

Alleviating the Burden of Labeling: Sentence Generation by Attention Branch Encoder–Decoder Network

Tadashi Ogura¹, Aly Magassouba¹, Komei Sugiura^{1,2}, Tsubasa Hirakawa³, Takayoshi Yamashita³,
Hironobu Fujiyoshi³, and Hisashi Kawai¹

Abstract—Domestic service robots (DSRs) are a promising solution to the shortage of home care workers. However, one of the main limitations of DSRs is their inability to interact naturally through language. Recently, data-driven approaches have been shown to be effective for tackling this limitation; however, they often require large-scale datasets, which is costly. Based on this background, we aim to perform automatic sentence generation of fetching instructions: for example, “Bring me a green tea bottle on the table.” This is particularly challenging because appropriate expressions depend on the target object, as well as its surroundings. In this paper, we propose the attention branch encoder–decoder network (ABEN), to generate sentences from visual inputs. Unlike other approaches, the ABEN has multimodal attention branches that use subword-level attention and generate sentences based on subword embeddings. In experiments, we compared the ABEN with a baseline method using four standard metrics in image captioning. Results show that the ABEN outperformed the baseline in terms of these metrics.

Index Terms—Novel Deep Learning Methods, Deep Learning for Visual Perception

I. INTRODUCTION

THE growth in the aged population has steadily increased the need for daily care and support. Domestic service robots (DSRs) that can physically assist people with disabilities are a promising solution to the shortage of home care workers [1]–[3]. This has boosted the need for standardized DSRs that can provide necessary support functions.

Nonetheless, one of the main limitations of DSRs is their inability to interact naturally through language. Indeed, most DSRs do not allow users to instruct them with diverse expressions. Recent studies have shown that data-driven approaches are effective for handling ambiguous instructions [4]–[7].

Unfortunately, these approaches often require large-scale datasets, and are time-consuming and costly. The main reason is the time that is required for human experts to provide sentences for images. Hence, methods to augment or generate

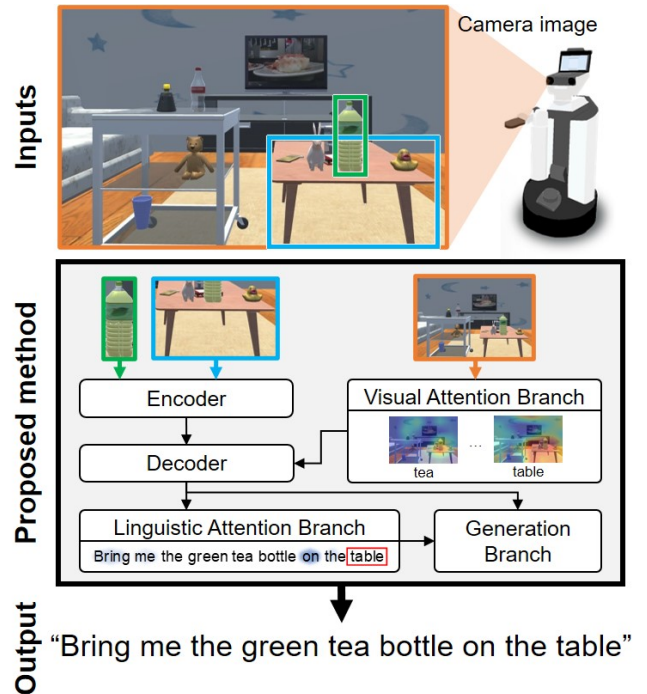


Fig. 1. Overview of the ABEN: the ABEN generates fetching instructions from given input images.

instructions automatically could drastically reduce this cost and alleviate the burden of labeling from human experts.

Based on this background, we aim to perform automatic sentence generation of “fetching instructions” (instructions to the DSR to fetch items). This task involves generating a natural fetching instruction, given a target object in an image: for example, “Bring me a green tea bottle on the table.” Such an instruction often includes a referring expression, such as “a green tea bottle on the table.” A referring expression is an expression in which an object is described with regard to a landmark, such as “table”. This is particularly challenging because of the many-to-many mapping between language and the environment.

In this paper, we propose the attention branch encoder–decoder network (ABEN), which generates fetching instructions from visual inputs. Fig. 1 shows a schematic diagram of the approach. The ABEN comprises a visual attention branch (VAB) and a linguistic attention branch (LAB), to attend both visual and linguistic inputs. An additional generation branch is introduced to generate sentences.

The ABEN extends the attention branch network (ABN) [8]

Manuscript received: February 22, 2020; Revised June 2, 2020; Accepted July 7, 2020. This paper was recommended for publication by Editor Dongheui Lee upon evaluation of the Associate Editor and Reviewers’ comments.

¹Authors are with the National Institute of Information and Communications Technology, 3-5 Hikaridai, Seika, Soraku, Kyoto 619-0289, Japan. firstname.lastname@nict.go.jp

²Author is with Keio University, 3-14-1 Hiyoshi, Kohoku, Yokohama, Kanagawa 223-8522, Japan. firstname.lastname@keio.jp

³Authors are with Chubu University, 1200 Matsumotocho, Kasugai, Aichi 487-8501, Japan. {hirakawa@mpg.cs, takayoshi@isc, fujiyoshi@isc}.chubu.ac.jp

Digital Object Identifier (DOI): see top of this page.

by introducing multimodal attention branches. In the ABN, attention maps are output by an attention branch, which highlights the most salient portions in the image, given a label to predict. In the ABEN, the attention map given by the VAB can serve as a visual explanation of the model, which is usually a black box. Similarly, the attention map provided by the LAB can serve as an explanation of subword relationships.

The ABEN was inspired by our previous approach, the multimodal ABN (Multi-ABN) [9], and shares its basic structure. The main differences between the ABEN and the Multi-ABN include the subword-level attention used in the LAB and BERT [10]-based subword embeddings used for sentence generation. A demonstration video is available at this URL¹.

The main contributions of this paper are as follows:

- The ABEN extends the Multi-ABN by introducing a linguistic branch and a generation branch, to model the relationship between subwords.
- The ABEN combines attention branches and BERT-based subword embedding, for sentence generation.

