# Uncertainty-based Visual Question Answering: Estimating Semantic Inconsistency between Image and Knowledge Base

Jinyeong Chae*
*OKESTRO Co., Ltd*
Seoul, Republic of Korea
jiny491@gmail.com

Jihie Kim
*Dept. of Artificial Intelligence*
*Dongguk University*
Seoul, Republic of Korea
jihie.kim@dgu.edu

*Abstract*—Knowledge-based visual question answering (KVQA) task aims to answer questions that require additional external knowledge as well as an understanding of images and questions. Recent studies on KVQA inject an external knowledge in a multi-modal form, and as more knowledge is used, irrelevant information may be added and can confuse the question answering. In order to properly use the knowledge, this study proposes the following: 1) we introduce a novel semantic inconsistency measure computed from caption uncertainty and semantic similarity; 2) we suggest a new external knowledge assimilation method based on the semantic inconsistency measure and apply it to integrate explicit knowledge and implicit knowledge for KVQA; 3) the proposed method is evaluated with the OK-VQA dataset and achieves the state-of-the-art performance.

*Index Terms*—knowledge-based visual question answering, semantic inconsistency, uncertainty, knowledge graph

## I. INTRODUCTION

Knowledge-based visual question answering (KVQA) task is to answer questions that require an understanding of images, questions, and additional external knowledge. The KVQA task is proposed with the aim of reaching human-level reasoning. Injecting huge knowledge related to the entities identified from images and questions in a multi-modal form is among the tasks being researched. However, as the knowledge base (KB) is often incomplete, when the context of the entities is not fully consistent with the KB, irrelevant information can be retrieved and confuse the question answering.

For the example in Fig. 1, the question can be answered with a full understanding of the image and question. However, the predicted answer can become yellow when we use related general knowledge, i.e., *(banana, HasProperty, yellow)*. In this case, there is a conflict between the image and the knowledge base. We define semantic inconsistency as such conflicts between the image context and the knowledge extracted based on the object in the image or question keywords. For KVQA, there have been a lot of approaches introduced to make use of external KB with the given image and the question. Recent studies [1] and [2] suggested a method of extracting external knowledge by using the object keywords

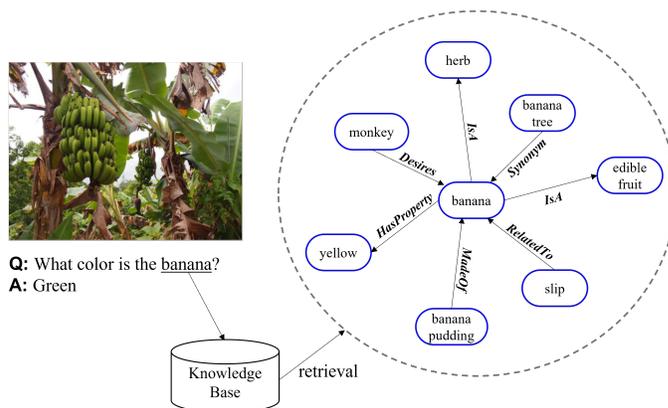*This work was done when J. Chae was with Dongguk University.



Fig. 1: An example of visual question answering, occurring a semantic inconsistency between the image and external knowledge. The knowledge graph is the external knowledge extracted according to the word of the question.

of the image and the words in the question. However, as shown above, these approaches can suffer from semantic inconsistency when the context information of the VQA is not well used given the new information from the KB. [3], [4], and [5] introduced graph-based approaches for the KVQA where the work focuses on how to extract the needed knowledge using graph algorithms. However, they lack considering how well the extracted knowledge match with the given image and the question, and the extracted knowledge can rather confuse the answer prediction. In such cases, we believe it is necessary to evaluate and adjust the amount of the external knowledge injected based on semantic inconsistency between the context of the image and the knowledge. In estimating semantic inconsistency, first we propose to make use of caption generation results which can indicate unusuality of the image. In addition, we develop a new uncertainty-based measure that uses knowledge context model pre-trained with commonsense knowledge. When generating the knowledge context for an image, the context that is inconsistent with the actual image context may incur uncertainty, and the semantic inconsistency

can be estimated through an uncertainty measure. Furthermore, the inconsistency can be also estimated by the similarity between the generated knowledge context and the image context. This study combines these into a new approach for measuring such inconsistencies and introduces a new way of assimilating external knowledge. This study is summarized as follows:

- We introduce a new semantic inconsistency measure based on caption generation, which is an ensemble of a) uncertainty of the caption and b) similarity between the caption generated with the KB and the ground-truth caption
- We propose an external knowledge assimilation method based on the proposed semantic inconsistency measure to control the use of external knowledge in KVQA.
- We apply the proposed method for combining explicit and implicit knowledge passed through Relational Graph Convolution Networks (RGCN) and VisualBERT, respectively in KVQA and achieve the state-of-the-art result when evaluated with the OK-VQA dataset.

## II. RELATED WORK

### A. KVQA approaches using pre-trained model

A lot of researches have studied image and text as a multi-modal form. By tokenizing the object in an image, an alignment between an object and text has been proposed to apply a self-attention model [6] [7]. In addition, [7] showed such models achieve better performance in various downstream tasks compared with other vision-language approaches [8]. Therefore, this study experiments with the approach suggested by [7] for extracting the implicit knowledge. Multi-modal approaches using image features from Faster R-CNN or ResNet and question embedding of pre-trained models are also proposed [9] [10]. [9] generated joint representation through Bilinear Attention Map. [10] extracted image-text joint representation by using image features and question embedding, and proposed a 3-way Tucker fusion method. In addition to using pre-trained models, there have also been studies trying to solve VQA tasks using additional external knowledge. [1] proposed ArticleNet using Wikipedia search API related to keywords of an image and words of the question. A method for extracting an external knowledge related to the objects in an image was also introduced [2]. [2] extracted the knowledge by using the object label output from the Faster R-CNN model. This study extracts more relevant knowledge by using not only image object keywords, but also words in the question.

