of tags tokens within the tokens of the output sentence. For each image, we utilize the top-100 tags with the highest cosine similarity which is outlined in the Tags Construction section, as a reference for assessment.

$$\text{TC} = \frac{|\{\text{tags tokens}\} \cap \{\text{caption tokens}\}|}{|\{\text{caption tokens}\}|}. \quad (7)$$

This metric is not independently assessed but relies on the prior knowledge of the pre-trained CLIP (Radford et al., 2021). The accuracy of the pseudo tags constructed by the pre-trained model CLIP

determines the realism of the granularity reflected by TC. However, it is essential to note that covering more tags in the generated sentences does not necessarily mean better quality; the inclusion of function words is also necessary for organizing sentence structures to ensure coherence and fluency. Additionally, some tags that may be misidentified by CLIP are also incorporated into the tags set. It effectively increases the fault tolerance of TC as humans also make mistakes sometimes. Therefore, TC may not reach an exceptionally high score (At least it won't get 100%). As we generate sentences that encompass more tags, the TC score increases, while the CIDEr score may decrease. It serves the purpose of ensuring that generated captions contain more details and should be considered alongside the CIDEr metric to evaluate sentence quality comprehensively. We apply NSC (Luo, 2020), a various form of self-critical sequence training (Rennie et al., 2017). During the reinforcement learning fine-tuning stage, we consider incorporating TC as a reward to balance the accuracy and granularity of the output captions:

$$\begin{aligned} R_{\text{CIDEr}} &= \text{CIDEr}(\hat{y}), \\ R_{\text{TC}} &= \text{TC}(\hat{y}), \end{aligned} \tag{8}$$

$$\begin{aligned} R(\hat{y}) &= \lambda_1 R_{\text{CIDEr}} + \lambda_2 R_{\text{TC}}, \\ \nabla_{\theta} L(\theta) &\approx -(R(\hat{y}) - b) \nabla_{\theta} \log p(\hat{y}|z), \end{aligned} \tag{9}$$

where $R(\cdot)$ is the rewards function, \hat{y} is a sampled caption, b denotes the reward of the baseline for \hat{y} acquired by sampling, and z is the latent variable in Section 3.3. The CIDEr reward is calculated according to the ground-truth sentence y^s , which encourages the model to preserve the content in the input image. The TC reward employs the above-mentioned in Eq. equation 7, encouraging the model to generate finer-grained information from the image. The coefficients λ_1 and λ_2 denote the weights of the two components, which are tunable hyper-parameters. The combined effect of both ensures the accuracy and details.

In addition, as discussed above, the quality of the tags generated by CLIP affects the model's caption generation performance. To address this, we further refine Equation (9) by incorporating the correlation between the generated tag tokens and the caption tokens when computing the tag consistency (TC). The revised formulation is as follows:

$$R'_{\text{TC}} = \text{TC}(\hat{y}) \times S(B(\text{caption tokens}), B(\text{tags tokens})), \tag{11}$$

where $S(\cdot)$ denotes the cosine similarity, and $B(\cdot)$ refers to the BERT-Small version (Turc et al., 2019) used for obtaining vector representations. R'_{TC} aims to allow the model to simultaneously consider the semantic consistency between the CLIP-generated caption and the ground-truth caption, as well as the richness of the generated caption.

3.5 COMPLEXITY ANALYSIS

In comparison to models with only a single encoder, although our approach introduces two encoders (grid encoder and tags encoder), the overall number of training parameters and duration has not increased significantly. This is attributed to the fact that the image segmentation and tag generation we employ are based on pre-trained large models. The operations of the two encoders can be completed during the preprocessing stage of our model, eliminating the need for redundant computations during training. In terms of the complexity of training parameters, the added parameter quantity is 768 x 512 for the grid part and 512 x 512 for the tag part (specific experimental parameters are detailed in Section 4.3). Hence, the complexity of our proposed method remains manageable.

