

A-CAP: Anticipation Captioning with Commonsense Knowledge

Duc Minh Vo

The University of Tokyo, Japan

vmduc@nlab.ci.i.u-tokyo.ac.jp

Akihiro Sugimoto

National Institute of Informatics, Japan

sugimoto@nii.ac.jp

Quoc-An Luong

The Graduate University for Advanced Studies, Japan

lqan@nii.ac.jp

Hideki Nakayama

The University of Tokyo, Japan

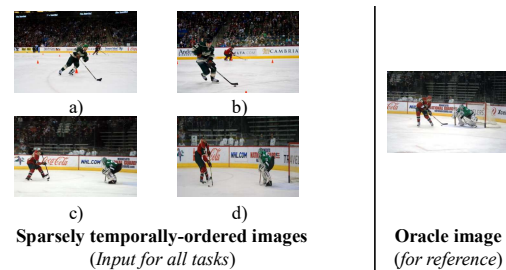
nakayama@ci.i.u-tokyo.ac.jp

Abstract

Humans possess the capacity to reason about the future based on a sparse collection of visual cues acquired over time. In order to emulate this ability, we introduce a novel task called **Anticipation Captioning**, which generates a caption for an unseen oracle image using a sparsely temporally-ordered set of images. To tackle this new task, we propose a model called A-CAP, which incorporates commonsense knowledge into a pre-trained vision-language model, allowing it to anticipate the caption. Through both qualitative and quantitative evaluations on a customized visual storytelling dataset, A-CAP outperforms other image captioning methods and establishes a strong baseline for anticipation captioning. We also address the challenges inherent in this task.

1. Introduction

When humans observe the real world, we not only capture visual information (e.g. objects), but also forecast the future from past and current observations. For example, in Fig. 1, given some photos of an attack in a hockey game, we can predict without a doubt that “the athlete will shoot the puck toward the goalie”. In fact, anticipatory ability aids us in surviving in a world of volatility. This ability necessitates a significant shift from visual to cognitive understanding, which extends far beyond the scope of tasks that primarily use visible visual data, such as object detection, action recognition, and existing image captioning. As a result, a variety of new tasks have been proposed to emulate humans’ anticipatory ability, such as generating future images [12,29], and action prediction [22,37]. Despite their great success, the aforementioned tasks frequently involve densely temporal information (i.e., video), which can be difficult to acquire at times, and their outcomes are not friendly to everyone, particularly those with visual impairments.



Task	Output(s) for input images	Output for oracle image
Image captioning	a) A man standing with a hockey stick. b) A man standing on an ice rink. c) A man holding a hockey stick. d) A group of men playing a game of hockey.	(N/A)
Story telling	a, b, c, d) This breakaway was the first threat to score. The wingman took the puck to the goal but a nice play by the goalie saved the goal. Finally, the other team gets the puck deep into red zone. He is now within 20 feet of the goal.	(N/A)
Anticipation captioning (Ours)	(N/A)	He shoots, he scores and the game ends one to nothing.

Figure 1. Given a set of sparsely temporally-ordered images (a, b, c, d), image captioning [38] and storytelling [35] tasks generate captions for those images, while our anticipation captioning task anticipates what happens afterward. To illustrate the potential future, we show their related oracle image. It should be noted that our task only receives the same inputs as others.

In this work, we hope to dislodge the time constraints imposed by previous tasks while also looking for a more user-friendly output format. Needless to say, textual description is a potential candidate because generating text from images has been successfully explored in a variety of ways [6, 14, 21, 33, 35, 38], showing a number of applications. Furthermore, we can easily leverage recent advances in text-to-image [28] or text-to-sound [36] as a flexible transformation that will benefit other downstream tasks,

allowing everyone to fully grasp our output in their own way. With this in mind, we go beyond the immediately visible content of the images, proposing a new task of image captioning problems, called **anticipation captioning**. Anticipation captioning is to generate a caption for an unseen image (referred to as the oracle image) that is the future of a given set of sparsely temporally-ordered images. The term “sparse” means that two consecutive images are not required to be as close in time as those in a video, allowing the scene to change freely as long as the change does not disrupt the information flow of the image sequence, as seen in Fig. 1. Our task is a new branch of the image captioning problems [6, 14, 21, 35, 38]; it is to predict only captions in the future. As an example, we depict the outputs of generic image captioning, visual storytelling, and our task in Fig. 1. The image captioning model [38] generates a caption for each individual image, whereas visual storytelling [35] connects all images in a story. Our task, on the other hand, produces a caption for the oracle image that is similar to human anticipation: “he shoots, he scores, and the game ends one to nothing”. Unlike [12, 22, 29, 37], anticipation captioning does not require strictly temporal information while producing a more informative output. In theory, the greater the success of this task, the greater the deductive ability of the intelligent system. Meanwhile, other applications such as incident prevention or behavior prediction can be launched.

Additionally, we propose a baseline model, A-CAP, to solve this new task rather than simply using current image captioning models, given their failures in predicting the future. We hypothesize that under common thinking, the future can be predicted from observable concepts (e.g., objects, events) appearing in the input images, implying that the future cannot be dramatically changed to the “football scene” from the “hockey scene”, for instance. As a result, we make full use of commonsense knowledge to connect all detected concepts in terms of a graph while expanding the graph toward forecasted ones, creating a knowledge graph. The term “forecasted concept” refers to a concept that is not visible in the given image but related to another concept visible in the image (we can infer the forecasted concept from the related concept using common thinking). Technically, each node in our constructed graph is either a detected concept in given inputs or a forecasted one explored using the ConceptNet [30], and nodes are connected if and only if they have corresponding ConceptNet relations. After aggregating all node information with a graph neural network, we use prompt learning [39, 40] to integrate the enriched nodes into a frozen pre-trained vision-language (VL) model, successfully generating the anticipated caption. The following are our primary contributions.