### B. Graph-based KVQA approaches

Besides using pre-trained models, studies using graph-based models were proposed [11] [12] [13] [14] [3]. [11] suggested the Neural State Machine based on a probabilistic graph for reasoning on VQA. [13] introduced a video scene graph and caption generation method, and applied them for reasoning on video-QA task. [12] studied a heterogeneous graph alignment network considering inter-and intra-modality for video-QA. [14] proposed a method to create graphs from visual,

linguistic, and numeric features and suggested an aggregator that combines the features. However, because the study focuses on the contents of the image, the method has a limitation in answering a question that requires additional knowledge. [3] suggested graph-based VQA for capturing the interrelationship between objects and entities of external knowledge by combining concept graph and scene graph. However, the scene graph relation is limited because only locational information is considered, and in the OK-VQA dataset, the extraction method for location-based scene graphs does not show significant performance improvements.

Moreover, studies using a pre-trained model and graph-based model have been suggested [15] [4] [5]. [15] introduced multi-modal graph networks for compositional generalization in VQA, but the method is evaluated with the VQA task that only requires object detection or recognition in answering questions for the object shape and the number of objects. [4] proposed a Knowledge Graph Augmented model using a pre-trained object detection model and graph-based method. However, the knowledge subgraph is generated by using the image object labels and the words of the question without considering the image-question context. [5] proposed to integrate image-text representation from the BERT-based model and graph information based on the concept of image objects and questions. However, when there are conflicts between the graph and the pre-trained model representation, use of knowledge can hinder the question answering, as described above. This study proposes a new approach that measures semantic inconsistencies between KB and the given problem, and moderates the use of knowledge based on the measurement.

## III. APPROACH

This section introduces a semantic inconsistency measure that makes use of uncertainty and semantic similarity modeling.

### A. Semantic inconsistency between an image and an external KB

In this study, we utilize caption generation to measure semantic inconsistency between an image and external KB. Inspired by [16], we adopt uncertainty model of caption generation and introduce a novel measure for estimating semantic inconsistency between the KB and the VQA context.

*1) Ensemble-based uncertainty estimation for KVQA:* In the existing image captioning, to generate a sentence $y$ when an input $x$ is given, the conditional distribution $p(y|x)$ is learned and tokens are continuously predicted from an autoregressive distribution.

$$p(y|x) = p(y_1|x) \prod_{i=2}^{k} p(y_i|x, y_1, \cdots, y_{i-1}) \qquad (1)$$

In Eq. (1), $y_i$ denotes the token corresponding to the index $i$ in sentence $y$, and the given set $\{x, y_1, \cdots, y_{i-1}\}$ denotes context $c_i$ for predicting the token corresponding to $i$. The

number of tokens that can be predicted is limited based on the given context. For example, the word "beach" cannot be generated when an image of a cat on a desk is given. When a set of words irrelevant to the context is denoted hallucinated word $V_h^{(c_i)}$, the following equation can be written

$$p(y_i \in V_h^{(c_i)}) = \sum_{v \in V_h^{(c_i)}} p(y_i = v|c_i) \tag{2}$$

In image captioning, token prediction in a given context is calculated with the following cross-entropy equation. The equation can be divided into two based on an entropy of the set of words relevant to the context and that of the set of words irrelevant to the context as

$$\begin{aligned} H(y_i|c_i) &= -\sum_{v \in V} p(y_i = v|c_i) log p(y_i = v|c_i) \\ &= -\sum_{v \in V \setminus V_h^{(c_i)}} p(y_i = v|c_i) log p(y_i = v|c_i) \\ &\quad -\sum_{v \in V_h^{(c_i)}} p(y_i = v|c_i) log p(y_i = v|c_i) \end{aligned} \tag{3}$$

The uncertainty that can be predicted by the Eq. (3) can be divided into two: 1) uncertainty that appears in selection of a token that describes the context; 2) uncertainty that appears due to the interference of words irrelevant to the context or an insufficient training system. The latter is directly related to calculating hallucinated words that are irrelevant to the given context. We make use of the latter in measuring uncertainty in KVQA, as described below. The latter can be decomposed into two: aleatoric uncertainty and epistemic uncertainty [17] [18] [19]. The uncertainties can be measured by an ensemble-based model [20] and calculated as follows:

$$\begin{aligned} u_{al}(y_i|c_i) &= \mathbb{E}_{q(w)}[H(y_i|c_i, w)] \\ &= \frac{1}{M} \sum_{m=1}^{M} H_m(y_i|c_i) \end{aligned} \tag{4}$$

$$\begin{aligned} u_{ep}(y_i|c_i) &= H(y_i|c_i) - \mathbb{E}_{q(w)}[H(y_i|c_i, w)] \\ &= H(y_i|c_i) - u_{al}(y_i|c_i) \end{aligned} \tag{5}$$

In Eq. (4), $w$ denotes the model weights and $q(w)$ denotes the posterior distribution of weights in the training data. If the weights are fixed, $H(y_i|c_i, w)$ represents the uncertainty related to the data. Aleatoric uncertainty can be written as $\mathbb{E}_{q(w)}[H(y_i|c_i, w)]$ and calculated by the mean of $H_m(y_i|c_i)$. Epistemic uncertainty can also be written by the difference between the entropy $H(y_i|c_i)$ of $p(y_i|c_i)$ and aleatoric uncertainty in Eq. (5).