4 EXPERIMENTS

4.1 DATASETS AND METRICS

We train and evaluate our model on the widely used benchmark MS-COCO (Lin et al., 2014) in the field of image captioning. For a fair comparison, we adopted the split method proposed by Karpathy (Karpathy & Fei-Fei, 2015), where each image is associated with a minimum of five manual annotations captions. Of these, 5000 images were used for validation, another 5000 for testing, and the rest for training. According to the standard evaluation have been proposed, mostly based on comparing generated captions to human ones, we utilized a comprehensive set of captioning metrics: BLEU (Papineni et al., 2002), METEOR (Banerjee & Lavie, 2005), ROUGE (Lin, 2004), CIDEr (Vedantam et al., 2015), SPICE (Anderson et al., 2016), along with the proposed TC.

Models	B-1	B-4	M	R	C	S	TC
Att2in (Rennie et al., 2017)	78.3	35.6	27.3	56.9	119.2	20.6	18.7
Up-Down (Anderson et al., 2018)	79.6	36.8	28.0	57.7	121.9	21.4	19.1
Transformer (Vaswani et al., 2017)	80.7	38.9	29.0	58.7	129.2	22.8	19.6
\mathcal{M}^2 (Cornia et al., 2020)	80.8	39.1	29.1	58.4	131.2	22.6	19.5
Xmodal-CTX (Kuo & Kira, 2022)	81.2	40.0	29.2	59.3	133.2	22.8	20.0
HAAV (Kuo & Kira, 2023)	80.3	37.1	28.2	58.0	127.4	21.8	19.3
BLIP-2 (Li et al., 2023)	81.4	39.5	29.6	59.8	133.6	22.3	21.7
FUSECAP (Rotstein et al., 2024)	80.5	39.2	29.2	59.5	133.1	21.9	22.5
Human annotations	/	/	/	/	/	/	17.7
Ours _(CIDEr)	<u>81.5</u>	39.9	29.5	59.4	<u>134.1</u>	<u>23.3</u>	20.0
Ours _(CIDEr+R_{TC})	81.1	39.4	29.1	59.1	132.2	23.2	<u>23.3</u>
Ours _(CIDEr+R_{TC})	81.9	40.6	30.2	60.5	135.2	23.7	23.9

Table 1: Comparison with baselines for image captioning results on the test set of MS-COCO Karpathy split. Human annotations refer to the ground-truths. B-N, M, R, C, S, and TC represent BLEU@N, METEOR, ROUGE-L, CIDEr, SPICE, and Tags Coverage metrics.

Strategy	B-1	B-4	M	R	C	S	TC
XE	76.7	35.0	28.3	56.8	117.8	21.7	20.3
NSC _(w/ CIDEr)	<u>81.5</u>	<u>39.9</u>	<u>29.5</u>	<u>59.4</u>	<u>134.1</u>	<u>23.3</u>	20.0
NSC _(w/ CIDEr+R_{TC})	81.1	39.4	29.1	59.1	132.2	23.2	<u>23.3</u>
NSC _(w/ CIDEr+R_{TC})	81.9	40.6	30.2	60.5	135.2	23.7	23.9

Table 2: The results of our models with and without the reward of TC on the test set. XE and NSC mean the cross-entropy loss training and new self-critical finetuning.

4.2 SELECTED BASELINES

We select the following baselines for comparison: Att2in (Rennie et al., 2017): It introduces a reinforcement learning method for optimizing image captioning systems. It directly optimizes the CIDEr metric, making it highly effective for image captioning tasks. Up-Down (Anderson et al., 2018): It proposes a novel combined bottom-up and top-down attention mechanism, achieving state-of-the-art performance in image captioning and VQA tasks. Transformer (Vaswani et al., 2017): It is a based model for image caption tasks. \mathcal{M}^2 (Cornia et al., 2020): It introduces a meshed transformer with memory for improved image captioning, achieving state-of-the-art performance. Xmodal-CTX (Kuo & Kira, 2022): It introduces a novel approach in visual captioning, utilizing auxiliary inputs to capture missing information and improving model grounding. HAAV (Kuo & Kira, 2023): It proposes an innovative image captioning approach that treats various visual and textual encodings as augmented views of the input image. BLIP-2 (Li et al., 2023): It bridges frozen image encoders and large language models using a lightweight Querying Transformer. FUSECAP (Rotstein et al., 2024) is a data-centric approach that enriches generic image captions by fusing outputs from frozen vision experts with original captions using a large language model.

