# INDOOR VISUAL UNDERSTANDING THROUGH IMAGE CAPTIONING

Dhomas Hatta Fudholi[*], Royan Abida N. Nayoan

Department of Informatics, Universitas Islam Indonesia, Yogyakarta, Indonesia

## Graphical abstract



## Abstract

Transformers have been widely used in image captioning tasks on English language datasets such as MSCOCO and Flickr. However, research related to image captioning in Indonesian is still rare and relies on machine translation to obtain the Indonesian dataset. In this study, the Transformer model is used to generate caption using the modified MSCOCO datasets to gain visual understanding in an indoor environment. We modified the MSCOCO dataset by creating new Indonesian text description based on the MSCOCO images. A few simple rules are made to create the Indonesian dataset by including the object's location, colour, and its characteristics. Experiments were carried out using several CNN pre-trained models to extract the image features before feeding them to the Transformer model. We also performed hyper-parameter settings on the models by assigning different values for batch size, dropouts, and attention heads to get the best model. BLEU-n, METEOR, CIDEr, and ROUGE-L are used to evaluate the model. From this study, by utilizing the EfficientNetB0 with a batch size of 128, dropouts of 0.2, and attention heads of 4, the model can get the best score in four different evaluation matrices. The EfficientNetB0 model reached the highest score on BLEU-4 with a score of 0.344, ROUGE-L of 0.535, METEOR of 0.264, and CIDEr of 0.492.

*Keywords*: image captioning, indoor, visual understanding, EfficientNet, Transformer

## 1.0 INTRODUCTION

Image captioning is an important and complex task that connects computer vision and language model. Image captioning aims to quickly understand the given image and generate textual description of the image automatically. Image captioning is not an easy task since it needs not only to identify objects in the image but to also perform semantic analysis to understand the information, relationship, and scene information of the image [1]. This technique has been widely used for different purposes in the real world, such as helping the visually impaired to comprehensively understand and sense their surroundings better [2], [3]. Most image captioning models [4]–[6] are based on the encoder-decoder framework. The encoder is responsible for extracting the image feature from the given image, while the decoder acts as the language model and is responsible to generate corresponding captions automatically. To perform these two tasks, image captioning in this paper uses a combination of deep learning models CNN as the encoder and Transformer as the decoder.

Deep learning (DL) is applied in solving the problem of image captioning. Deep learning is a subset of machine learning that has been improved both in terms of algorithms and the efficiency of its pre-processing techniques [7]. In recent decades, DL has produced several state-of-the-art models compared to other traditional algorithms that can solve complicated problems. Some of the most popular models are Recurrent Neural Networks (RNN) and Convolutional Neural Network (CNN), where RNN usually works in solving sequential problems and CNN is suitable for computer vision tasks [8], [9].

In the last few decades, image captioning has been using the encoder-decoder approach. CNN is typically used as the

encoder then followed by an RNN based model as the decoder. While it is the dominating architectures in image captioning, RNN based models and its variants are known to have difficulties in memorizing the inputs [10]. To solve this issues, previous works in image captioning changed the decoder to Transformer [11], [12]. Not only transformer solves RNN's limitations, but it also has attention mechanism that encourages the model to explore the relationships between the detected entities and requires less time in training the model [10]. Since Transformer was originally designed for the machine translation tasks, Transformer requires its input and output to be both sequences of text. But in the image captioning task, Transformer needs an image as an input [12].

Most of image captioning tasks or applications are based on English captions. For Indonesian image captioning, previous works translated large-scale datasets such as MS COCO and Flickr text description by using google translate or professional English-Indonesian machine translation to get the Indonesian dataset [13]–[15]. In this study we created an Indonesian image captioning to assist visually impaired people in an indoor environment to help them get a better sense of their surroundings. To reach that goal, we changed the text description of the MS COCO dataset. We created an Indonesian text description based on the MS COCO images by including the object's location, shape, color, and its characteristics.