A recent study shows that the model pre-trained with a large amount of image captioning data incorporates commonsense knowledge that is implicit in the data [21]. We use such a pre-trained model (with commonsense knowledge) to generate captions including knowledge context from the KVQA image data, and predict the uncertainty of the knowledge for the given VQA using the above ensemble model. The proposed method

is illustrated in Fig. 2. As shown in Fig 2, when the image that birds are flying over the sand beach is given, the generated caption with commonsense knowledge is that two birds flying over a beach with a ship in the water. The caption reflects a general knowledge that ships are on a beach.

*2) Measuring similarity between caption sentences:* In addition to the above uncertainty model, this study proposes a novel measure that predicts the uncertainty of the knowledge based on the similarity between the generated and the ground-truth caption. That is, if the generated caption with the pre-trained model is much different from the ground-truth caption, the generated commonsense knowledge may be not much of use for the given problem. The S-BERT sentence embedding method [22] is used to calculate the caption similarity. The similarity between the caption embeddings is calculated as follows:

$$sim^{cap}(S_g, S_t) = \frac{f(S_g) \cdot f(S_t)}{\|f(S_g)\| \cdot \|f(S_t)\|}, f : encoder \tag{6}$$

In Eq. (6), $S_g$ and $S_t$ denote the generated caption and the ground-truth caption, respectively. $f$ is an encoder for extracting a representation. The similarity is calculated from dot product between the representations of the generated caption and the ground-truth.

### B. Knowledge-based visual question answering

Based on the above uncertainty measures, we present a new approach that integrates implicit knowledge and explicit knowledge external KB into KVQA.

*1) Use of knowledge based on semantic consistency:* As shown in Fig. 3, we adjust the use of the given KB based on the above mentioned uncertainty measures. In KVQA, when the uncertainty is high and the similarity is low, the scores of the metrics mean that the meaning of generated caption with commonsense knowledge is far from that of the image. Therefore, the system tends to attend to the content of image-question information. Otherwise, the external knowledge is more attended.

$$\begin{aligned} v^{score} &= \sigma(W_v * [sim^{cap}, u]) \\ g^{score} &= \sigma(W_g * [sim^{cap}, u]) \\ \mathbf{z}_v^{implicit} &= v^{score} * \mathbf{z}^{implicit} \\ \mathbf{z}_g^{explicit} &= g^{score} * \mathbf{z}^{explicit} \end{aligned} \tag{7}$$

As shown in Fig. 3 and Eq. (7), $v^{score} \in \mathbb{R}$ and $g^{score} \in \mathbb{R}$ are calculated through a fully connected layer and a sigmoid function $\sigma$ after concatenating the similarity $sim^{cap} \in \mathbb{R}$ and the uncertainty $u \in \mathbb{R}$. Each of the score values is finally represented by $\mathbf{z}_v^{implicit}$ and $\mathbf{z}_g^{explicit}$ through a dot product between the pre-calculated representation $\mathbf{z}^{implicit} \in \mathbb{R}^{d_{zi}}$ extracted from a vision-language model and the knowledge representation $\mathbf{z}^{explicit} \in \mathbb{R}^{d_{ze}}$, where $\mathbf{z}^{explicit}$ is extracted from a knowledge graph. The use of image-question information and KB are adjusted based on the score.
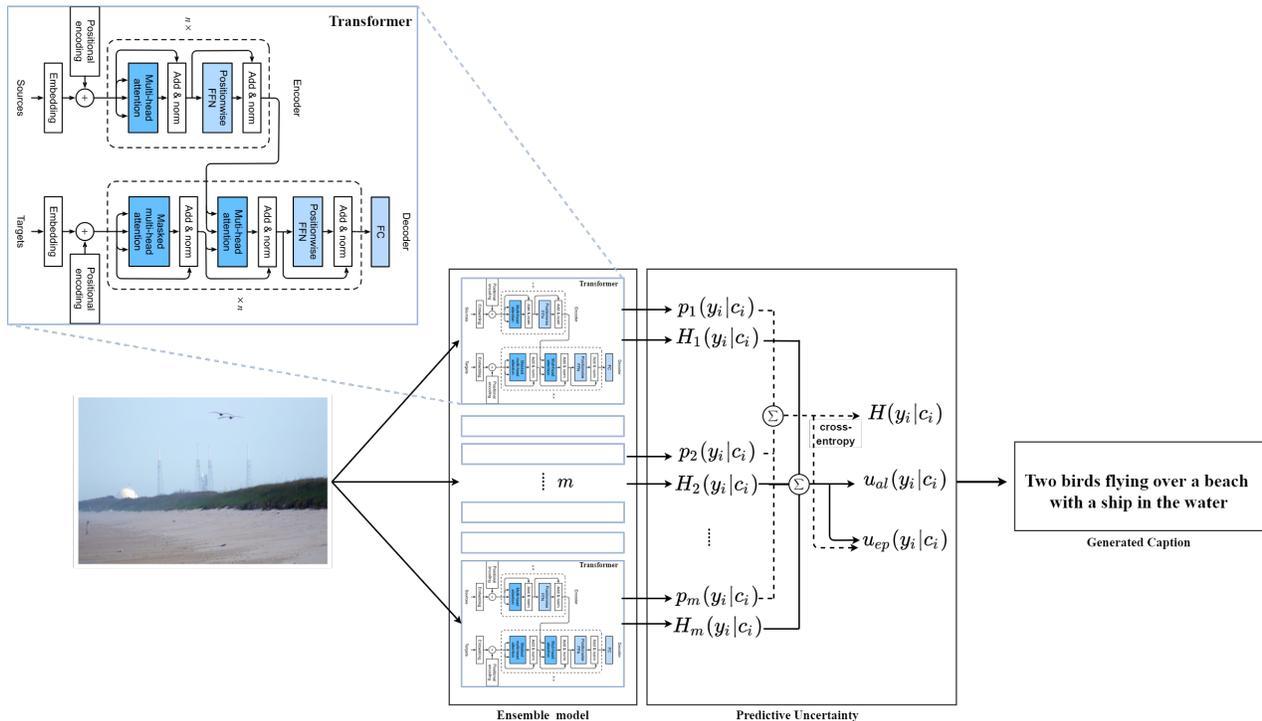
Fig. 2: Ensemble-based uncertainty estimation based on caption generation. The given image shows birds flying over a beach.
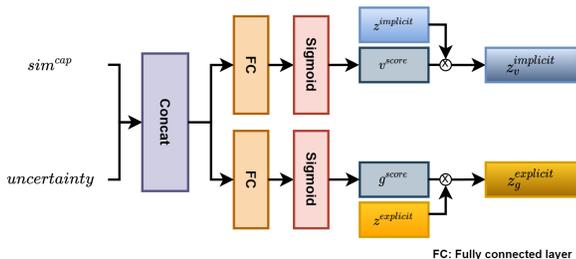