4.3 EXPERIMENT SETTINGS

We tune the top-k parameter on the validation set and find that the performance saturates at $k = 9$ in the tag construction. The objects features we use follow Xmodal-CTX (Kuo & Kira, 2022). The grids and tags features extractors we used are frozen pre-trained CLIP. The dimensions of the extracted object features, grid features, and tag features are 2054, 768, and 512 respectively. After the multi-modal projector, they are fused to the latent variable z with a dimension of 512. We train our model with Adam optimization and the Reduce-LR-On-Plateau method on A800. The model is initially trained using cross-entropy XE for 25 epochs with the learning rate of 1×10^{-4} , followed by fine-tuning with NSC (Luo, 2020) reinforcement learning with appropriate rewards for another 15 epochs. All comparisons among experimental methods were conducted under fair conditions. More experimental details are provided in the Appendix section.

4.4 COMPARISON FOR IMAGE CAPTIONING RESULTS

Firstly, we compare our model with the trained-from-scratch methods as shown Table 1. Except for \mathcal{M}^2 (Cornia et al., 2020), which is obtained from their published model, the other models are

432
433
434
435
436
437
438

Asymmetric Attention	XE	NSC	B-1	B-4	C	S	TC
w/o attn	✓		76.9	34.7	115.8	21.3	20.3
w/ attn	✓		76.7	35.0	117.8	21.7	20.3
w/ multi-attn	✓		77.3	35.3	118.4	21.9	20.2
w/ attn		✓	81.9	40.6	135.2	23.7	23.9
w/ multi-attn		✓	82.3	41.2	135.9	24.1	24.3

439
440
441
442

Table 3: Ablation study for the proposed asymmetric attention in the multi-modal projector. The first row without attention means features after extractors are directly fully connected to the next part in the model. The multi-attn refers to multi-head attention with heads of 8 here.

443
444
445
446
447
448

reproduced in the codes framework (Luo et al., 2018). We mark the best scores in bold and the second with the underline. Our method achieves the best overall performance across all metrics, including 135.2 CIDEr, 23.7 SPICE, and 23.9 TC. Compared to the strongest baseline FUSECAP, our model improves CIDEr by 2.1, SPICE by 1.8, and TC by 1.4. It also outperforms BLIP-2 by 1.6 in CIDEr, 1.4 in SPICE, and 2.2 in TC, while maintaining stronger or comparable performance on BLEU, METEOR, and ROUGE-L metrics.

449
450
451
452
453
454
455
456
457

Secondly, we observe that Ours_(CIDEr) performs well on standard evaluation metrics, achieving 134.1 in CIDEr and 23.3 in SPICE, while maintaining a moderate TC score of 20.0. This outcome aligns with its design objective of prioritizing caption accuracy, though it may generate less detailed descriptions. In comparison, Ours_(CIDEr+R_{TC}) shows a slight decrease in CIDEr (132.2) and SPICE (23.2), but achieves a higher TC score, 3.3 points above FUSECAP and 5.6 points above human annotations, suggesting improved caption richness. Furthermore, Ours_(CIDEr+R'_{TC}) achieves the highest scores across all metrics, indicating that jointly optimizing for semantic consistency and richness leads to more balanced captions in terms of both accuracy and detail. These results demonstrate the effectiveness of incorporating a dual-objective reward in guiding caption generation.

458
459

4.5 STRATEGY ANALYSIS OF THE PROPOSED METHOD

460
461
462
463
464
465

From the results in Table 2, it is observed that the model optimized through CIDEr reward may incur some losses in terms of fine-grained metric. Specifically, to achieve higher CIDEr scores, the model tends to generate results closer to the ground truths, limiting the generated vocabulary within a certain range defined by human annotations. Our TC reward encourages the model to explore more possibilities, leading to a broader coverage of tags details in the generated vocabulary. Thus, using TC reward does help to cover more detailed tags in image captioning.