II. RELATED WORK

There have been many attempts to construct communicative robots for manipulation tasks [3]. Recently, in some studies [6], [7], [11]–[13], linguistic inputs were processed along with visual information to handle the many-to-many mapping between language and the environment.

These studies have often used data-driven approaches that were originally proposed in the natural language processing (NLP) and computer vision communities. For instance, [12] proposed a method for predicting target objects from natural language in a pick-and-place task environment, using a visual semantic embedding model. Similarly, [13] tackled the same type of problem using a two-stage model to predict the likely target from the language expression and the pairwise relationships between different target candidates. More recently, in a context related to DSRs, [7] proposed the use of both the target and source candidates to predict the likely target in a supervised manner. In [6], the placing task was addressed through a generative adversarial network (GAN) classifier that predicted the most likely destination from the initial instruction.

Nonetheless, these methods required large-scale datasets, which are seldom available in a DSR context because they require substantial labeling effort from human experts. Datasets such as RefCOCO [14] or MSCOCO [15] are widely used in visual captioning, however they are not specifically designed for robots. The Room-to-Room dataset [4] is a dataset designed for multimodal language understanding for navigation, however manipulation is not handled. Conversely, a pick-and-place dataset such as PFN-PIC [12], contains top-view images only, and does not handle furniture, and is therefore not suitable for DSRs.

To address this problem, a promising solution involves generating synthetic instructions to label unseen visual inputs to augment such datasets. Moreover, such a method enables real-time task generation in simulators, where DSRs are instructed to fetch everyday objects in randomly generated environments.



Fig. 2. Left: Typical scene in which the DSR is observing everyday objects. Right: The camera image recorded from the DSR’s position shown in the left-hand panel. The blue and green boxes represent the target (blue glass) and its source (metal wagon). Typical instructions include “Bring me the blue glass next to the teddy bear” and “Bring me the blue glass on the same level as the teddy bear on the metal wagon.”

Most studies use rule-based approaches to generate sentences (e.g. [16]),

however, they cannot fully capture and reproduce the many-to-many mapping between language and the physical world. Indeed, handling natural sentences that include referring expressions is particularly challenging. Conversely, an end-to-end approach was used in [17] for estimating spatial relations to describe an object in a sentence. Nonetheless, the set of spatial relations was limited to six and was hand-crafted. Unlike these studies, we target an end-to-end approach that does not require hand-crafted or rule-based methods. In our previous work [9], we introduced the Multi-ABN, which generates fetching instructions by using a multimodal attention branch mechanism [8]. The Multi-ABN is a long short-term memory (LSTM) that is enhanced by visual and linguistic attention branches.

This study extends the Multi-ABN by introducing subword-level attention, which has the benefit of interpretability, unlike the linguistic attention in [9]. Our approach can model the relationship between the generated subwords. Furthermore, unlike most sentence generation methods, our approach generates sentences via a BERT-based subword embedding [10] model, which was shown by [7] to perform better than a word embedding model.

Therefore, the architecture of the ABEN extends the ABN with multimodal attention. Multimodal attention has been widely investigated in image captioning. Recent studies in multimodal language understanding have shown that both linguistic and visual attention are beneficial for question-answering tasks [18], [19] or visual grounding [20], [21]. Similarly, [22] introduces an attention method that performs a weighted average of linguistic and image inputs. In contrast to these attention mechanisms, attention branches are based on class activation mapping (CAM) networks [23]. CAM focuses on the generation of masks that, overlaid onto an image, highlight the most salient area given a label. In the ABN, such a structure is introduced through an attention branch that generates attention maps to improve the prediction accuracy of the base network. In the ABEN, visual and linguistic attention maps are generated to mask the visual input and the sequence of generated subwords.

Subwords have been widely used for machine translation [24]–[26] as well as being used in most recent language models such as BERT, ALBERT [27] or XLNet [28]. These

¹<https://youtu.be/H7vsGmJaE6A>

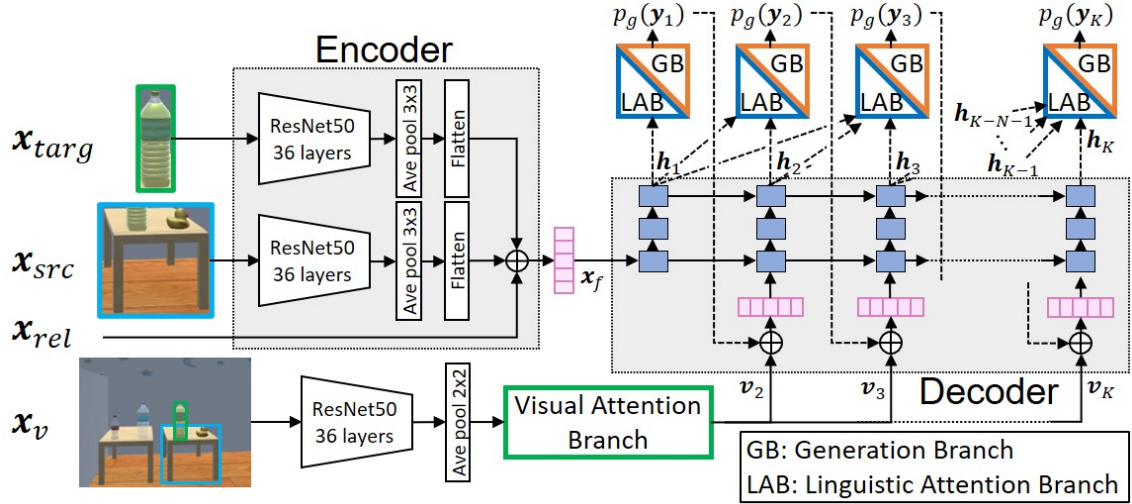


Fig. 3. Structure of the ABEN. The ABEN comprises an encoder, a decoder, a visual attention branch, a linguistic attention branch, and a generation branch.

methods achieve state-of-the-art performance in many natural language understanding tasks. Combining such a method with the attention branch architecture to enrich sentence generation, is one of the main novelties of the ABEN. Indeed, traditionally in robotics, very simple embedding and language models have been used. For example, recent studies used simple skip-gram (e.g. [12], [13], [29], [30]).