Fig. 3: Use of knowledge adjusted based on uncertainty measures.

*2) KVQA with semantic consistency model:* For KVQA, this study proposes a semantic consistency model that relies on the uncertainty measures described above. The model relies on two types of knowledge sources inspired by [5]: 1) explicit knowledge and 2) implicit knowledge. The former is the knowledge extracted from RGCN that has the external KB as input [23]. The latter is a vision-language embedding extracted from VisualBERT trained with a large-scale data. Furthermore, the use of the explicit and implicit knowledge is adjusted based on the uncertainty estimation, as described in section III-A.

**Explicit knowledge extraction**: Explicit knowledge is created by extracting relevant knowledge from the external KB, using the objects recognized in the image. In this study, about 4000 image keywords including objects, places, and attributes of objects are extracted with the following models: 1) ResNet-152 (ImageNet [24]); 2) ResNet-18 (Place365

[25]); 3) Faster R-CNN (VisualGenome [26]); 4) Mask-RCNN (LVIS [27]). External KB used are as follows: 1) DBPedia (categorical information) [28]; 2) ConceptNet (commonsense knowledge) [29]; 3) VisualGenome (spatial relationship) [26]; 4) hasPartKB (part relationship) [30]. The relevant knowledge is retrieved with image keywords and question words. As a result, a total of 36,000 edges and 8,000 nodes are extracted. For integrating knowledge graphs, we use RGCN that distinguishes types and directions of edges in this study. The followings are used as RGCN inputs: 1) keyword presence that indicates words in the question with filtered words with one-hot matrix; 2) an image keyword probability extracted from a pre-trained model; 3) Word2vec representation of each keyword or average Word2vec representation of multiple words [31]; 4) implicit knowledge representation $\mathbf{z}^{implicit}$ extracted from VisualBERT. The extracted explicit and implicit knowledge are integrated into KVQA as described above.

**Implicit knowledge extraction**: Transformer-based language models trained with a large-scale corpus are known to learn commonsense. Therefore, we use the VisualBERT model to make use of the implicit knowledge generated from the image and the question [7], as shown in Fig. 4. Although there are various studies that align images and sentences together, we apply the appropriate model to our task using experiments in [8]. The question representations are extracted by the pre-trained BERT model with BookCorpus dataset and English Wikipedia, and we use the representations as the input to the VisualBERT model. Furthermore, the visual representations are extracted from the Faster R-CNN model pre-trained with
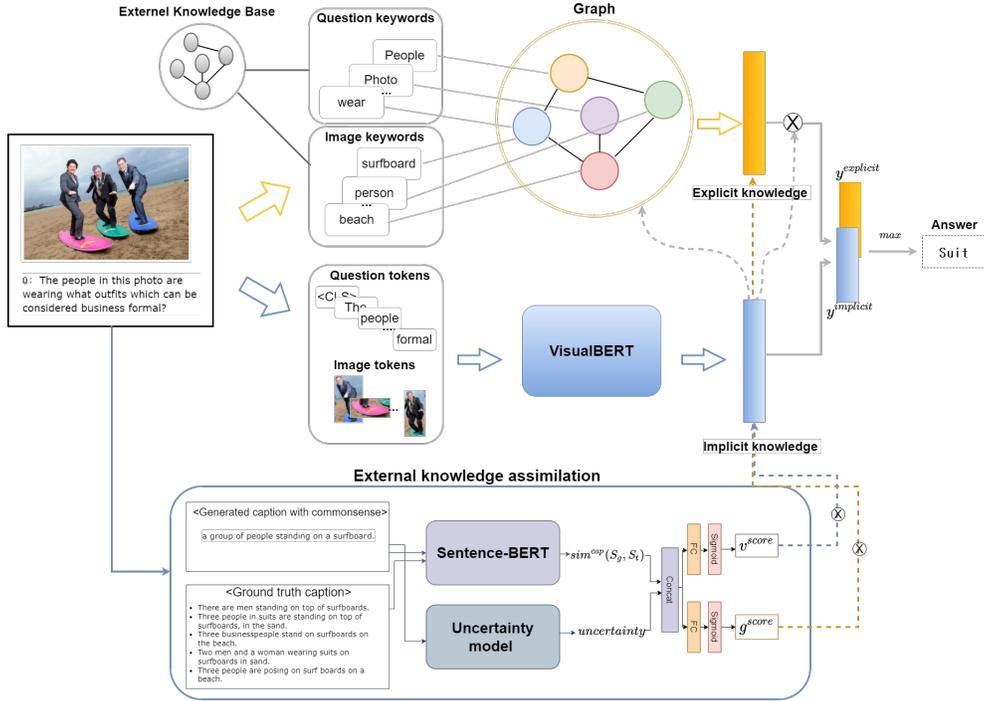
Fig. 4: Overall model architecture. In the external knowledge assimilation, when the image is given, the generated caption with commonsense is a group of people standing on a surfboard. The five ground-truth captions of the image are as follows: 1) there are men standing on top of surfboards; 2) three people in suits are standing on top of surfboards, in the sand; 3) three business people stand on surfboards on the beach; 4) two men and a woman wearing suits on surfboards in sand; 5) three people are posing on surfboards on a beach. A final answer is predicted with explicit knowledge and implicit knowledge combining the external knowledge assimilation.