466

4.6 ABLATION STUDY

467
468
469
470
471
472
473
474
475
476
477
478

We conducted an ablation analysis on our asymmetric attention multi-modal projector, as shown in the first row “w/o attn” of Table 3. In this setting, we directly input the extracted objects features, grids features, and tags features through the Block into \mathcal{M}^2 for training. It can be observed that the multi-modal projector with attention improves 2.0 on CIDEr when using the strategy of XE. After applying the asymmetric attention projection across modalities, features from different modalities can better integrate information to meet the requirement of the downstream task of generating more fine-grained descriptions. To verify it, we further enhanced attention by incorporating a multi-head attention mechanism with $h = 8$. Compared with attention, we observe a further performance improvement - 0.7 in CIDEr and 0.4 in TC with multi-head attention after reinforcement learning fine-tuning, indicating the effectiveness of our multi-modal projector. Due to the higher computational cost of multi-head attention mechanism, we did not opt for this set.

479

5 CONCLUSION

480
481
482
483
484
485

In this paper, we address the challenge of balancing precision and richness in image captioning. We propose a CLIP-based model that captures fine-grained tags by extracting object, grid, and tag features and integrating them with an asymmetric attention projector. To encourage accurate yet detailed captions, we introduce TC, a fine-grained evaluation metric, into the reinforcement learning reward. Experiments show that our method generates captions with both precision and richness.

REFERENCES

- 486
487
488 Peter Anderson, Basura Fernando, Mark Johnson, and Stephen Gould. Spice: Semantic proposi-
489 tional image caption evaluation. In *Proceedings of the European conference on computer vision*,
490 pp. 382–398. Springer, 2016.
- 491 Peter Anderson, Xiaodong He, Chris Buehler, Damien Teney, Mark Johnson, Stephen Gould, and
492 Lei Zhang. Bottom-up and top-down attention for image captioning and visual question answer-
493 ing. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp.
494 6077–6086, 2018.
- 495 Satanjeev Banerjee and Alon Lavie. Meteor: An automatic metric for mt evaluation with improved
496 correlation with human judgments. In *Proceedings of the ACL workshop on intrinsic and extrinsic*
497 *evaluation measures for machine translation and/or summarization*, pp. 65–72, 2005.
- 499 Shashank Bujimalla, Mahesh Subedar, and Omesh Tickoo. B-scst: Bayesian self-critical sequence
500 training for image captioning. *arXiv:2004.02435*, 2020.
- 501 Marcella Cornia, Matteo Stefanini, Lorenzo Baraldi, and Rita Cucchiara. Meshed-memory trans-
502 former for image captioning. In *Proceedings of the IEEE/CVF conference on computer vision*
503 *and pattern recognition*, pp. 10578–10587, 2020.
- 505 Jiuxiang Gu, Jianfei Cai, Gang Wang, and Tsuhan Chen. Stack-captioning: Coarse-to-fine learn-
506 ing for image captioning. In *Proceedings of the association for the advancement of artificial*
507 *intelligence*, volume 32, 2018.
- 508 Jack Hessel, Ari Holtzman, Maxwell Forbes, Ronan Le Bras, and Yejin Choi. Clipscore: A
509 reference-free evaluation metric for image captioning. *arXiv:2104.08718*, 2021.
- 511 Lun Huang, Wenmin Wang, Jie Chen, and Xiao-Yong Wei. Attention on attention for image cap-
512 tioning. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 4634–
513 4643, 2019.
- 514 Andrej Karpathy and Li Fei-Fei. Deep visual-semantic alignments for generating image descrip-
515 tions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp.
516 3128–3137, 2015.
- 517 Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie
518 Chen, Yannis Kalantidis, Li-Jia Li, David A Shamma, et al. Visual genome: Connecting lan-
519 guage and vision using crowdsourced dense image annotations. *International journal of computer*
520 *vision*, 123:32–73, 2017.
- 522 Chia-Wen Kuo and Zsolt Kira. Beyond a visual object detector: Cross-modal textual and visual
523 context for image captioning. In *Proceedings of the IEEE/CVF conference on computer vision*
524 *and pattern recognition*, pp. 17969–17979, 2022.
- 525 Chia-Wen Kuo and Zsolt Kira. Haav: Hierarchical aggregation of augmented views for image cap-
526 tioning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*,
527 pp. 11039–11049, 2023.
- 528 Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image
529 pre-training with frozen image encoders and large language models. In *International conference*
530 *on machine learning*, pp. 19730–19742. PMLR, 2023.
- 532 Chin-Yew Lin. Rouge: A package for automatic evaluation of summaries. In *Text summarization*
533 *branches out*, pp. 74–81, 2004.
- 534 Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr
535 Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Proceedings of*
536 *the European conference on computer vision*, pp. 740–755. Springer, 2014.
- 538 Jiasen Lu, Caiming Xiong, Devi Parikh, and Richard Socher. Knowing when to look: Adaptive
539 attention via a visual sentinel for image captioning. In *Proceedings of the IEEE conference on*
computer vision and pattern recognition, pp. 375–383, 2017.