- We introduce a novel task of anticipation captioning, which predicts a caption for the future from a given set of sparsely temporally-ordered images.

- For anticipation captioning, we establish a strong baseline model A-CAP, which incorporates commonsense knowledge into a pre-trained VL model.

We evaluate the effectiveness of A-CAP in both qualitative and quantitative ways, using a customized VIST dataset [14]. Extensive experiments show that A-CAP successfully generates captions for oracle images that are more accurate, descriptive, and reasonable than those generated by other captioning methods [35, 38].

2. Related work

Future forecasting has long been studied in computer vision. Some attempts [12, 16, 29, 34] have been made to generate future images/frames from a given video (i.e., dense time-series images). Meanwhile, some methods [22, 37] use past observations to predict future events. These methods heavily rely on the dense temporal-structure to learn visual representations, implying that such representations are different from those for sparsely temporally-ordered images. Furthermore, generated images/frames are not always of high quality [12, 16, 29, 34], and the set of predicted future events is limited [22, 37], making them difficult to apply to downstream tasks. Our method, on the other hand, accepts only sparsely temporal information as long as we can detect objects/events. Furthermore, our method is designed to generate textual descriptions that are easier to interpret than outputs by other methods [12, 16, 22, 29, 34, 37].

In NLP, there are also several approaches to predict the future: story ending generation [7, 18], temporal order anticipation [23, 24]. Though those methods use texts as inputs while our method uses images, we can think of story ending generation as an indirect way to solve our problem because we can generate a story first and then predict its ending.

Image captioning is a long-standing problem with numerous methods developed to address various purposes. Captioning models [6, 21] in an early stage aim to generate generic descriptions for given images. They are then evolved in various directions to generate dense captions [15], novel object captions [33], controllable captions [9], or visual story telling [8, 14, 35]. Anticipation captioning belongs to the image captioning family, with the exception that we predict a caption for the future. Furthermore, our method is based on recent methods [33, 38], which use a vision-language model to generate better captions. Rather than fine-tuning or retraining the model, we use prompt learning [39, 40] to replace the object tags used in the concatenated sequence of words—object tags—ROIs of VinVL [38] with our detected and forecasted concepts.

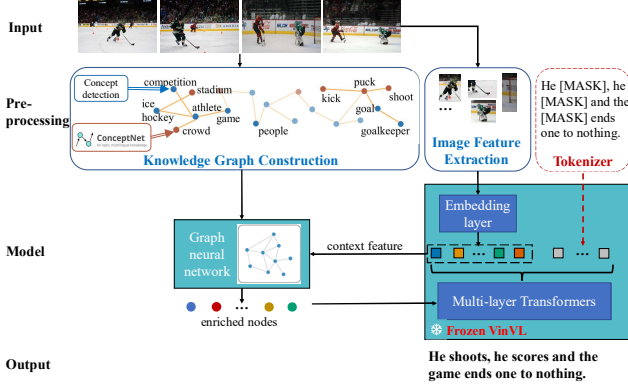


Figure 2. The overall pipeline of our proposed A-CAP. The pre-processing step is used to build the knowledge graph, extract image features and tokenize the input words. In the knowledge graph construction, blue nodes represent the detected concepts obtained from concept detection while brown nodes represent the forecasted concepts obtained from the ConceptNet. Our network consists of a trainable graph neural network and a frozen pre-trained VinVL [38]. The outputs of the graph neural network are the enriched nodes of the knowledge graph. During inference time, the dash-dotted red part is removed.

3. Our approach

3.1. Problem statement

Our input is a set of k sparsely temporally-ordered images I_1, I_2, \dots, I_k . It is worth noting that I_i and I_{i+1} are not necessarily strongly temporal as illustrated in Fig. 1. We assume that an image I_{k+1} is an oracle image that continues the set of k images, and that a caption C_{k+1} corresponds to I_{k+1} which is a future of I_1, I_2, \dots, I_k . Obviously, the oracle image is sparsely temporally-ordered with respect to the input images as we intentionally seek to anticipate the future.

Our task is to generate caption C_{k+1} using given k images. The task is formally defined as follows:

$$C_{k+1} = \text{CAPTION}(I_1, I_2, \dots, I_k), \quad (1)$$

where $\text{CAPTION}(\cdot)$ is a captioning system that will be discussed later. Note that we produce neither captions for each input image I_1, \dots, I_k nor oracle image I_{k+1} .

3.2. Proposed A-CAP

3.2.1 Design of A-CAP

Given the progress of vision-language models in image captioning tasks, we choose VinVL [38] as our base architecture. VinVL takes a concatenated sequence of words-concepts-ROIs as input (note that words are not used during inference time; object tags are used instead of concepts in the original paper [38]). The core idea is the usage of

concepts, which allows better alignment between the vision and language spaces. The above observation suggests that incorporating forecasted concepts into VinVL is critical in allowing the model to generate the anticipated caption. However, simply using VinVL is not wise because it detects only concepts appearing in images. We thus find forecasted concepts based on the detected concepts. Under normal circumstances, forecasted concepts should be related to current observable concepts. Therefore, to retrieve forecasted concepts, we use commonsense knowledge, which consists of many popular concepts and their relationships.

VinVL [38] is trained on a very large dataset, making fine-tuning or re-training difficult. To avoid this difficulty, we use the prompt learning technique to train the concept embeddings only while other parameters are fixed. In what follows, we detail our model.