VisualGenome/COCO dataset and the result becomes the input of VisualBERT. To produce $\mathbf{z}^{implicit}$ representation, we use mean-pooling with outputs extracted from the VisualBERT model.

To get a final answer, we predict the answer within a set of vocabulary of answers $V \in \mathbb{R}^v$ where $v$ is the size of vocabulary. The final implicit score $y^{implicit}$ and explicit score $y_i^{explicit}$ are calculated to predict the answer from the set $V \in \mathbb{R}^v$ as follows.

$$y^{implicit} = \sigma(W * \mathbf{z}_v^{implicit} + b) \qquad (8)$$

$$y_i^{explicit} = \sigma((W_{ge} * \mathbf{z}_{i,g}^{explicit} + b_{ge})^T \\ (W_{vi} * \mathbf{z}_v^{implicit} + b_{vi})) \qquad (9)$$

In Eq. (8)-(9), $y^{implicit}$ is calculated with a fully connected layer with weight $W$ and bias $b$, and a sigmoid function. In addition, $y_i^{explicit}$ is a score of word $i$ corresponding to $V \in \mathbb{R}^v$, computed with linear transformations of $\mathbf{z}_{i,g}^{explicit}$ and $\mathbf{z}_v^{implicit}$. The final answer is selected by choosing the highest value from both $y^{implicit}$ and $y^{explicit}$. The model is trained with binary cross-entropy.

TABLE I: Table of OK-VQA dataset.

| Dataset | # of images | # of questions |
|---------|-------------|----------------|
| Train   | 8,998       | 9,009          |
| Test    | 5,033       | 5,046          |
| Total   | 14,031      | 14,055         |

## IV. EXPERIMENTS AND RESULTS

### A. Dataset and baseline

We use the OK-VQA dataset [1] which is a popular KVQA benchmark dataset. The dataset consists of a total of 14,031 images and 14,055 questions. The detailed dataset sizes for training and testing are shown in Table I. For the validation dataset, we also use 1/3 of the training dataset based on the number of questions.

MSCOCO dataset [32] is used to pre-train baseline models that generate captions. The dataset size is shown in Table II. In addition, Att2in [33], BuDn [34], and Transformer [35] are selected as the baseline models for caption generation, which are the representative image captioning models, and are used to generate captions of the OK-VQA dataset.

TABLE II: Table of MSCOCO dataset.

| Dataset | # of images | # of captions |
|---------|-------------|---------------|
| Train | 82,783 | 413,915 |
| Validation | 40,504 | 202,520 |
| Test | 40,775 | 379,249 |
| Total | 164,062 | 995,684 |

TABLE III: Table of pearson correlation with the uncertainty and the similarity. $sim^{cap}$ represents a caption similarity. $un^{al}$ and $un^{ep}$ represent aleatoric uncertainty and epistemic uncertainty, respectively.

| | Corr |
|---|------|
| $sim^{cap}$ & $un^{al}$ | -0.1907 |
| $sim^{cap}$ & $un^{ep}$ | -0.1653 |
| $un^{al}$ & $un^{ep}$ | 0.4518 |

### B. Metrics

In this study, a standard evaluation metric used in VQA challenge [36] is employed to evaluate the performance with the OK-VQA dataset. Furthermore, we evaluate the generated caption with BLEU [37], CIDER [38], METEOR [39], and ROUGE-L [40] metrics.

### C. Uncertainty-based caption generation

Table IV shows the performances of the baseline model for caption generation with the OK-VQA dataset. When we compared the image caption performance of the Att2in, BuDn, and Transformer models with the OK-VQA dataset, overall, the Transformer model shows better performance than others, and our study uses the Transformer model for uncertainty modeling. Fig. 7 shows aleatoric uncertainty and epistemic uncertainty of the word in the generated caption, and the word for uncertain actions and unusual objects in the image shows higher uncertainty than the average uncertainty of the sentence.

Table III illustrates the pearson correlation between uncertainty and caption similarity. The caption similarity and aleatoric uncertainty have a negative correlation of -0.1907, and the correlation between similarity and epistemic uncertainty is -0.1653. The correlation between aleatoric uncertainty and epistemic uncertainty shows a positive correlation, with a value of 0.4518. The correlation analysis indicates that there are relations between caption similarity and uncertainty, as we expected. In Fig. 5, the distributions of (a) aleatoric uncertainty and (b) epistemic uncertainty are right-skewed, while in (c) caption similarity distribution presents a left-skewed shape. Since there are extreme values in distributions, we believe that the semantic inconsistency can be identified with the uncertainties of the caption and the caption similarity.