III. PROBLEM STATEMENT

A. Task description

The study targets natural language generation for DSRs. Hereinafter, we call this task the fetching instruction generation (FIG) task. A typical scenario is shown in Fig. 2, in which the target to fetch is described in the instruction “Bring me the blue glass on the same level as the teddy bear on the metal wagon.” This example emphasizes the challenges of the FIG task, because each instruction should describe the targeted object uniquely.

To avoid ambiguity, it is necessary to generate sentences including referring expressions, because there may be many objects of the same type. Referring expressions allow the targeted object to be characterized uniquely with respect to its surrounding environment. In Fig. 2, the referring expression “on the same level as the teddy bear on the metal wagon” is needed to disambiguate the targeted object from others.

This is particularly challenging because appropriate expressions depend on the target itself, as well as its surroundings. For instance, the target object in Fig. 2 can be described as “glass near the bear doll” and “blue tumbler glass on the second level of the wagon”, in addition to many other candidate expressions. Therefore, it is necessary to handle the many-to-many mapping between language and the physical world.

The FIG task is characterized by the following:

- **Input:** RGB image of an observed scene.
- **Output:** the most likely generated sentence for a given target and source.

The inputs of the ABEN are explained in detail in Section IV.

We define the terms used in this paper as follows:

- **Target:** an everyday object, e.g., bottles or fruit, that is to be fetched by the robot.
- **Source:** the origin of the target, e.g., furniture, such as shelves or drawers.

In the FIG task, we assume that the two-dimensional bounding boxes of the target and the source are defined in advance. Furthermore, referring expressions related to depth perception (e.g., “behind” or “in front”) are not addressed because no three-dimensional information is available.

The evaluation of the generated sentences is based on several standard metrics—BLEU [31], ROUGE [32], METEOR [33], and CIDEr [34]—that are commonly used for image captioning studies. Although they are imperfect, by using several metrics, we may overcome their limitation and assess better the quality of the sentences. In [33], BLEU and METEOR were reported to have a correlation of 0.817 and 0.964, respectively, with human evaluation. Furthermore, these metrics also allow us to compare our approach to existing methods.

A simulated environment (see Fig. 2) is used to collect the image inputs. Indeed, because we aim to generate sentences in a wide range of configurations, using a simulated environment is effective for addressing these situations at a low cost. Moreover, using a simulation has the advantage that the experimental results can be reproduced.

As the simulated robot platform, we use a standardized DSR, namely Human Support Robot (HSR) [35]. Our simulator is based on SIGVerse [36], [37], which is an official simulator for HSR that provides a three-dimensional environment based on the Unity engine.

In the data collection phase, HSR [35] navigates in procedurally generated environments with everyday objects. Thereafter, RGB images of target and source candidates are recorded using the camera, with which HSR is equipped.

IV. PROPOSED METHOD

A. Novelty

To generate fetching instructions, the ABEN extends the Multi-ABN [9] by introducing a subword generation architecture using BERT [10] embedding in addition to subword-level

TABLE I
DIFFERENCE BETWEEN (A) TYPICAL WORD-TOKENS WITH
PRE-PROCESSING FOR RARE AND/OR MISSPELLED WORDS AND
(B) SUB-WORD TOKENIZATION.

Expression	(a)	(b)
topright object	topright, object	top, right, object
sprayer	<UNK>	spray, er
grayis bottle	<UNK>, bottle	gray, is, bottle

attention. For this purpose, as shown in Fig. 3, the ABEN comprises an encoder (base network), a decoder, a VAB, and a LAB. The following characteristics of the ABEN should be emphasized:

- Unlike the ABN [8], which comprises a base network coupled with attention and perception branches, the ABEN follows an encoder–decoder structure (i.e., there is no perception branch) based on an LSTM network.
- Unlike the Multi-ABN, fetching instructions are generated from a sequence of subwords with BERT encoding.
- The ABEN introduces the novel structures of linguistic attention branches and generation branches to allow a subword-level attention mechanism. Hence, the ABEN attention is fully interpretable, unlike that of the Multi-ABN which uses latent-space linguistic attention.

B. Input and Subword Tokenization

Fig. 3 shows the network structure of the ABEN. For a scene i , let us define our set of inputs \mathbf{x}_i as:

$$\mathbf{x}_i = \{\mathbf{x}_v(i), \mathbf{x}_{src}(i), \mathbf{x}_{targ}(i), \mathbf{x}_{rel}(i)\}. \quad (1)$$

For readability, we omit the index i so that \mathbf{x}_i is simply written as \mathbf{x} . Given this, the inputs are defined as follows:

- \mathbf{x}_v : the input scene as an RGB image.
- \mathbf{x}_{targ} : the cropped image of the target in \mathbf{x}_v
- \mathbf{x}_{src} : the cropped image of the source in \mathbf{x}_v
- \mathbf{x}_{rel} : the relational features between \mathbf{x}_v , \mathbf{x}_{targ} , and \mathbf{x}_{src} .

\mathbf{x}_{rel} comprises the position and size features of (a) the target relative to the source, (b) the target relative to the full image, and (c) the source relative to the full image. Each of these relations is characterized by:

$$\mathbf{r}_{l/m} = \left[\frac{x_l}{W_m}, \frac{y_l}{H_m}, \frac{w_l}{W_m}, \frac{h_l}{H_m}, \frac{w_l h_l}{W_m H_m} \right], \quad (2)$$

where x_l , y_l , w_l , and h_l denote the horizontal and vertical positions and the width and height, respectively, of the component l . W_m and H_m denote the width and height, respectively, of the component m . Consequently, the relation features are defined as $\mathbf{x}_{rel} = \{\mathbf{r}_{targ/src}, \mathbf{r}_{targ/v}, \mathbf{r}_{src/v}\}$ with dimension 15.