3.2.2 Network architecture

We base A-CAP on the VinVL [38] architecture. As discussed above, we use concepts as a prompt to allow the model to generate a desired caption. We can then focus on learning the embeddings for all detected and forecasted concepts. To this end, we first retrieve the forecasted concepts using the detected ones and then construct the knowledge graph that connects all concepts. This is because the graph structure is effective for learning the interactions between concepts. We use an undirected graph for simplicity where two concepts are connected as long as their relationship exists. The concept embeddings are then enhanced using a graph neural network. Next, the enriched concept embeddings are injected into a frozen VinVL to generate the caption. Fig. 2 depicts our simple yet effective A-CAP.

3.2.3 Modules of A-CAP

Pre-processing. The input images are pre-processed to (i) construct the knowledge graph and (ii) extract image features. We also tokenize the ground-truth captions used to train the model during training. We obtain N features (ROIs) with the size of 1×2054 each after image feature extraction using Faster-RCNN [27] trained on the COCO dataset. Each image feature is fed into VinVL’s embedding layer to reduce its size to 1×768 . We then take the average of all image features $\bar{\mathbf{f}} = \frac{1}{N} \sum_{i=1}^N \mathbf{f}_i$ to construct a context feature (1×768) which will be used later. Simultaneously, we obtain L word embeddings of the caption $\{\mathbf{w}_i\}_{i=1}^L$, each of which has a size of 1×768 . For more information on image feature extraction and tokenizer, see VinVL [38].

We now detail knowledge graph construction. We follow Chen et al. [8] to detect concepts for each input image. Specifically, we use clarifai [1] to obtain the top-ten concepts $\{c_i\}_{i=1}^{10}$ for each image. As a result, we detect $k \times 10$ concepts in total. Then, using ConceptNet [30], we use each

detected concept as a query to heuristically retrieve forecasted concepts with 2-hop neighbors of the query. Since the number of forecasted concepts is large (> 400) and many of them are unrelated to input images, we employ a filtering process to retain only the informative concepts.

Let c_i^f be a forecasted concept. Using a pre-trained language model RoBERTa [20], we compute a relevance score between the forecasted concept and image context as:

$$\rho_{c_i^f} = f_{\text{head}}(f_{\text{enc}}([\bar{\mathbf{f}}; \mathbf{c}_i^f])),$$

where $\mathbf{c}_i^f = \text{BERT}(c_i^f)$ is an embedding vector of the concept c_i^f extracted by a pre-trained BERT [10], $[\cdot; \cdot]$ denotes the concatenation operator, f_{enc} is the encoder part of the language model while f_{head} is a softmax layer. This score indicates the probability of c_i^f related to $\bar{\mathbf{f}}$.

We keep M forecasted concepts having high relevance scores. In total, we have $k \times 10$ detected concepts $\{c_i\}_{i=1}^{k \times 10}$ and M forecasted concepts $\{c_i^f\}_{i=1}^M$ in our knowledge graph ($k \times 10 + M$ nodes). If two concepts are related in the ConceptNet [30], an undirected edge is given to connect them. For simplicity, we do not use a specific relation (e.g., has, IsA). Furthermore, a concept in I_i is connected to its related concepts in the adjacent images I_{i-1} and I_{i+1} to ensure information flow and the awareness of the temporal order of the images. Hereafter, we use the same notation to refer to detected and forecasted concepts $\{c_i\}_{i=1}^{k \times 10 + M}$.

Graph neural network is used to update the node embeddings through iterative messages passing between neighbors on the graph. We use graph attention network [32] to build our graph neural network. To produce the input for the graph network, we first employ pre-trained BERT [10] to embed each concept into an embedding with the size of 1×768 . To be more specific, each node embedding is calculated as $\mathbf{e}_i = \text{BERT}(c_i)$. To strengthen the connection between concepts and image context, we concatenate the node embedding and the context feature as $\mathbf{e}_i = [\mathbf{e}_i; \bar{\mathbf{f}}]$. Brevity, we summarize the entire computation in each graph layer:

$$\{\tilde{\mathbf{e}}_1^{(l)}, \dots, \tilde{\mathbf{e}}_{k \times 10 + M}^{(l)}\} = \text{GNN}(\{\mathbf{e}_1^{(l-1)}, \dots, \mathbf{e}_{k \times 10 + M}^{(l-1)}\}),$$

where l indicates the current graph layer while $l-1$ does the previous one, $\text{GNN}(\cdot)$ represents a graph layer. In detail, each node is updated by:

$$\begin{aligned} \hat{\alpha}_{ji} &= (\mathbf{e}_i^{(l-1)} \mathbf{W}_q)(\mathbf{e}_j^{(l-1)} \mathbf{W}_k)^\top, \\ \alpha_{ji} &= \text{SOFTMAX}(\hat{\alpha}_{ji} / \sqrt{D}), \\ \hat{\mathbf{e}}_i^{(l-1)} &= \sum_{j \in \mathcal{N}_i \cup \{i\}} \alpha_{ji} (\mathbf{e}_j^{(l-1)} \mathbf{W}_v), \\ \tilde{\mathbf{e}}_i^{(l)} &= \text{LAYERNORM}(\mathbf{e}_i^{(l-1)} + \hat{\mathbf{e}}_i^{(l-1)} \mathbf{W}_o), \end{aligned}$$

where $\mathbf{W}_q, \mathbf{W}_k, \mathbf{W}_v, \mathbf{W}_o \in \mathbb{R}^{D \times D}$ are learnable matrices, \mathcal{N}_i represents the neighbors of node i , $D = 768 + 768 =$

1536, SOFTMAX and LAYERNORM are the softmax function and the batch normalization, respectively. We note that $\mathbf{e}_i^{(0)}$ is the initial node embedding (i.e., $[\mathbf{e}_i; \bar{\mathbf{f}}]$).