We propose an Indonesian image captioning model using pretrained CNN and vanilla Transformer model. Here, an Indonesian dataset is created to help visually impaired people in getting a better indoor visual understanding. Pre-trained CNN model is used in the encoding stage to extract image features. While Transformer is used in the decoding stage to generate the caption. The contributions of this study can be summarized as follows: (1) creates an Indonesian dataset for indoor objects that includes the object's location, shape, characteristics, and its color; (2) compares a few pre-trained CNN to extract image features (DenseNet169, EfficientNetB0, IncepResNet V2, Inception V3, and Xception); and (3) builds a novel image captioning model for indoor visual understanding using CNN-Transformer architecture.

This paper is organized as follows: the related work is discussed in Section 2. Methodology is proposed in Section 3. Experiments and evaluation results are detailed in Section 4. Finally, the concluding remarks and future work are summarized in Section 5.

## 2.0  LITERATURE REVIEW

Image captioning is an artificial intelligence research area that combines image understanding and language description. Image captioning needs to recognize objects within the image and generates well-formed image description. Early image captioning approach uses retrieval-based where image captioning is treated as a retrieval task [16]. Then in the same year Vinyals introduced an image captioning using a deep neural network approach using encoder-decoder framework in his paper titled "Show and Tell: A Neural Image Caption Generator" [17].

Deep learning is known for its capability in handling large dataset such as images or videos. In deep neural network-based approach for image captioning, encoder-decoder framework is used. This technique allows the model to learn

and extract the features automatically from the training dataset. CNN is usually used for feature learning before finally feeding the extracted feature to the RNN to generate the image description [18]. Now, most models use deep learning-based approach for various language image captioning since they can generate semantically accurate captions [1], [12], [19].

Before the Transformer became popular, RNN was a well-known approach in sequence problems. RNN has a problem in capturing long-term dependencies since it decreases the long-term information at every time step [20]. To overcome this limitation, Transformer replaces the recurrent neural network as the decoder in encoder-decoder architecture using attention mechanism [21]. Transformer uses multi-headed self-attention that allows parallel training which makes it more efficient than RNN based architectures. Previous studies show that applying transformer-based models gave promising results and high-quality image captioning [10], [12], [22]. In recent years, several upgrades have also been made to the transformer model after the replacement of RNN to the transformer model in encoder-decoder framework. Previous research uses Dual Global Enhanced Transformer (DGET) using Global Enchanced Encoder (GEE) and Global Enchanced Decoder (GED) to incorporate global information in the encoding and decoding stages [23]. Another research uses LATGeO to learn the geometry relationship among objects [24], using transformer encoder to avoid convolution operation that has a limitation in global context modelling [25], and mask-cross-entropy strategy to improve the diversity of the generated captions and explore uncommon word relations [26].

Most image captioning task usually generated English captions, since publicly available dataset such as MSCOCO [27] and Flickr [28] is only available in English. Lately image captioning research that generates different language has started to emerge such as Burmese, Chinese, and Indonesian image captioning. This research relied on machine translation to translate the English captions into their own language [19], [29]. Previous Indonesian image captionings also use machine translation such as google translate to get the Indonesian caption [15]. But since google translate still faces difficulty in getting the correct translation based on the context, the author had to correct a lot of incorrectly translated captions [13], [14]. Until recently, Indonesian image captioning had only performed image captioning task using the encoder-decoder approach by applying CNN and RNN-based models.

We propose an encoder-decoder framework on Indonesian dataset to create visual understanding. We use the publicly available dataset MSCOCO to create the indoor object dataset (bed, chair, couch, hanging lamp, oven, sink, tv, washing machine, windows, and potted plant) and change the text description for each image by including object's location, characteristics, color, its shape. We experimented using different pretrained CNN models to extract image features and performed a hyper-parameter tuning to obtain the best result. We also replaced the RNN-based model using vanilla Transformer to act as the decoder in the encoder-decoder framework.

To note, this paper is a continuation of our previous work [30]. The previous work focused on developing an image captioning model on an Indonesian dataset using a combination of pre-trained CNN as the encoder and vanilla

Transformer as the decoder then compare it to a merge encoder-decoder using CNN and RNN based model. The difference between our initial work and this work is we created 5 Indonesian instead of 3 captions. We also created a more rigid Indonesian dataset where we only mention the main objects in the images and include the name of the object, the colour, the shape, the characteristics, and nothing else. For further explanation regarding how we created the Indonesian caption see Section 3.1 Data Collection. Furthermore, we experimented using several different pre-trained CNN such as DenseNet169, IncepResNet V2, EfficientNetB0, InceptionV3, and Xception to extract the image features before feeding them to the vanilla Transformer.