In contrast to most methods for sentence generation, BERT-based subword embeddings, instead of classic word-based embedding, are used as the ground truth. BERT was pretrained on 3.5 billion words and is therefore robust against data sparseness regarding rare words. In our previous work on multimodal language understanding, we introduced BERT-based subword embedding; this was one of the earliest applications of BERT in robotics [7]. It has been reported that BERT-based subword embedding functioned better than simple word-based embedding for the PFN-PIC dataset [12]. In many NLP

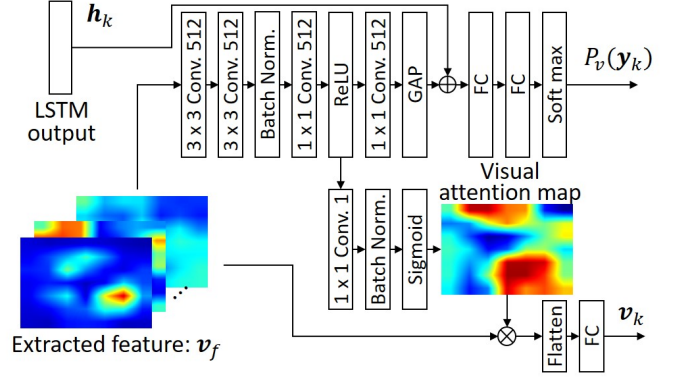


Fig. 4. Architecture of the visual attention branch.

studies, BERT and other Transformer-based approaches have been applied successfully to challenging tasks. For domain adaptation, we can additionally fine-tune a BERT-based model that is pre-trained on a large-scale dataset.

Furthermore, subword tokenization [38] is robust against the misspelling words. Indeed, a matching is still possible in subword units. As illustrated in Table I, a the word ‘grayish’ misspelled as ‘grayis’ can still be matched with the subword ‘gray’ which is impossible with classic word embedding. As a result, the subword tokenization and generation handle more word variations because there is no need to perform stopword replacement or stemming (e.g., for conjugated verbs).

C. Structure

1) *Encoder*: The encoder transforms visual information into a latent space feature that is later decoded as a sentence by the decoder. The inputs of the encoder are the target \mathbf{x}_{targ} , source \mathbf{x}_{src} and relation features \mathbf{x}_{rel} as illustrated in Fig. 3. A feature \mathbf{x}_f is generated by the encoder. To do so, the target and source images are both encoded by a convolutional neural network (CNN). In this study, we use ResNet-50 [39] as the backbone neural network. The encoding process involves extracting the output of the 36th layer of ResNet-50, which is followed by a global average pooling (GAP) and a flattening process for dimension reduction. Feature \mathbf{x}_f is then obtained as the concatenation of the two encoded visual features with the relation feature \mathbf{x}_{rel} .

2) *Decoder*: The decoder generates a sequence of latent-space features $\mathbf{H} = \{\mathbf{h}_1, \mathbf{h}_2, \dots, \mathbf{h}_K\}$, for each step k , from the encoded feature \mathbf{x}_f by using a multi layer LSTM. These latent-space features allow the linguistic attention and generation branches to predict a sequence of subwords corresponding to the fetching instruction. For that purpose, each cell of the LSTM, at step k , is initialized with the embedding vector of the previous subword predicted \mathbf{y}_{k-1} , as well as a visual feature \mathbf{v}_k obtained from the VAB. Feature \mathbf{v}_k is described below with the VAB structure. Thereafter, the hidden state of the LSTM propagates as shown in Fig. 3 and the output \mathbf{h}_k is generated for each step k .

3) *Visual Attention Branch*: Fig. 4 shows the structure of the VAB. The VAB used in this study is based on the method proposed in [9]. From the VAB, informative regions of features extracted from the image \mathbf{x}_v are emphasized to predict the subword \mathbf{y}_k . Similarly to the encoder, the input \mathbf{x}_v is processed by the 36th layer of ResNet-50 and generates feature maps

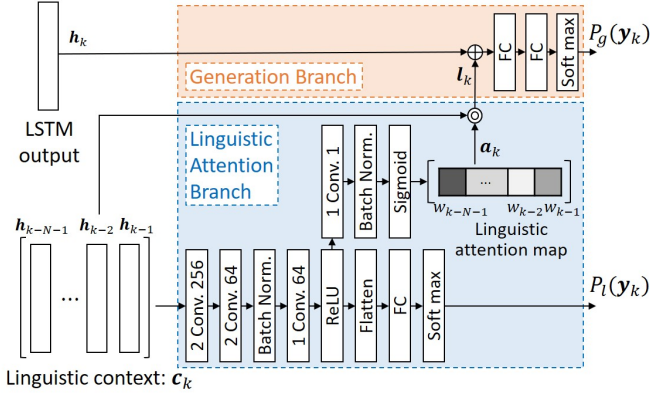


Fig. 5. Architecture of the linguistic and generation branches. The output of the LAB is obtained from Equation (6) which is represented as ‘ \odot ’.

\mathbf{v}_f . These feature maps have a dimension $7 \times 7 \times 512$ and are input to the VAB after being processed with 2×2 average pooling. The visual feature maps are encoded through four convolutional layers before being processed by a GAP. The likelihood $P_v(\mathbf{y}_k)$ of the current subword \mathbf{y}_k is then predicted. In parallel, a visual attention map is created by an additional convolution and sigmoid normalization of the third convolutional layer of the visual attention branch. This attention map focuses selectively on certain parts of an image related to the predicted sequence. The VAB outputs visual features \mathbf{v}_k that are weighted by the attention mask. A cross-entropy loss L_v is minimized by the VAB.

4) *Linguistic Attention Branch*: Fig. 5 shows the network structure of the LAB. The LAB takes, as input, the last N outputs of the LSTM as a linguistic context \mathbf{c}_k . The parameter N is fixed and is described in the experimental section. We define a linguistic context \mathbf{c}_k as follows:

$$\mathbf{c}_k = \{\mathbf{h}_{k-N-1}, \mathbf{h}_{k-N}, \dots, \mathbf{h}_{k-1}\}, \quad (3)$$

where \mathbf{h}_k is the LSTM output at step k . The linguistic context \mathbf{c}_k has dimension $N \times d$, where d is the dimension of the LSTM hidden state. Thus, the LAB aims to produce an attention map of dimension $1 \times N$ that weights each component of \mathbf{c}_k . To this end, \mathbf{c}_k is processed by three one-dimensional convolutional layers enhanced by batch normalization (BN) and ReLU. Thereafter, the subword \mathbf{y}_k is predicted from the following fully connected (FC) and softmax layer. In parallel the attention map \mathbf{a}_k is obtained by connecting the second convolutional layer to a convolutional layer with size 1×1 , followed by BN and Sigmoid functions. The attention map \mathbf{a}_k has dimension $1 \times N$ and can be expressed as:

$$\mathbf{a}_k = \{w_{k-N-1}, w_{k-N}, \dots, w_{k-1}\}, \quad (4)$$

where each parameter w_k is the weight of the corresponding hidden state \mathbf{h}_k . The output \mathbf{l}_k of the LAB is the weighted linguistic context given by:

$$\mathbf{l}_k = \{\mathbf{o}_{k-N-1}, \mathbf{o}_{k-N}, \dots, \mathbf{o}_{k-1}\}, \quad (5)$$

where \mathbf{o}_k can be expressed as:

$$\mathbf{o}_k = (1 + w_k)\mathbf{h}_k. \quad (6)$$

Similarly to the VAB, a cross-entropy loss L_l is minimized based on the likelihood $P_l(\mathbf{y}_k)$ of the predicted subword.