In practice, we use 2 graph layers. After the graph attention network, we add two more fully connected layers to reduce the size of each $\tilde{\mathbf{e}}_i$ to 1×768 .

Frozen VinVL. As discussed above, the concept embeddings learned from the graph neural network are used as a prompt to generate the caption. To this end, we inject all $\{\tilde{\mathbf{e}}_i\}_{i=1}^{k \times 10 + M}$ into a frozen pre-trained VinVL [38]. As a result, the input of VinVL is changed to $\{\mathbf{w}_1, \dots, \mathbf{w}_L, [\text{SEP}], \tilde{\mathbf{e}}_1, \dots, \tilde{\mathbf{e}}_{k \times 10 + M}, [\text{SEP}], \mathbf{f}_1, \dots, \mathbf{f}_N\}$. We note that [SEP] is a special token used to distinguish different types of tokens. We do not feed \mathbf{w}_i to the network during inference time, but instead, create $L \times [\text{MASK}]$ as pseudo words. Formally, Eq. 1 becomes

$$C_{k+1} = \text{A-CAP}(\tilde{\mathbf{e}}_1, \dots, \tilde{\mathbf{e}}_{k \times 10 + M}, [\text{SEP}], \mathbf{f}_1, \dots, \mathbf{f}_N).$$

Loss function. Following previous works, we simply use cross entropy between the generated and the ground-truth captions to train the network. We do not use CIDEr optimization because the pre-trained VinVL has been well-trained on a large text-image corpus.

4. Experiments

4.1. Dataset and training details

Dataset. We use the visual storytelling dataset (VIST) [14] with a modification to evaluate our method because there is no dataset tailored for our task. The original VIST includes 210,819 photos from 10,117 Flickr albums. Given five input temporally ordered images from the same event, the corresponding five human-annotated sentences are provided as ground-truths. There are 4,098, 4,988, and 5,050 samples for training, validation, and test sets, respectively. We use the first four images of each sample as input ($k = 4$) and the last sentence of each sample as the ground-truth caption. We keep the last image of each sample as an oracle image for reference. The training, validation, and test sets all have the same number of samples as the original dataset.

Dataset verification. We investigate the correlation between C_{k+1} and C_1, C_2, \dots, C_k (corresponding captions to I_1, I_2, \dots, I_k) in two ways. First, we compute the sentence cosine similarity $\text{sim}(S(C_{k+1}), S(C_i))$ ($i = 1, \dots, k$) and then test whether those similarities monotonically increase (i.e., $\text{sim}(S(C_{k+1}), S(C_i)) < \text{sim}(S(C_{k+1}), S(C_{i+1}))$) ($S(\cdot)$ is a pre-trained SentenceTransformer model [2], outputting an embedding vector for a given sentence). We confirm that 72.69% of samples follow monotonic increasing, 10.32% have only one sentence similarity that violates monotonic increasing, and only 4.4% do not comply with the monotonicity. As the second, we use a pre-trained BERT model [10] to figure out whether C_{i+1} is the next







Sparsely temporally-ordered images	Oracle image (for reference)	VinVL	VinVL + Oracle image	AREL + BART	A-CAP	Ground-truth
		after the ceremony, the teams got to eat outside.	the defense was able to close out the game and had a great time.	i was getting ready to leave the game and i took a picture of the players on the field before the game.	the goalie caught the puck as it passed the goalie.	he shoots, he scores and the game ends one to nothing.
		she let the crowd ask questions in the end.	we got to meet the people behind the company's logo.	he welcomed to the stage his new assistant	at the end of the show, the audience enjoyed themselves.	they were all about preserving the internet
		the llamas were very curious.	the competition ended with a bang.	they had a great time.	it was a great time for the horse racers.	he thought he was going to cry
		the vice president closed the meeting by thanking all the workers of the company.	the party went on well into the night.	everyone was having a great time.	they ended the night with a speech.	eventually the winner was announced, and he was very grateful

Figure 3. Examples of generated captions obtained by all compared methods. We show the oracle images and ground-truth captions for reference purposes. VinVL [38] generates captions that are out of context with the input images. VinVL [38] + Oracle image sometimes generates reasonable captions. AREL [35] + BART [17] tends to generate a general ending for the sequence of images. In contrast, our method A-CAP predicts more accurate, descriptive, and plausible captions than others.

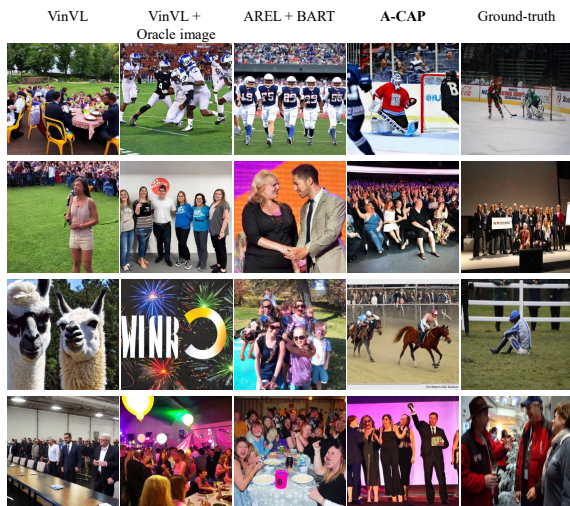


Figure 4. The generated images obtained by using stable diffusion model [28] to generate an image from each generated caption in Fig. 3. The order of images is the same as the order of captions in Fig. 3. The images generated using our captions are close to the ground-truth ones while those by other methods are not.

sentence of C_i . We see that 77.34% of the samples satisfy the next sentence condition (i.e., C_{i+1} is always the next

sentence of C_i for all sentences in the sequence), 17.78% have only one sentence that does not meet the condition, and 0.06% do not satisfy the condition (i.e., C_{i+1} is never the next sentence of C_i). The above verification shows that the VIST dataset mostly meets our assumption.