We also analyzed the uncertainty relationship according to the number of hallucinated objects in the generated caption as shown in Fig. 6. In Fig. 6, the x-axis categories mean
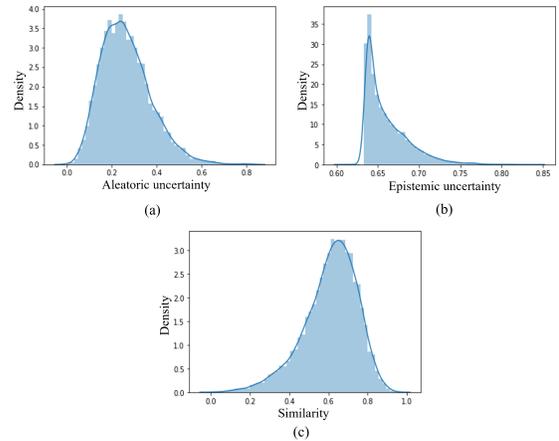


Fig. 5: The distribution of the uncertainty of generated caption and the similarity between the caption and the ground-truth caption.
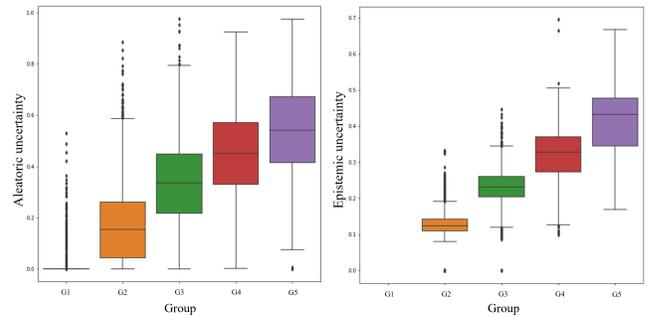


Fig. 6: Boxplot of the uncertainty value according to the number of hallucinated objects.

the following ranges, and the range means the ratio of the number of hallucinated words among the generated caption: 1) $0 \leq$ G1 $< 0.2$; 2) $0.2 \leq$ G2 $< 0.4$; 3) $0.4 \leq$ G3 $< 0.6$; 4) $0.6 \leq$ G4 $< 0.8$; 5) $0.8 \leq$ G5 $\leq 1.0$. The proportion of hallucinated objects of generated captions is calculated according to a synonym criteria of [41]. After synonym filtering of the generated caption, the number of hallucinated objects in the generated caption is counted. We divide the ratio of the number of words of the hallucinated objects among the caption words into 5 groups. We calculate the average uncertainty of the caption over the average uncertainty of the hallucinated objects. As shown in Fig. 6, the more hallucinated objects in the caption, the higher aleatoric and epistemic uncertainty. We also performed a qualitative analysis, as shown in Fig. 7. For the example shown in Fig. 7, the generated caption contains uncertain words with higher aleatoric and epistemic uncertainty than $m$ the average aleatoric uncertainty and the average epistemic uncertainty in a sentence.

### D. KVQA with semantic inconsistency

*1) Comparison with state-of-the-art approaches:* We compare our proposed semantic inconsistency model with the

TABLE IV: Performances of image captioning with commonsense knowledge on the OK-VQA dataset.

| Model | BLEU-1 | BLEU-2 | BLEU-3 | BLEU-4 | CIDER | METEOR | ROUGE-L |
|---|---|---|---|---|---|---|---|
| Att2in [33] | 0.7843±0.00005 | 0.6077±0.0002 | 0.4508±0.00032 | 0.3302±0.00038 | 1.0833+0.0016 | 0.2604±0.00018 | 0.5561±0.00019 |
| BuDn [34] | 0.8123±0.00015 | 0.6516±0.00009 | 0.5017±0.00003 | 0.3786±0.0004 | 1.2527±0.00039 | 0.2858±0.00002 | 0.5859±0.00005 |
| Transformer [35] | **0.8290±0.00028** | **0.6828±0.00036** | **0.5410±0.0004** | **0.4216±0.0004** | **1.3864±0.0012** | **0.2997±0.00013** | **0.6043±0.00023** |



Caption: A fire **(0.40, 0.68)** hydrant on the side of a street, m = (0.14, 0.64)

(a)

Caption: A white bird is standing on the top **(1.12, 0.82)** of an oven, m = (0.36, 0.66)

(b)

Caption: A person **(0.812, 0.75)** is holding a teddy bear, m = (0.24, 0.67)

(c)

Fig. 7: Image captioning results on OK-VQA dataset. Values in bracket are aleatoric uncertainty and epistemic uncertainty, respectively, and $m$ represents an average aleatoric uncertainty and an average epistemic uncertainty in a sentence, respectively.



Q: Which item in this room is usually to wash hands?
A: sink
Baseline: bed
Ours: sink

(a)

Q: What is the purpose of these objects?
A: [decoration, art]
Baseline: tell time
Ours: decoration

(b)

Q: What is the purpose of the logos on this truck?
A: identification
Baseline: car
Ours: safety

(c)

Fig. 8: Comparison with the predicted answers of the proposed method and the baseline model on OK-VQA dataset.

following state-of-the-art approaches including those with pre-trained methods and a combination of graph-based and pre-trained methods: 1) BAN [9]: Bilinear attention network which uses co-attention module with question features and image features from pre-trained models; 2) BAN+AN [1]: The model incorporates the external knowledge into BAN by using ArticleNet; 3) BAN+KG-Aug [4]: The model incorporates knowledge graph augmented model into BAN by using late augmentation scheme; 4) MUTAN [10]: Multimodal tucker fusion network which focuses on image and textual features extracted from pre-trained models based on tucker decomposition; 5) MUTAN+AN [1]: Similarly with BAN+AN, this method also incorporates the external knowledge into MUTAN by using ArticleNet; 6) KA [3]: The model uses image features, question features, and concept graphs with the external knowledge; 7) KRISP [5]: The model integrates image-text representation extracted from the BERT-based model and graph information based on external knowledge. In general, the methods using the knowledge information show better performance. As shown in Table VI, the model with both explicit knowledge, implicit

knowledge, and semantic inconsistency measure achieves the state-of-the-art performance.