TABLE II
PARAMETER SETTINGS OF THE ABEN

Opt. method	Adam (Learning rate = 1.0×10^{-4} , $\beta_1 = 0.7$, $\beta_2 = 0.99999$)
Backbone CNN	ResNet-50
LSTM	3 layers, 768-dimensional cell
N	10
Generation Branch	FC: 768, 768
Batch size	32

5) *Generation branch*: Fig. 5 shows the structure of the generation branch, which builds the sequence of subwords that compose the fetching instruction. The inputs \mathbf{h}_k and \mathbf{l}_k are concatenated and processed by FC layers, from which the likelihood of the next subword $p_g(\mathbf{y}_k)$ is predicted. A cross-entropy loss L_g is minimized in the generation branch.

6) *Loss functions*: The global loss function L_{ABEN} of the network is given by:

$$L_{\text{ABEN}} = L_v + L_l + L_g, \quad (7)$$

where L_v , L_l , and L_g denote cross-entropy losses based on the VAB, LAB, and generation branch, respectively. Using L as a generic notation for L_v , L_l and L_g , the cross-entropy loss is expressed as follows:

$$L = - \sum_n \sum_m y_{nm}^* \log p(y_{nm}), \quad (8)$$

where y_{nm}^* denotes the label given to the m -th dimension of the n -th sample, and y_{nm} denotes its prediction. It should be emphasized that the same labels are used for $P_l(\mathbf{y}_k)$ in the LAB and for $P_g(\mathbf{y}_k)$ in the generation branch.

V. EXPERIMENTS

A. Dataset

The dataset was collected in simulated home environments as described in Section III. The robot patrolled the environment automatically and collected images of designated areas. The environment was procedurally generated with everyday objects and furniture. Each image collected was labeled automatically with the bounding boxes of the sources and targets extracted from the simulator. These images were then annotated by three different labelers due to the limited size of the dataset. Each of them was instructed to provide a fetching instruction for each target. It should be noted that each image may contain multiple candidate targets and sources. Overall, we collected a dataset of 2,865 image–sentence pairs from 308 unique images and 1,099 unique target candidates. The dataset is available at this URL².

Standard linguistic pre-processing was performed on the instructions. The characters were converted to lowercase, and sentence periods were removed. Stopword replacement and stemming were not performed because subword tokenization and generation were able to handle word variations.

The dataset was split into 80%, 10%, and 10% parts for the training, validation, and test sets, respectively. After removing invalid samples, we could obtain 2,295 training samples, 264 validation samples, and 306 test samples. Because there was no overlap between the training, validation, and test sets, the test set was considered to be unseen.

²<https://keio.box.com/s/cbup2rttf1gkf487sgad34fq01wa5r0>

TABLE III

QUANTITATIVE RESULTS OF FIG. THE RESULTS ARE THE AVERAGE OVER 5 TRIALS. FOR READABILITY, THE METRICS ARE MULTIPLIED WITH 100. "ABEN w/o BBSE" USES SIMPLE SKIP-GRAM INSTEAD OF BERT-BASED SUBWORD EMBEDDINGS. "ABEN (SS)" AND "ABEN (TF)" USE SCHEDULED SAMPLING AND TEACHER FORCING, RESPECTIVELY.

Method	Evaluation metric						
	BLEU-1	BLEU-2	BLEU-3	BLEU-4	ROUGE	METEOR	CIDEr
VSE [40]	43.9±1.5	29.7±1.3	19.0±1.7	11.5±1.8	35.7±1.2	14.3±0.7	21.3±4.2
Multi-ABN [9]	49.1±0.9	35.4±1.8	24.0±2.3	16.0±2.4	37.8±1.4	19.9±1.1	27.5±5.0
ABEN w/o BBSE	58.2±1.4	38.5±1.7	23.7±2.6	13.9±2.1	42.8±1.2	17.9±0.6	38.0±2.7
ABEN (SS)	61.8±2.8	46.6±3.3	34.0±3.1	24.7±2.9	47.7±1.4	22.0±1.5	54.5±6.8
ABEN (TF)	60.2±1.9	45.1±1.8	33.5±2.2	24.9±2.4	48.2±1.3	22.8±1.7	57.6±4.6

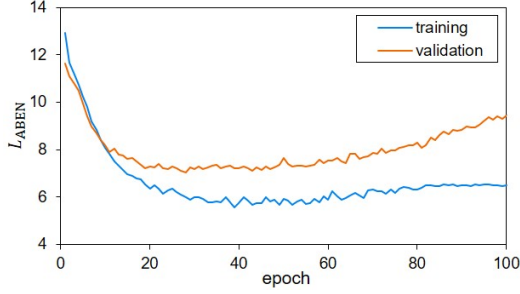


Fig. 6. Training and validation losses for 100 epochs.

B. Parameter settings

The parameter settings of the ABEN are shown in Table II. We used Adam as the optimizer, with a learning rate of 1×10^{-4} . The dimension of the BERT-based embedding vector was 768. We used a three-layer LSTM in the decoder (see Fig. 3) where each cell had a dimension of 768. The parameter N which characterizes the size of the linguistic context \mathbf{c}_k was set to 10. More specifically, we considered the 10 previous output of the LSTM to infer the linguistic attention map \mathbf{a}_k . As a result, in the early steps ($k < 10$), the linguistic context \mathbf{c}_k was initialized with the output of the encoder \mathbf{x}_f for all hidden states that were not available. The generation branch had two FC layers, each of which had 768 nodes. Each dimension of \mathbf{x}_{rel} was standardized so that its mean and standard deviation became 0 and 1, respectively. The visual inputs \mathbf{x}_{targ} , \mathbf{x}_{src} and \mathbf{x}_v were resized as $224 \times 224 \times 3$ images before being input to ResNet-50.