Training details. We set the length of the word sequence $L = 35$, the number of ROIs $N = 100$ (25 ROIs for each image), the number of forecasted concepts $M = 60$ (the number of concepts is $4 \times 10 + 60 = 100$ in total).

We build A-CAP using PyTorch, in which we use the pre-trained VinVL model published by its authors [3]. We remark that we freeze all the parameters of VinVL during training time. Given the small size of our used dataset, we train the model for only 10 epochs with a batch size of 16 and a learning rate of $3e-5$. It takes four hours to train our model on a single GTX-3090 GPU.

4.2. Compared methods and evaluation metrics

Compared methods. We carefully design methods that can be straightforwardly applied to our task. For a fair comparison, all compared methods are fine-tuned on VIST. To avoid over-tuning, we only train the methods for a few epochs and select their best checkpoints.

VinVL [38] is a cutting-edge image captioning model. We strictly adhere to its settings, but instead of a single im-

age, we use the input as our method. Comparing our method to VinVL will demonstrate the advancement of our method over the conventional image captioning model.

VinVL [38] + Oracle image is the method where VinVL uses the ground-truth oracle image in training and testing. Since we do not successfully generate oracle images using existing methods, we may regard this method as a method that sequentially generates the oracle image and caption.

AREL [35] + BART [17] is a combination of visual storytelling (AREL [35]) and story ending generation (BART [17]). Particularly, we generate a story for the input and then generate the ending sentence for that story. We compare the ending sentence to the caption by our method.

Evaluation metrics. Since our problem is an open domain generation like dialogue generation, we follow [11] to use automatic metrics to quantitatively evaluate all the methods in two aspects: *accuracy* and *descriptiveness*. For accuracy evaluation, we report referenced metrics including BLEU [25], CIDEr [31]. Since those metrics are sensitive to the whole sentence structure [19], we also report SPICE [5], CLIPScore, and RefCLIPScore [13] to overcome the structural dependency. For descriptiveness evaluation, we adopt a self-retrieval strategy, drawing on prior work. This strategy is based on the observation that more descriptive captions with significant details frequently lead to more precise self-retrieval, i.e., retrieving the target image from a set of similar images given the generated caption. We report the refined R@1, R@5, and R@10 scores using CLIP [26] as the retriever.

4.3. Qualitative comparisons

In Fig. 3, we show some randomly selected examples of captions generated by our method as well as others. Despite its enormous success in image captioning, VinVL [38] is unable to generate the expected captions. We can see that the captions generated by VinVL are completely out of context with the input images. This observation suggests that the current image captioning model is inadequate for our task. VinVL [38] + Oracle image generates reasonable captions to some extent when the oracle images are close enough to the input images (see first and second samples). However, if the temporal information is too sparse as in the third and fourth samples, it fails to generate captions that are linked to the inputs. These results imply that even if we can generate a high-quality unseen oracle image, the model struggles to complete the task. We notice that AREL [35] + BART [17] generates a general ending for the story (e.g., having a great time). On the contrary, our method produces more accurate and reasonable captions that reflect the inputs’ future. In most cases, we can see that our method accurately predicts what is likely to happen, which is close to the ground-truth captions. When we examine the third sample in greater detail, we can see that our caption is in-

correct because we failed to detect the concept “falling” in the second image. However, we believe that the generated caption is still plausible under ordinary situations.

To have a better understanding of the generated captions, we use the stable diffusion model [28] implemented on the Huggingface platform [4] with the default settings to generate an image from each generated caption, and choose the first generated image for each method as shown in Fig. 4. The images obtained from our generated captions are similar to the ground-truth ones, indicating that our method generates correct anticipated captions. Furthermore, Fig. 4 demonstrates the benefits of our task to downstream tasks, specifically future image generation in this case.

4.4. Quantitative comparisons

The quantitative scores are summarized in Table 1, first four rows. We first assess all methods based on their accuracy. All of the results in Table 1 support the advantage of our method over the other methods. Though our method obtains the highest scores, we notice that it does not significantly outperform the other methods on referenced metrics (BLEU and CIDEr). The reason for this observation is that those metrics are calculated using ground-truth captions. Because our task is an open-domain generation, it is difficult to generate a caption that is nearly identical to the ground-truth one. However, based on the qualitative comparison in Figs. 3 and 4, we can conclude that our method outperforms the others. SPICE and the unreferenced metrics (CLIPScore, RefCLIPScore) also justify our conclusion. We see substantial improvements in these metrics, indicating that our generated captions accurately reflect the oracle images. Notably, as shown in Fig. 3, our generated captions are, without a doubt, the future of input images.

The descriptiveness of generated captions is then assessed using R@1, R@5, and R@10 scores. In comparison to VinVL [38] and AREL [35] + BART [17], our method outperforms them significantly. This is thanks to the fact that captions generated by our method are close to the ground-truth images, whereas those obtained by the other methods are not. Our method and VinVL [38] + Oracle image achieve the same level. This is not surprising, given that VinVL [38] + Oracle image generates captions directly from oracle images.