- 540 Ruotian Luo. A better variant of self-critical sequence training. *arXiv:2003.09971*, 2020.
- 541
- 542 Ruotian Luo, Brian Price, Scott Cohen, and Gregory Shakhnarovich. Discriminability objective for
543 training descriptive captions. *arXiv:1803.04376*, 2018.
- 544
- 545 Lin Ma, Zhengdong Lu, Lifeng Shang, and Hang Li. Multimodal convolutional neural networks for
546 matching image and sentence. In *Proceedings of the IEEE international conference on computer
547 vision*, pp. 2623–2631, 2015.
- 548 Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic
549 evaluation of machine translation. In *Proceedings of the annual meeting of the Association for
550 Computational Linguistics*, pp. 311–318, 2002.
- 551 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,
552 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual
553 models from natural language supervision. In *Proceedings of the International conference on
554 machine learning*, pp. 8748–8763. PMLR, 2021.
- 555 Marc’Aurelio Ranzato, Sumit Chopra, Michael Auli, and Wojciech Zaremba. Sequence level train-
556 ing with recurrent neural networks. *arXiv:1511.06732*, 2015.
- 557
- 558 Steven J Rennie, Etienne Marcheret, Youssef Mroueh, Jerret Ross, and Vaibhava Goel. Self-critical
559 sequence training for image captioning. In *Proceedings of the IEEE conference on computer
560 vision and pattern recognition*, pp. 7008–7024, 2017.
- 561
- 562 Noam Rotstein, David Bensaid, Shaked Brody, Roy Ganz, and Ron Kimmel. Fusecap: Leveraging
563 large language models for enriched fused image captions. In *Proceedings of the IEEE/CVF winter
564 conference on applications of computer vision*, pp. 5689–5700, 2024.
- 565 Iulia Turc, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Well-read students learn better:
566 On the importance of pre-training compact models. *arXiv preprint arXiv:1908.08962v2*, 2019.
- 567
- 568 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,
569 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Proceedings of the Advances
570 in neural information processing systems*, 30, 2017.
- 571 Ramakrishna Vedantam, C Lawrence Zitnick, and Devi Parikh. Cider: Consensus-based image
572 description evaluation. In *Proceedings of the IEEE conference on computer vision and pattern
573 recognition*, pp. 4566–4575, 2015.
- 574 Oriol Vinyals, Alexander Toshev, Samy Bengio, and Dumitru Erhan. Show and tell: A neural
575 image caption generator. In *Proceedings of the IEEE conference on computer vision and pattern
576 recognition*, pp. 3156–3164, 2015.
- 577
- 578 Yufei Wang, Zhe Lin, Xiaohui Shen, Scott Cohen, and Garrison W Cottrell. Skeleton key: Im-
579 age captioning by skeleton-attribute decomposition. In *Proceedings of the IEEE conference on
580 computer vision and pattern recognition*, pp. 7272–7281, 2017.
- 581 Jun Yu, Jing Li, Zhou Yu, and Qingming Huang. Multimodal transformer with multi-view visual
582 representation for image captioning. *IEEE transactions on circuits and systems for video technol-
583 ogy*, 30(12):4467–4480, 2019.
- 584
- 585 Sanqiang Zhao, Piyush Sharma, Tomer Levinboim, and Radu Soricut. Informative image captioning
586 with external sources of information. *arXiv:1906.08876*, 2019.
- 587
- 588
- 589
- 590
- 591
- 592
- 593