*2) Ablation study:* An ablation study is performed with three values of caption similarity, aleatoric uncertainty, and epistemic uncertainty with the weights in Eq. (7). In Table V, the baseline model that makes use of both explicit and implicit knowledge shows an accuracy of 31.15%. When caption similarity is added, the accuracy increases by 0.4%. In addition, when aleatoric and epistemic uncertainty are added, respectively, it shows further improvement. Also, when the similarity and epistemic uncertainty are added, the accuracy increases by 0.49%. The best performance of 32.45% in accuracy is achieved when the caption similarity and the aleatoric uncertainty are concatenated. As shown in Table III, since the similarity and aleatoric uncertainty has a higher correlation than between the similarity and epistemic uncertainty, the concatenated model seems to provide the best performance. When the three values of caption similarity, aleatoric uncertainty, and epistemic uncertainty are used together, the accuracy is 31.19%, which is only slightly better than the baseline. These results indicate that the caption similarity captures the

semantic inconsistency relatively well and when the value which has a high correlation with the similarity is given to the model, it can predict correct answers better.

TABLE V: An ablation study of the external knowledge assimilation methods with the OK-VQA dataset.

| Model | Accuracy |
|---|---|
| Baseline | 31.15 |
| Baseline + $sim^{cap}$ | 31.55 |
| Baseline + $uncertainty^{al}$ | 31.28 |
| Baseline + $uncertainty^{ep}$ | 31.93 |
| Baseline + $sim^{cap}$ + $uncertainty^{ep}$ | 31.64 |
| Baseline + $sim^{cap}$ + $uncertainty^{al}$ | **32.45** |
| Baseline + $sim^{cap}$ + $uncertainty^{ep}$ + $uncertainty^{al}$ | 31.19 |

TABLE VI: Results with the OK-VQA dataset, comparing our work with the state-of-the-art approaches. * represents results from a re-implementation with the author's code and parameter setting using three experiments.

| Model | Accuracy |
|---|---|
| Q-Only | 14.93 |
| BAN [9] | 25.17 |
| BAN + AN [1] | 25.61 |
| MUTAN [10] | 26.41 |
| BAN + KG-Aug [4] | 26.71 |
| MUTAN + AN [1] | 27.84 |
| KA [3] | 29.03 |
| KRISP* [5] | 31.15 |
| Ours | **32.45** |

*3) Qualitative results:* We also present a qualitative analysis of the model in Fig. 8. We compare the prediction from our model with the baseline's. For (a) and (b), our model selects the correct answer. In addition, for (c) our model predicts an answer that is more similar to the correct answer than the baseline model. Also, in Fig. 8, the proposed method predicts correct answers even when the image shows a part of the sink (a), and with an unusual combination of objects and the background (b).

## V. Conclusion and future work

In this study, we propose a novel semantic inconsistency measure through uncertainty modeling and semantic similarity for KVQA that can make use of diverse KBs more effectively. As KBs are often incomplete or incompatible with the given problem, the use of knowledge should be moderated. With the proposed model, we achieve the state-of-the-art results on KVQA. As a future work, we plan to further explore diverse ways of using KBs based on the characteristics of the KB and the given problem.

## VI. Acknowledgements

## References

[1] Kenneth Marino, Mohammad Rastegari, Ali Farhadi, and Roozbeh Mottaghi. Ok-vqa: A visual question answering benchmark requiring external knowledge. In *Proc. of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3195–3204, 2019.

[2] Liyang Zhang, Shuaicheng Liu, Donghao Liu, Pengpeng Zeng, Xiangpeng Li, Jingkuan Song, and Lianli Gao. Rich visual knowledge-based augmentation network for visual question answering. *IEEE Transactions on Neural Networks and Learning Systems*, 32(10):4362–4373, 2021.

[3] Maryam Ziaeefard and Freddy Lécué. Towards knowledge-augmented visual question answering. In *Proc. of the 28th International Conference on Computational Linguistics*, pages 1863–1873, 2020.

[4] Guohao Li, Xin Wang, and Wenwu Zhu. Boosting visual question answering with context-aware knowledge aggregation. In *Proc. of the 28th ACM International Conference on Multimedia*, pages 1227–1235, 2020.

[5] Kenneth Marino, Xinlei Chen, Devi Parikh, Abhinav Gupta, and Marcus Rohrbach. Krisp: Integrating implicit and symbolic knowledge for open-domain knowledge-based vqa. In *Proc. of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 14111–14121, 2021.

[6] Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. Vilbert: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2019.

[7] Liunian Harold Li, Mark Yatskar, Da Yin, Cho-Jui Hsieh, and Kai-Wei Chang. Visualbert: A simple and performant baseline for vision and language. *arXiv preprint arXiv:1908.03557*, 2019.

[8] Amanpreet Singh, Vedanuj Goswami, and Devi Parikh. Are we pre-training it right? digging deeper into visio-linguistic pretraining. *arXiv preprint arXiv:2004.08744*, 2020.

[9] Jin-Hwa Kim, Jaehyun Jun, and Byoung-Tak Zhang. Bilinear attention networks. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 1564–1574, 2018.

[10] Hedi Ben-Younes, Rémi Cadene, Matthieu Cord, and Nicolas Thome. Mutan: Multimodal tucker fusion for visual question answering. In *Proc. of the IEEE International Conference on Computer Vision (ICCV)*, pages 2612–2620, 2017.

[11] Drew A. Hudson and Christopher D. Manning. Learning by abstraction: The neural state machine. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2019.