We conclude that our method is more promising than the other methods in solving the anticipation captioning task. Furthermore, the experiments highlight the shortcomings of using image captioning and story ending models in our task.

4.5. Detailed analysis

Ablation study. To validate the plausibility of our model design, we investigate two ablated models: A-CAP w/o GNN and A-CAP w/o context. A-CAP w/o GNN denotes the model that does not use a graph neural network (instead,

Table 1. Quantitative comparison against other methods. For accuracy evaluation, we report referenced metrics (BLEU [25] (B-1, B-4), CIDEr [31]), SPICE [5], and unreferenced metrics (CLIPScore and RefCLIPScore [13]). For descriptiveness evaluation, we report top-1, top-5 and top-10 retrieval accuracy (R@1, R@5, R@10, respectively). Our method outperforms others on all metrics. Higher scores are better. Gray background indicates results obtained by our method, and Δ indicates the improvement over compared methods.

Method	Accuracy						Descriptiveness		
	B-1	B-4	CIDEr	SPICE	CLIPScore	RefCLIPScore	R@1	R@5	R@10
VinVL [38]	31.7	3.1	2.6	13.8	40.7	42.8	1.3	6.5	10.8
VinVL [38] + Oracle image	34.9	3.8	4.3	16.9	57.9	61.3	8.1	17.2	31.1
AREL [35] + BART [17]	30.9	2.0	3.1	11.4	37.8	39.7	1.1	5.9	9.3
A-CAP	37.2	6.9	4.7	20.1	65.2	70.2	8.7	18.9	31.5
A-CAP w/o GNN	34.8	5.2	3.7	14.5	38.2	47.3	3.6	8.7	15.4
A-CAP w/o context	36.1	6.2	4.2	13.9	39.8	46.9	4.1	9.5	16.1
Δ	2.3 \uparrow	3.1 \uparrow	0.4 \uparrow	3.2 \uparrow	7.3 \uparrow	8.9 \uparrow	0.6 \uparrow	1.7 \uparrow	0.4 \uparrow

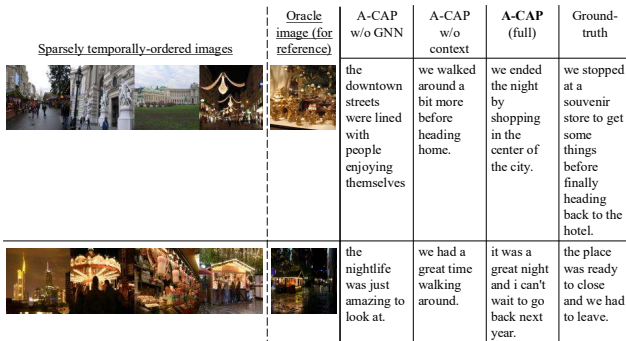


Figure 5. Examples of generated captions by two ablated models: A-CAP w/o GNN, A-CAP w/o context, and full model A-CAP. We select two inputs where the detected concepts almost overlap. A-CAP w/o GNN generates captions that most likely describe the inputs. A-CAP w/o context generates captions that are far from the inputs and similar to each other.

we directly feed the concept embeddings $\tilde{e}_i = \text{BERT}(c_i)$ to the pre-trained VinVL). A-CAP w/o context is the model in which we do not concatenate the node embeddings and the context feature (we instead use only the node embeddings as graph neural network inputs). We also drop the two fully connected layers on top of the graph neural network because reducing the size of embeddings is no longer required.

The last two rows of Table 1 quantify the performance of the two ablated models. When we simplify the model, the performance scores are degraded. In the case of A-CAP w/o GNN, the concept embeddings are insufficient to guide the model to generate the expected caption. As a result, the caption most likely describes the inputs as depicted in Fig. 5. The graph neural network enriches and connects concept embeddings, making them more powerful as a prompt to the model. Similarly, A-CAP w/o context breaks the connections between concepts and the context of images in gen-

Table 2. Impact of the number of forecasted concepts on the performance of our model. Using either a large number of concepts or no concepts drops the performance drastically.

Number of forecasted concepts	Accuracy			Descriptiveness		
	SPICE	CLIPScore	RefCLIPScore	R@1	R@5	R@10
$M = 400$	5.8	15.3	12.1	1.1	3.7	7.6
$M = 200$	5.4	16.7	13.0	0.9	4.2	7.1
$M = 100$	15.7	48.6	52.4	6.2	15.7	26.6
$M = 60$ (used model)	20.1	65.2	70.2	8.7	18.9	31.5
$M = 0$	14.2	43.1	44.7	1.9	7.3	11.2

eral, resulting in captions that are far from the inputs and similar to each other if the detected concepts are similar (Fig. 5). This indicates that the context feature compensates for the concepts in order to make the correct prediction. In contrast, the full model generates plausible captions.

We do not investigate the model where all the parameters are trainable since the training collapsed despite our best efforts. The reason for this failure is that the training data is too small in comparison with the one used to train VinVL.

Impact of the number of forecasted concepts. As stated above, when we search for concepts on ConceptNet, we usually have more than 400 forecasted concepts. We empirically retain $M = 60$ forecasted concepts to eliminate irrelevant concepts and balance the number of concepts and image features. We now investigate how the number of forecasted concepts affects the captions generated. To this end, we run our method through a series of scenarios using the number of forecasted concepts at 400, 200, 100, and 0.