We trained the ABEN with the aforementioned dataset. The training was conducted on a machine equipped with a Tesla V100 with 32 GB of GPU memory, 768 GB RAM and an Intel Xeon 2.10 GHz processor. The ABEN was trained for 100 epochs, which was sufficient for loss convergence in pilot experiments.

C. Quantitative results

Table III shows the quantitative results, where standard metrics scores, used in image captioning, are reported. We conducted five experimental runs for each method. The table shows the mean and standard deviation for each metric. The BLEU-N column shows the standard BLEU score based on N-grams, where $N \in \{1, 2, 3, 4\}$. The CIDEr score was averaged over the same N-grams as BLEU. Additionally, we used ROUGE-L [41] which is based on the longest common subsequence. ROUGE-L did not use N-grams. METEOR

was computed from unigrams only, but endows a paraphrase dictionary.

We compared the ABEN with two baseline methods: visual semantic embedding (VSE) [40] and Multi-ABN [9]. Based on the standard method for model selection in deep neural network (DNN), we selected the best model as the one that maximized the METEOR score of the validation set. This is because METEOR has a paraphrase dictionary, which is more suitable for handling natural language. Fig. 6 depicts the training and validation loss of a typical run.

The results show that the ABEN outperformed the Multi-ABN and VSE for all four metrics. In particular, the CIDEr score was drastically improved by 30.1 points relative to the Multi-ABN and by 36.3 points relative to, VSE on average. Additionally, the t-test showed that the difference from VSE was statistically significant for all the metrics ($p < 0.001$). The difference from the Multi-ABN was also statistically significant ($p < 0.05$). Therefore, the ABEN significantly outperformed these baseline methods for the FIG task.

We conducted an ablation study on word embedding. In the ablation, we compared simple skip-gram and BERT-based subword embedding. In the table, "ABEN w/o BBSE" uses simple skip-gram instead of BERT-based subword embedding. The BERT-based subword embedding has better performance than skip-gram. The t-test showed that the results were statistically significant ($p < 0.001$) for all the metrics except BLEU-1.

Additionally, we tested two approaches in training: teacher forcing (TF) and scheduled sampling (SS) [42]. We adopted the standard SS setup with a linear decay $\epsilon = (\text{max_epoch} - \text{epoch}) / \text{max_epoch}$, where ϵ is the probability of using the label for training.

In the table, the results of these approaches are shown as "ABEN (SS)" and "ABEN (TF)". The t-test showed that the p-values for all the metrics are $p > 0.1$. Therefore, there was no statistically significant difference between teacher forcing and scheduled sampling. This indicates that TF did not cause the performance to deteriorate significantly in this task.

D. Qualitative results

For more insight into the performance of the ABEN, we analyzed the generated sentences qualitatively, as shown in Fig. 7. The top panels of the figures show the input scenes, and the middle and bottom tables show the reference sentences (Ref1, Ref2, and Ref3) and the sentences generated by the ABEN and the Multi-ABN. In the input image, the targets and sources are highlighted by green and blue boxes, respectively.


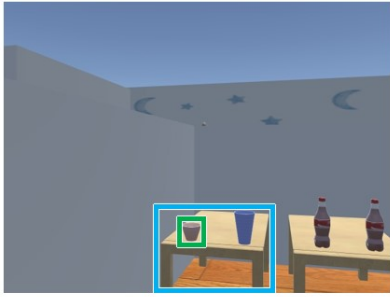
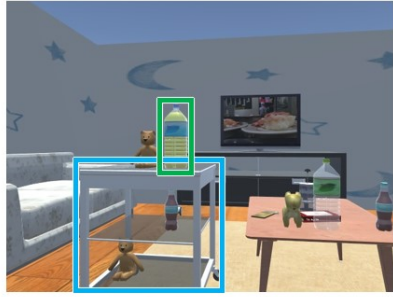
		
Ref1: “give me the duck on the shelf”	Ref1: “pick up the pink cup please”	Ref1: “get me the big bottle on the metal wagon”
Ref2: “let me have a duck item on the middle shelf of cabinet”	Ref2: “bring me a coffee cup on the left-sided small table”	Ref2: “bring me a green tea plastic bottle on the top of wagon, please”
Ref3: “please bring a red beak yellow duck toy”	Ref3: “could you grab a pink mug that is placed on the far left”	Ref3: “grasp the bottle containing the contents in the stage on the white wagon”
Multi-ABN: “i to toys shelf and take the yellow”	Multi-ABN: “please catch a the cup on the left-cup small table”	Multi-ABN: “grab the bottle on the i”
ABEN(Ours): “go to the shelf and take the yellow toy”	ABEN(Ours): “fetch a pink cup on the left-hand table”	ABEN(Ours): “can you bring me a green tea plastic bottle on the top”

Fig. 7. Three typical samples of qualitative results. Top figures show input images with bounding boxes of targets and sources. Middle row tables show three reference sentences annotated by labelers. Bottom tables show sentences generated by the Multi-ABN and ABEN.

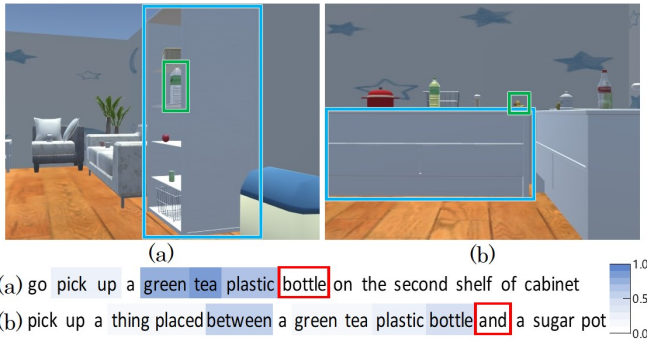


Fig. 8. Two typical qualitative results for linguistic attention. The blue and green boxes represent the source and the target. The predicted words are shown in the red boxes, and the attention values are overlaid in blue. In the left case (a), “green”, “tea” and “plastic” were attended for predicting “bottle”. In the right case (b), “between” and “bottle” were attended for predicting “and”.