[12] Pin Jiang and Yahong Han. Reasoning with heterogeneous graph alignment for video question answering. *Proc. of the AAAI Conference on Artificial Intelligence*, 34(07):11109–11116, Apr. 2020.

[13] Noa García and Yuta Nakashima. Knowledge-based video question answering with unsupervised scene descriptions. In *Proc. of of the European Conference on Computer Vision(ECCV)*, 2020.

[14] Difei Gao, Ke Li, Ruiping Wang, Shiguang Shan, and Xilin Chen. Multimodal graph neural network for joint reasoning on vision and scene text. In *Proc. of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 12743–12753, 2020.

[15] Raeid Saqur and Karthik Narasimhan. Multimodal graph networks for compositional generalization in visual question answering. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 33, pages 3070–3081, 2020.

[16] Yijun Xiao and William Yang Wang. On hallucination and predictive uncertainty in conditional language generation. *arXiv preprint arXiv:2103.15025*, 2021.

[17] Armen Der Kiureghian and Ove Ditlevsen. Aleatory or epistemic? does it matter? *Structural Safety*, 31(2):105–112, 2009. Risk Acceptance and Risk Communication.

[18] Stefan Depeweg, Jose-Miguel Hernandez-Lobato, Finale Doshi-Velez, and Steffen Udluft. Decomposition of uncertainty in Bayesian deep learning for efficient and risk-sensitive learning. In *Proc. of the International Conference on Machine Learning (ICML)*, pages 1184–1193, 2018.

[19] Alex Kendall and Yarin Gal. What uncertainties do we need in bayesian deep learning for computer vision? In *Advances in Neural Information Processing Systems (NeurIPS)*, page 5580–5590, 2017.

[20] Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable predictive uncertainty estimation using deep ensembles. In *Advances in Neural Information Processing Systems (NeurIPS)*, page 6405–6416, 2017.

[21] Weijie Su, Xizhou Zhu, Yue Cao, Bin Li, Lewei Lu, Furu Wei, and Jifeng Dai. Vl-bert: Pre-training of generic visual-linguistic representations. In *Proc. of the International Conference on Learning Representations (ICLR)*, 2020.

[22] Nils Reimers and Iryna Gurevych. Sentence-bert: Sentence embeddings using siamese bert-networks. *arXiv preprint arXiv:1908.10084*, 2019.

[23] Michael Schlichtkrull, Thomas N Kipf, Peter Bloem, Rianne Van Den Berg, Ivan Titov, and Max Welling. Modeling relational data with graph convolutional networks. In *European Semantic Web Conference*, pages 593–607. Springer, 2018.

[24] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Li Kai, and Fei-Fei Li. Imagenet: A large-scale hierarchical image database. In *Proc. of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 248–255, 2009.

[25] Alejandro López-Cifuentes, Marcos Escudero-Viñolo, Jesús Bescós, and Álvaro García-Martín. Semantic-aware scene recognition. *Pattern Recognition*, 102:107256, 2020.

[26] Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A. Shamma, Michael S. Bernstein, and Fei-Fei Li. Visual genome: Connecting language and vision using crowdsourced dense image annotations. *International Journal of Computer Vision*, 123(1):32–73, 2017.

[27] Agrim Gupta, Piotr Dollár, and Ross B. Girshick. Lvis: A dataset for large vocabulary instance segmentation. In *Proc. of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5351–5359, 2019.

[28] Sören Auer, Christian Bizer, Georgi Kobilarov, Jens Lehmann, Richard Cyganiak, and Zachary Ives. Dbpedia: A nucleus for a web of open data. In *The Semantic Web*, pages 722–735. Springer Berlin Heidelberg, 2007.

[29] H. Liu and P. Singh. Conceptnet — a practical commonsense reasoning tool-kit. *BT Technology Journal*, 22(4):211–226, 2004.

[30] Sumithra Bhakthavatsalam, Kyle Richardson, Niket Tandon, and Peter Clark. Do dogs have whiskers? A new knowledge base of haspart relations. *arXiv preprint arXiv:2006.07510*, 2020.

[31] Tomás Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient estimation of word representations in vector space. In *Proc. of the International Conference on Learning Representations (ICLR)*, 2013.

[32] Xinlei Chen, Hao Fang, Tsung-Yi Lin, Ramakrishna Vedantam, Saurabh Gupta, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco captions: Data collection and evaluation server. *arXiv preprint arXiv:1504.00325*, 2015.

[33] Steven J Rennie, Etienne Marcheret, Youssef Mroueh, Jerret Ross, and Vaibhava Goel. Self-critical sequence training for image captioning. In *Proc. of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 7008–7024, 2017.

[34] 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 *Proc. of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6077–6086, 2018.

[35] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 5998–6008, 2017.

[36] Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C Lawrence Zitnick, and Devi Parikh. Vqa: Visual question answering. In *Proc. of the IEEE International Conference on Computer Vision (ICCV)*, pages 2425–2433, 2015.

[37] Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic evaluation of machine translation. In *Proc. of the Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 311–318, 2002.

[38] Ramakrishna Vedantam, C Lawrence Zitnick, and Devi Parikh. Cider: Consensus-based image description evaluation. In *Proc. of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4566–4575, 2015.

[39] Satanjeev Banerjee and Alon Lavie. Meteor: An automatic metric for mt evaluation with improved correlation with human judgments. In *Proc. of the ACL workshop on intrinsic and extrinsic evaluation measures for machine translation and/or summarization*, pages 65–72, 2005.

[40] Chin-Yew Lin. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*, pages 74–81, 2004.

[41] Anna Rohrbach, Lisa Anne Hendricks, Kaylee Burns, Trevor Darrell, and Kate Saenko. Object hallucination in image captioning. *arXiv preprint arXiv:1809.02156*, 2018.