Table 2 shows the results of all tested scenarios on accuracy and descriptiveness. We can see that retrieving a large number of concepts ($M = 400$ or $M = 200$) degrades performance. The reason is obvious because when we include a larger number of irrelevant concepts, the input becomes too noisy, preventing the model from selecting essential information. The model with $M = 100$ forecasted concepts



Sparsely temporally-ordered images	Oracle image (for reference)	A-CAP	Ground-truth
		the bride and groom are about to cut the cake.	night settles on this wonderful day and everyone heads home.

Figure 6. A case study of samples with low scores. Though our method generates a plausible caption, it is far from the ground-truth caption. The reason is that the oracle image changes significantly from the inputs.

comes close to our best performance ($M = 60$). Finally, we examine an extreme case where no forecasted concept is employed ($M = 0$). The performance drops to the same level as that of VinVL [38] (first row in Table 1). This is due to the fact that the inputs to the two models are nearly identical. This experiment confirms that the number of forecasted concepts has an effect on our performance, implying that retrieving a sufficient number of concepts results in improved effectiveness.

A case study of samples with low scores. While our method produces promising quantitative results, we notice a relatively small number of samples with low scores when delving into each sample in detail. We thus manually check those samples, as shown in Fig. 6. Given what is happening in the inputs, our generated caption is reasonable because the next step of the wedding party is “cutting a wedding cake”. The ground-truth caption, in contrast, is completely different because the scene shifts from “wedding” to “night-time”. We recall that our hypothesis is that the scene does not change significantly, but in this case, it does. Though our method fails to predict the far future, it does correctly predict the near future. We may ignore such failures because they contradict our hypothesis. In fact, when we exclude those failure samples from quantitative comparison, our outperformance becomes more significant than before.

Limitations. First, our method is heavily reliant on concept detection (here, clarifai). When we are unable to detect important concepts, our method is unable to predict the correct caption, as seen in Fig. 4, third example. Second, as shown in Table 2, the performance of our method is dependent on the number of forecasted concepts from commonsense knowledge. We use a simple filtering process in this paper, namely, computing the relevance score between concept and image context and empirically retaining $M = 60$ forecasted concepts. Our strategy is effective, but it may not be optimal. To improve this issue, it is necessary to learn how to determine a suitable number of concepts. One possible solution is to learn concept selection while training the model. This is left for our future work.

5. Discussions

We now discuss the potential negative societal impacts of our task. While we believe our introduced task will push

more applications to make our lives safer and benefit downstream tasks, we have noticed that it has the potential to be abused. One of the concerns is that it will be used to predict behavior for nefarious purposes, such as criminal activity.

Our task still has some difficulties. First, to the best of our knowledge, no suitable dataset exists to serve as a benchmark. Though our used VIST dataset [14] is useful to some extent, it is originally designed for the visual storytelling task, so it does not always meet task requirements, as already seen. As a result, a new dataset for this task is required, which should cover various scenarios such as near future, far future, abnormal thinking, and rationale. We should note that owing to the labor cost of creating a dataset, we are currently using the customized VIST to assess the performance of our method. Second, evaluating the task is difficult. Although appropriate evaluation metrics for the open domain are still unavailable, our used metrics are partially effective in our task. This is because, as we do not account for the diversity of potential futures, generating a caption close to the ground-truth (BLEU, CIDEr) is a valid indicator of the model’s predictive capability. Moreover, considering the dataset that we employed, CLIP-based scores are suitable for evaluating the degree of similarity between the generated captions and the oracle images, which are presumed to represent the future of the input images. In fact, our experiments show that the current metrics cannot evaluate the task thoroughly. User study may compensate for the automatic metrics, but it is expensive and subjective, as is customary. We believe that new metrics for this task can capitalize on the advantages of the vision-language space, such as CLIP [26]. Furthermore, new metrics should emphasize the rationale, which explains the reason why the model generates that caption but not another.

6. Conclusion

We introduced a new task, called anticipation captioning, that generates a caption for an unseen oracle image, given a sparsely temporally-ordered set of images. For this new task, we proposed a baseline model (A-CAP), which incorporates commonsense knowledge into the off-the-shelf vision-language model VinVL. We evaluated A-CAP on a customized VIST dataset, showing that A-CAP outperforms other image captioning methods. We also addressed the potential positive and negative impacts of the task as well as its challenges, in order to encourage further research.

Acknowledgement. This work was supported by the Institute of AI and Beyond of the University of Tokyo, JSPS/MEXT KAKENHI Grant Numbers JP19H04166, JP22H05015, and 22K17947, and the commissioned research (No. 225) by the National Institute of Information and Communications Technology (NICT), Japan.