Dimension Reduction	B-1	B-4	C	S	TC
CN	76.0	34.0	114.0	21.1	20.4
FC	76.7	35.0	117.8	21.7	20.3
SP	77.1	35.1	116.9	21.4	20.3

Table 4: Different ways of dimension reduction in asymmetric cross attention. CN, FC, and SP refer to convolution, fully connect, and sample, see text for detailed descriptions of three ways.

A APPENDIX

Several additional experiments and analyses are provided in the appendix, including different refinements on tag features, studies on the balance between accuracy and fine-grained rewards, evaluations under varying top- k settings, and qualitative results with visualizations comparing model outputs and human annotations. We commit to releasing the complete version of the dataset and code at the acceptance stage.

A.1 DIFFERENT TYPES OF REFINEMENTS ON TAGS FEATURES

We performed refinements on the processed tags features to reduce computational costs and enhance attention efficiency. Specifically, we explored three dimension reduction strategies in the asymmetric cross-attention module: convolution (CN), fully connected (FC), and sampling (SP). Table 4 presents the performance comparison across these approaches using standard evaluation metrics.

The FC method shows the highest performance in CIDEr (117.8), SPICE (21.7), and BLEU-4 (35.0), indicating its effectiveness in integrating information from high-dimensional tag features. This can be attributed to the capacity of fully connected layers to model comprehensive feature interactions, which contributes to improved caption accuracy. The SP method achieves slightly lower scores than FC in CIDEr and BLEU-4 but attains the highest BLEU-1 score (77.1), suggesting it may help preserve more diverse n-gram expressions by selectively retaining tag features. Although its TC score (20.3) is marginally lower than that of CN (20.4), the performance remains competitive. The CN approach yields the lowest CIDEr score (114.0) among the three, but it performs best in the TC metric (20.4), implying that convolutional operations may help retain fine-grained details relevant to tag coverage, despite producing slightly less accurate captions overall.

In general, FC is recommended when accuracy is prioritized in caption generation, while SP can be a viable option for scenarios that emphasize diversity and flexibility. CN serves as a balanced alternative, offering modest performance in both accuracy and richness.

A.2 BALANCE WEIGHTS BETWEEN ACCURACY AND FINE-GRAINED DESCRIPTIONS

We examine the impact of varying the reward weights between CIDEr and TC (R_{TC} and R'_{TC}) during model fine-tuning, as shown in Table 5.

When the weight of TC increases while keeping the CIDEr weight fixed at 1, a moderate trade-off is observed: the CIDEr score decreases slightly (e.g., from 134.1 to 132.2 with R_{TC}), while the TC score improves notably (from 20.0 to 23.3). Besides, in the case of the refined reward R'_{TC} , both accuracy and richness improve together, with CIDEr increasing to 135.2 and TC reaching 23.9 when the weights are balanced. These findings indicate that incorporating both semantic consistency and richness into the reward formulation can help produce captions that maintain accuracy while providing more detailed descriptions. Careful tuning of reward weights is therefore important for achieving a balance between precision and informativeness in image captioning.