The left-hand sample in Fig 7 shows a sentence that was generated successfully, semantically and syntactically, by the ABEN, in contrast to that generated by the Multi-ABN. Indeed, the target can be uniquely identified from the sentence “go to the shelf and take the yellow toy”, which is a valid fetching instruction. Conversely, the sentence generated by the Multi-ABN, refers somehow to the source (‘shelf’) and the target (‘toys’, ‘yellow’) but is incorrect syntactically. Such a sentence would require additional review by a human expert in the targeted use case of generating datasets of image–sentence pairs. Furthermore, the sentence generated by the ABEN is totally different from the reference sentences; this suggests that the many-to-many mapping between language and the environment was captured successfully.

Similarly, the second sample in the middle column illustrates the successful generation of a referring expression, which was used to disambiguate the source. The ABEN generated the sentence “fetch a pink cup on the left-hand

table” to refer to the target. In particular, the source was described correctly (“on the left-hand table”) even though the scene contains another similar source. Conversely, the baseline method generated “please catch a the cup on the left-cup small table”, which included erroneous syntax about the source (“left-cup” instead of “left-hand”). Additionally, over-generation appeared, as the phrase “a the”.

The right-hand sample illustrates ambiguity about the target. In this scene, there were three bottles. Therefore, the sentence should include referring expressions to determine the target uniquely. The sentence generated by the ABEN was able to disambiguate the target, which was referred to as “green tea plastic bottle”. However, the source description was incomplete. Indeed, a more exhaustive source description such as “on the top of the white wagon”, could be expected. Nonetheless, this fetching instruction remains understandable to human experts. Conversely, the baseline method generated a sentence that was syntactically incorrect but also ambiguous. The target was simply referred to as “bottle”, from which the target cannot be identified.

Additionally, we analyzed the relationship between the subwords in the sentence generation process through the linguistic attention maps shown in Fig. 8. The lower part illustrate the salient subwords that were used to predict the subword marked with a red frame. In Fig. 8(a), to predict “bottle”, the most salient subwords were “green”, “tea”, and “plastic”. In Fig. 8(b), to predict the subword “and” in the sentence, the most salient words were “bottle” but also “between”. These results indicate that the ABEN handles subword relationships in a representation that is understandable to humans.

Overall, these results emphasize that the ABEN generates more natural sentences than the baseline method, through the contribution of our proposed LAB architecture and the subword generation strategy.

VI. CONCLUSIONS

Most data-driven approaches for multimodal language understanding require large-scale datasets. However, building such a dataset is time-consuming and costly. In this study, we proposed the ABEN, which generates fetching instructions from images. Target use cases include generating and augmenting datasets of image-sentence pairs.

The following contributions of this study can be emphasized:

- The ABEN extends the Multi- ABN by introducing a linguistic branch and a generation branch, to model the relationship between subwords.
- The ABEN combines attention branches and BERT-based subword embedding for sentence generation.

Future studies will investigate the application of the ABEN to real-world settings.

ACKNOWLEDGMENT

This work was partially supported by JSPS KAKENHI Grant Number 20H04269, JST CREST, SCOPE, and NEDO.