References

- [1] <https://github.com/soodoku/clarifai>. 3
- [2] <https://www.sbert.net/>. 4
- [3] <https://github.com/microsoft/Oscar>. 5
- [4] <https://huggingface.co/spaces/stabilityai/stable-diffusion>. 6
- [5] Peter Anderson, Basura Fernando, Mark Johnson, and Stephen Gould. Spice: Semantic propositional image caption evaluation. In *ECCV*, 2016. 6, 7
- [6] 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 *CVPR*, 2018. 1, 2
- [7] Gang Chen, Yang Liu, Huanbo Luan, Meng Zhang, Qun Liu, and Maosong Sun. Learning to generate explainable plots for neural story generation. *IEEE/ACM Trans. Audio, Speech and Lang. Proc.*, 2020. 2
- [8] Hong Chen, Yifei Huang, Hiroya Takamura, and Hideki Nakayama. Commonsense knowledge aware concept selection for diverse and informative visual storytelling. In *AAAI*, 2021. 2, 3
- [9] Long Chen, Zhihong Jiang, Jun Xiao, and Wei Liu. Human-like controllable image captioning with verb-specific semantic roles. In *CVPR*, 2021. 2
- [10] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In *NAACL*, 2019. 4
- [11] Sarah E. Finch and Jinho D. Choi. Towards unified dialogue system evaluation: A comprehensive analysis of current evaluation protocols. In *Proceedings of the 21th Annual Meeting of the Special Interest Group on Discourse and Dialogue*, pages 236–245, 1st virtual meeting, July 2020. Association for Computational Linguistics. 6
- [12] Danijar Hafner, Timothy Lillicrap, Jimmy Ba, and Mohammad Norouzi. Dream to control: Learning behaviors by latent imagination. *ICLR*, 2019. 1, 2
- [13] Jack Hessel, Ari Holtzman, Maxwell Forbes, Ronan Le Bras, and Yejin Choi. CLIPScore: a reference-free evaluation metric for image captioning. In *EMNLP*, 2021. 6, 7
- [14] Ting-Hao K. Huang, Francis Ferraro, Nasrin Mostafazadeh, Ishan Misra, Jacob Devlin, Aishwarya Agrawal, Ross Girshick, Xiaodong He, Pushmeet Kohli, Dhruv Batra, et al. Visual storytelling. In *NAACL*, 2016. 1, 2, 4, 8
- [15] Justin Johnson, Andrej Karpathy, and Li Fei-Fei. Densecap: Fully convolutional localization networks for dense captioning. In *CVPR*, 2016. 2
- [16] Vincent Le Guen and Nicolas Thome. Disentangling physical dynamics from unknown factors for unsupervised video prediction. In *CVPR*, 2020. 2
- [17] Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. In *ACL*, 2020. 5, 6, 7
- [18] Zhongyang Li, Xiao Ding, and Ting Liu. Story ending prediction by transferable bert. In *IJCAI*, 2019. 2
- [19] Lixin Liu, Jiajun Tang, Xiaojun Wan, and Zongming Guo. Generating diverse and descriptive image captions using visual paraphrases. In *ICCV*, 2019. 6
- [20] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized BERT pretraining approach. *CoRR*, 2019. 4
- [21] Jiasen Lu, Jianwei Yang, Dhruv Batra, and Devi Parikh. Neural baby talk. In *CVPR*, 2018. 1, 2
- [22] Lukáš Neumann, Andrew Zisserman, and Andrea Vedaldi. Future event prediction: If and when. In *CVPRW*, 2019. 1, 2
- [23] Qiang Ning, Zhili Feng, Hao Wu, and Dan Roth. Joint reasoning for temporal and causal relations. In *ACL*, 2018. 2
- [24] Qiang Ning, Hao Wu, Rujun Han, Nanyun Peng, Matt Gardner, and Dan Roth. TORQUE: A reading comprehension dataset of temporal ordering questions. In *EMNLP*, 2020. 2
- [25] Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic evaluation of machine translation. In *ACL*, 2002. 6, 7
- [26] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision. In *ICML*, 2021. 6, 8
- [27] Shaoqing Ren, Kaiming He, Ross B. Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In *NIPS*, 2015. 3
- [28] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *CVPR*, 2022. 1, 5, 6
- [29] Julian Schrittwieser, Ioannis Antonoglou, Thomas Hubert, Karen Simonyan, Laurent Sifre, Simon Schmitt, Arthur Guez, Edward Lockhart, Demis Hassabis, Thore Graepel, et al. Mastering atari, go, chess and shogi by planning with a learned model. *Nature*, 2020. 1, 2
- [30] Robyn Speer, Joshua Chin, and Catherine Havasi. Conceptnet 5.5: An open multilingual graph of general knowledge. In *AAAI*, 2017. 2, 3, 4
- [31] Ramakrishna Vedantam, C. Lawrence Zitnick, and Devi Parikh. Cider: Consensus-based image description evaluation. In *CVPR*, 2015. 6, 7
- [32] Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua Bengio. Graph attention networks. *ICLR*, 2017. 4
- [33] Duc Minh Vo, Hong Chen, Akihiro Sugimoto, and Hideki Nakayama. NOC-REK: novel object captioning with retrieved vocabulary from external knowledge. In *CVPR*, 2022. 1, 2
- [34] Carl Vondrick and Antonio Torralba. Generating the future with adversarial transformers. In *CVPR*, 2017. 2
- [35] Xin Wang, Wenhui Chen, Yuan-Fang Wang, and William Yang Wang. No metrics are perfect: Adversarial reward learning for visual storytelling. In *ACL*, 2018. 1, 2, 5, 6, 7
- [36] Dongchao Yang, Jianwei Yu, Helin Wang, Wen Wang, Chao Weng, Yuexian Zou, and Dong Yu. Diffsound: Discrete dif-

- fusion model for text-to-sound generation. *arXiv e-prints*, 2022. [1](#)
- [37] Kuo-Hao Zen, William B. Shen, De-An Huang, Min Sun, and Juan Carlos Niebles. Visual forecasting by imitating dynamics in natural sequences. In *ICCV*, 2017. [1](#), [2](#)
- [38] Pengchuan Zhang, Xiujun Li, Xiaowei Hu, Jianwei Yang, Lei Zhang, Lijuan Wang, Yejin Choi, and Jianfeng Gao. Vinvl: Making visual representations matter in vision-language models. In *CVPR*, 2021. [1](#), [2](#), [3](#), [4](#), [5](#), [6](#), [7](#), [8](#)
- [39] Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. Conditional prompt learning for vision-language models. In *CVPR*, 2022. [2](#)
- [40] Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. Learning to prompt for vision-language models. *IJCV*, 2022. [2](#)