A.3 TOP-K EXPERIMENTAL RESULTS

We further evaluated our model by varying the hyper-parameter top- k in *cross-entropy training* (XE). As depicted in Table 7, the optimal results in terms of CIDEr, SPICE, and Tags Coverage were achieved when utilizing top-9 tags features as the parameter. The accuracy of generated captions, measured by metrics like CIDEr, gradually improved with an increase in top- k and reached optimal

648
649
650
651
652
653
654
655
656
657

Weight of CIDEr	Weight of R_{TC}	B-1	B-4	C	S	TC
1	0	81.5	39.9	134.1	23.3	20.0
1	0.5	81.6	40.0	133.3	23.2	21.3
1	1	81.1	39.4	132.2	23.2	23.3

Weight of CIDEr	Weight of R_{TC}	B-1	B-4	C	S	TC
1	0	81.5	39.9	134.1	23.3	20.0
1	0.5	81.7	40.1	134.7	23.4	22.5
1	1	81.9	40.6	135.2	23.7	23.9

658
659

Table 5: Image captioning results of different rewards on MS-COCO Karpathy test split.

660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675

Hyperparameter	XE	NSC	Note
num_layers	3	*	Number of encoder and decoder layers
encoder_size	512	*	encoder embedding size
ff_size	2048	*	feed-forward network size
h	8	*	head of multi-head attention
dropout_rate	0.1	*	dropout rate
max_lenth	20	*	max length of captions
vocab_size	9487	*	length of vocabulary
max_epochs	25	+15	training epochs
lr	1e-4	5e-6	learning rate
optimizer	Adam	*	Adam optimizer
objs_size	2054	*	objects features size
grids_size	768	*	grids features size
tags_size	512	*	tags features size
tags_num	100	*	number of tags to calculate Tags Coverage

676
677

Table 6: Hyperparameters for cross entropy (XE) training and new self-critical (NSC) training. The values for untuned parameters are inherited from the base image captioning model.

678
679
680
681
682
683
684
685
686
687
688
689
690

	B-1	B-4	M	R	C	S	TC
3	76.9	34.6	27.9	56.6	116.1	21.2	20.3
4	76.6	34.5	27.9	56.5	116.3	21.5	20.2
5	76.6	34.5	28.0	56.6	116.3	21.4	20.5
6	76.8	34.9	27.9	56.8	116.7	21.7	20.2
7	76.8	34.6	27.9	56.6	115.9	21.5	20.2
8	76.9	34.4	27.9	56.5	115.9	21.4	20.5
9	76.7	35.0	28.3	56.8	117.8	21.7	20.3
10	76.9	35.0	28.0	56.5	117.1	21.4	20.3
11	77.1	35.0	28.1	56.8	117.1	21.5	20.2
12	76.7	34.9	27.8	56.6	115.7	21.3	20.0

691
692
693

Table 7: Captioning results (trained with XE loss) of our model with top- k tags features used: from top-3 to top-12. B- N , M, R, C, S, and TC represent BLEU@ N , METEOR, ROUGE-L, CIDEr, SPICE, and Tags Coverage metrics.

694
695
696
697
698
699
700
701

performance around top-9. Therefore, in the main paper, we conducted experiments using top-9 as our parameter of the number of tags features used in our model. The table illustrates that, as the number of tag features increases, the model can glean more knowledge from the rich tag information, thereby enhancing the accuracy of the model in generating captions and improving the Tags Coverage metric. However, it’s crucial to note that a higher top- k parameter is not necessarily better. On one hand, we cannot ensure that all tags inputted into the model are perfectly accurate,

702		Ours A cat standing on a sidewalk next to <u>a flock of pigeons</u>		Ours A group of children standing in a field with <u>soccer balls</u>
703		Human A very cute cat near a bunch of birds		Human A group of young children standing around a field
704				
705				
706				
707				
708				
709		Ours A <u>white</u> horse <u>pulling a carriage with people</u> in the grass		Ours A crowd of people walking down a street with <u>a stop sign</u>
710		Human A horse and buggy that is on a grassy field		Human People stand in a city street at a rally
711				
712				
713				
714				
715		Ours A woman <u>wearing sunglasses</u> looking out of an airplane window		Ours <u>A group of surfers</u> riding a wave on surfboards <u>in the ocean</u>
716		Human A woman sleeping on a plane with a window view of the wing		Human A person is riding the waves on a surfboard.
717				
718				
719				
720				
721				
722		Ours A <u>black and white</u> cat sitting on a desk next to a computer		Ours A crowd of people walking in front of <u>a building with a clock tower</u>
723		Human A can laying on a desk in front of a computer		Human A crowd of people walking in an outdoor fair
724				
725				
726				
727				

Figure 5: Examples of image captioning services generated by our model and the ground-truths.

as label construction may have inherent errors. On the other hand, an excess of tags may introduce noise, preventing the model from focusing on the essential information it should learn.