REFERENCES

- [1] L. Piyathilaka and S. Kodagoda, "Human Activity Recognition for Domestic Robots," in *Field and Service Robotics*, 2015, pp. 395–408.
- [2] C.-A. Smarr *et al.*, "Domestic Robots for Older Adults: Attitudes, Preferences, and Potential," *International Journal of Social Robotics*, vol. 6, no. 2, pp. 229–247, 2014.
- [3] L. Iocchi, D. Holz, J. Ruiz-del Solar, K. Sugiura, and T. Van Der Zant, "RoboCup@ Home: Analysis and Results of Evolving Competitions for Domestic and Service Robots," *Artificial Intelligence*, vol. 229, pp. 258–281, 2015.
- [4] P. Anderson, Q. Wu, D. Teney, J. Bruce, M. Johnson, N. Sünderhauf, I. Reid, S. Gould, and A. van den Hengel, "Vision-and-Language Navigation: Interpreting Visually-Grounded Navigation Instructions in Real Environments," in *IEEE CVPR*, 2018, pp. 3674–3683.
- [5] X. Wang, Q. Huang, A. Celikyilmaz, J. Gao, D. Shen, Y.-F. Wang, W. Y. Wang, and L. Zhang, "Reinforced Cross-Modal Matching and Self-Supervised Imitation Learning for Vision-Language Navigation," in *IEEE CVPR*, 2019, pp. 6629–6638.
- [6] A. Magassouba, K. Sugiura, and H. Kawai, "A Multimodal Classifier Generative Adversarial Network for Carry and Place Tasks from Ambiguous Language Instructions," *IEEE RAL*, vol. 3, no. 4, pp. 3113–3120, 2018.
- [7] A. Magassouba, K. Sugiura, A. T. Quoc, and H. Kawai, "Understanding Natural Language Instructions for Fetching Daily Objects Using GAN-Based Multimodal Target-Source Classification," *IEEE RAL*, vol. 4, no. 4, pp. 3884–3891, 2019.
- [8] H. Fukui, T. Hirakawa, T. Yamashita, and H. Fujiyoshi, "Attention Branch Network: Learning of Attention Mechanism for Visual Explanation," in *IEEE CVPR*, 2019, pp. 10705–10714.
- [9] A. Magassouba, K. Sugiura, and H. Kawai, "Multimodal Attention Branch Network for Perspective-Free Sentence Generation," *CoRL*, 2019.
- [10] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding," in *NAACL-HLT*, 2019, pp. 4171–4186.
- [11] K. Sugiura, N. Iwahashi, H. Kashioka, and S. Nakamura, "Active Learning of Confidence Measure Function in Robot Language Acquisition Framework," in *IEEE/RSJ IROS*, 2010, pp. 1774–1779.
- [12] J. Hatori *et al.*, "Interactively Picking Real-World Objects with Unconstrained Spoken Language Instructions," in *IEEE ICRA*, 2018, pp. 3774–3781.
- [13] M. Shridhar and D. Hsu, "Interactive Visual Grounding of Referring Expressions for Human-Robot Interaction," in *RSS*, 2018.
- [14] S. Kazemzadeh *et al.*, "ReferItGame: Referring to Objects in Photographs of Natural Scenes," in *EMNLP*, 2014, pp. 787–798.
- [15] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, "Microsoft COCO: Common Objects in Context," in *ECCV*, 2014, pp. 740–755.
- [16] L. Kunze, T. Williams, N. Hawes, and M. Scheutz, "Spatial Referring Expression Generation for Hri: Algorithms and Evaluation Framework," in *AAAI Fall Symposium Series*, 2017.
- [17] F. I. Dogan, S. Kalkan, and I. Leite, "Learning to Generate Unambiguous Spatial Referring Expressions for Real-World Environments," in *IEEE/RSJ IROS*, 2019, pp. 4992–4999.
- [18] D.-K. Nguyen and T. Okatani, "Improved Fusion of Visual and Language Representations by Dense Symmetric Co-attention for Visual Question Answering," in *IEEE CVPR*, 2018, pp. 6087–6096.
- [19] J. Lei, L. Yu, T. L. Berg, and M. Bansal, "TVQA+: Spatio-Temporal Grounding for Video Question Answering," in *arXiv preprint arXiv:1904.11574*, 2019.
- [20] H. Akbari, S. Karaman, S. Bhargava, B. Chen, C. Vondrick, and S.-F. Chang, "Multi-level Multimodal Common Semantic Space for Image-phrase Grounding," in *IEEE CVPR*, 2019, pp. 12476–12486.
- [21] L. Yu, Z. Lin, X. Shen, J. Yang, X. Lu, M. Bansal, and T. L. Berg, "MAAttNet: Modular Attention Network for Referring Expression Comprehension," in *IEEE CVPR*, 2018, pp. 1307–1315.
- [22] C. Hori, T. Hori, T. K. Marks, and J. R. Hershey, "Early and Late Integration of Audio Features for Automatic Video Description," in *IEEE ASRU*, 2017, pp. 430–436.
- [23] B. Zhou, A. Khosla, A. Lapedriza, A. Oliva, and A. Torralba, "Learning Deep Features for Discriminative Localization," in *IEEE CVPR*, 2016, pp. 2921–2929.
- [24] M. Schuster and K. Nakajima, "Japanese and Korean Voice Search," in *IEEE ICASSP*, 2012, pp. 5149–5152.
- [25] R. Sennrich, B. Haddow, and A. Birch, "Neural Machine Translation of Rare Words with Subword Units," in *ACL*, 2016, pp. 1715–1725.
- [26] T. Kudo, "Subword Regularization: Improving Neural Network Translation Models with Multiple Subword Candidates," in *ACL*, 2018, pp. 66–75.
- [27] Z. Lan, M. Chen, S. Goodman, K. Gimpel, P. Sharma, and R. Soricut, "ALBERT: A Lite BERT for Self-supervised Learning of Language Representations," in *ICLR*, 2020.
- [28] Z. Yang, Z. Dai, Y. Yang, J. Carbonell, R. R. Salakhutdinov, and Q. V. Le, "Xlnet: Generalized Autoregressive Pretraining for Language Understanding," in *NIPS*, 2019, pp. 5754–5764.
- [29] D. Matthews, S. Kriegman, C. Cappelle, and J. Bongard, "Word2vec to Behavior: Morphology Facilitates the Grounding of Language in Machines," in *IEEE/RSJ IROS*, 2019, pp. 4153–4160.
- [30] V. Cohen *et al.*, "Grounding Language Attributes to Objects using Bayesian Eigenobjects," in *IEEE/RSJ IROS*, 2019, pp. 1187–1194.
- [31] K. Papineni, S. Roukos, T. Ward, and W.-J. Zhu, "BLEU: A Method for Automatic Evaluation of Machine Translation," in *ACL*. Association for Computational Linguistics, 2002, pp. 311–318.
- [32] C.-Y. Lin, "ROUGE: A Package for Automatic Evaluation of Summaries," in *Text Summarization Branches Out*. Association for Computational Linguistics, 2004, pp. 74–81.
- [33] S. Banerjee and A. Lavie, "METEOR: An Automatic Metric for MT Evaluation with Improved Correlation with Human Judgments," in *the ACL Workshop on IEEM for MTS*, 2005, pp. 65–72.
- [34] R. Vedantam *et al.*, "CIDER: Consensus-based Image Description Evaluation," in *IEEE CVPR*, 2015, pp. 4566–4575.
- [35] T. Yamamoto *et al.*, "Development of Human Support Robot as the research platform of a domestic mobile manipulator," *ROBOMECH journal*, vol. 6, no. 1, p. 4, 2019.
- [36] T. Inamura, J. T. C. Tan, K. Sugiura, T. Nagai, and H. Okada, "Development of Robocup@Home Simulation towards Long-term Large Scale HRI," in *Robot Soccer World Cup*, 2013, pp. 672–680.
- [37] Y. Mizuchi and T. Inamura, "Cloud-based Multimodal Human-robot Interaction Simulator Utilizing ROS and Unity Frameworks," in *IEEE/SICE SII*, 2017, pp. 948–955.
- [38] Y. Wu *et al.*, "Google's Neural Machine Translation System: Bridging the Gap Between Human and Machine Translation," *arXiv preprint arXiv:1609.08144*, 2016.
- [39] K. He, X. Zhang, S. Ren, and J. Sun, "Deep Residual Learning for Image Recognition," in *IEEE CVPR*, 2016, pp. 770–778.
- [40] O. Vinyals *et al.*, "Show and Tell: A Neural Image Caption Generator," in *IEEE CVPR*, 2015, pp. 3156–3164.
- [41] C.-Y. Lin and F. J. Och, "Automatic Evaluation of Machine Translation Quality using Longest Common Subsequence and Skip-bigram Statistics," in *ACL*, 2004, pp. 605–612.
- [42] S. Bengio, O. Vinyals, N. Jaitly, and N. Shazeer, "Scheduled Sampling for Sequence Prediction with Recurrent Neural Networks," in *NIPS*, 2015, pp. 1171–1179.