A.4 IMPLEMENTATION DETAILS

We provide a list of hyper-parameters including their values during cross-entropy training(XE) and new self-critical training (NSC) in Table 6. Others not present are following the works before. For cross-entropy training, the model can be trained with three Nvidia 3090 GPUs. For NSC training, the model can be trained with a single A800 GPU.

A.5 QUALITATIVE RESULTS AND VISUALIZATION

Fig. 5, Fig. 6 and Fig. 7 present the qualitative results obtained by our model and the original human annotations. The portions of the captions that subjectively represent detailed tags information are highlighted in purple and underlined. For instance, the colors of the horse and the cat in images are caught by our model. On average, our model covers more fine-grained details and object relationships, producing descriptions with both high accuracy and details.

756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809

	<p>Ours Two people in the ocean one is riding a wave on a surfboard</p> <p>Human A person riding a surf board on a body of water</p>		<p>Ours A large airplane parked at an airport with people walking by</p> <p>Human A airplane that is sitting on a tarmac</p>
	<p>Ours A large passenger jet sitting on top of an airport tarmac</p> <p>Human An airplane parked on the tarmac at an airport</p>		<p>Ours A boat filled with people floating on top of a body of water</p> <p>Human A bunch of people are on a small boat</p>
	<p>Ours A person holding a glass filled with sugar covered donuts</p> <p>Human A cup of food in a persons hand</p>		<p>Ours A white cat sitting in a bathtub next to a red and white rug</p> <p>Human A cat sitting in a bathtub behind the curtain</p>
	<p>Ours A truck parked on the beach with people walking on the sand</p> <p>Human A live guard truck parked on a beach</p>		<p>Ours A display case in a bakery filled with donuts and pastries</p> <p>Human a bakery with boxes of donuts and bread</p>
	<p>Ours A man is petting an elephant that is standing next to a log</p> <p>Human A man is reaching over to an elephant</p>		<p>Ours A group of people standing outside of a double decker bus</p> <p>Human A group of people are on the grass by busses</p>

Figure 6: Examples of image captioning services generated by our model and the ground-truths.

810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863

	<p>Ours A young boy laying on a couch holding a nintendo wii game controller</p> <p>Human A sleeping child holding a Wii controller in hand</p>		<p>Ours A pizza sitting on a wooden cutting board with a bottle of wine</p> <p>Human A pizza sitting on top of a wooden cutting board</p>
	<p>Ours Three dogs sitting next to a fruit stand with fruits and vegetables</p> <p>Human Three dogs sitting side by side in the street</p>		<p>Ours A laptop computer sitting on a desk with a glass of orange juice</p> <p>Human A laptop on a wooden table near a drink</p>
	<p>Ours Two young boys sitting in a baseball dugout with baseball bats</p> <p>Human A couple of men standing next to each other</p>		<p>Ours A man jumping in the air on a skateboard over a fire hydrant</p> <p>Human A man that is in the air with a skateboard</p>
	<p>Ours A construction truck with two traffic lights on a bridge</p> <p>Human A street scene with a large truck driving by</p>		<p>Ours A red and yellow train on the tracks in front of a building</p> <p>Human A photo of a train passing by a building</p>
	<p>Ours A woman standing in a field with a herd of sheep and a dog</p> <p>Human A man is with some sheep in a field</p>		<p>Ours A group of people standing in the snow holding snowboards</p> <p>Human A group of four people standing next to each other in the snow</p>

Figure 7: Examples of image captioning services generated by our model and the ground-truths.