## 3.0 METHODS

The aim of our research is building an image captioning model which can help in giving visual understanding, especially for indoor environment. Figure 1 shows the research methodology that is done. Data collection is the first step that collects related data for the dataset. The collected data will be preprocessed before going to feature extraction and modeling step. It is done to build clean dataset. Finally, evaluation is done to see the performance of the developed model. Each process is elaborated as follows.
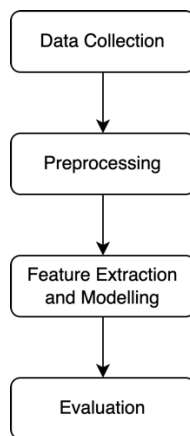
**Figure 1** Research methodology

### 3.1 Data Collection

Image captioning datasets were collected from the publicly available MSCOCO that provides more than 150k images [27]. MSCOCO allows user to retrieve images based on the selected object. In this research, 10 different indoor objects are selected from MSCOCO, namely, bed, chair, couch, hanging lamp, oven, sink, tv, washing machine, windows, and potted plant. Each object contains 60 images with 5 different image captions.

To achieve the goal in creating a visual understanding for an indoor environment, it is necessary to make changes to the dataset. MSCOCO already has its own text description for their images in English. In this paper, we created new text descriptions in Indonesian that includes object's location, characteristics, color, and its shape to help visually impaired people in recognizing indoor objects. A total of 600 images and 3000 text descriptions are randomly selected into training and

testing set with a ratio of 8:2 [31]. The number is quite low since Indonesian can be considered as low resource language.

We made a few rules in writing the images captions following the previous research for image captioning dataset [27] (1) Since the goal is to describe indoor space surroundings to achieve visual understanding, the information of their surrounding objects and location information of each object is needed whether the objects are located on the left side/right side of the room. (2) If the objects in the image are of a close distance, the location information of the object is not needed. (3) Describe only the main objects within image. (4) Describe the colour of the objects and their characteristics as it could be beneficial in distinguishing each object [32]. An example of the Indonesian caption along with the English translated caption from a sample image in Figure 2 can be seen Table 1. The current caption given in Table 1 follows the four caption rules mentioned before.



**Figure 2** An image dataset extracted from MSCOCO

**Table 1** A comparison of our current caption and previous caption from our initial work based on the image in Figure 2 [30]

| No. | Current Caption | Previous Caption |
|---|---|---|
| 1 | *"Di sisi kiri terdapat tempat tidur dengan selimut berwarna biru."* (On the left side there is a bed with a blue blanket.) | *"Di samping kiri terdapat tempat tidur dengan sprei berwarna biru."* (On the left there is a bed with blue sheets.) |
| 2 | *"Di sisi kiri terdapat meja rias dari kayu dengan cermin berbentuk oval."* (On the left is a wooden dresser with an oval mirror.) | *"Di depan terdapat jendela di antara rak buku dan meja rias."* (In front there is a window between the bookcase and dresser.) |
| 3 | *"Di depan terdapat rak buku tinggi dari kayu."* (In the front there is a tall wooden bookshelf.") | *"Di samping tempat tidur terdapat meja rias."* (There is a dresser beside the bed.) |
| 4 | *"Di depan terdapat jendela dari kaca tanpa tirai."* (In the front there is an uncovered glass window.) | - |
| 5 | *"Di bagian bawah terdapat pot tanaman berwarna putih."* (At the bottom there is a white potted plant.") | - |

## 3.2 Preprocessing

After the MSCOCO images are collected, it is necessary to do preprocessing on the image and text description before passing them as inputs to our model. For text description, we performed a few preprocessing steps, (1) Changing the sentence into lower case, (2) Removing punctuations such as !"#$%&()*+.,-/:;=?@[\]^_`{|}~ ', (3) Removing single characters, (4) Removing numeric values, and (5) Adding <start> and <end> tag to mark the beginning and end of sentences, (6) Adding <pad> tag for sentences that are less than 25 words and <unk> tag for unknown words, and, (7) Tokenization to build a vocabulary. For image preprocessing, the images are resized depending on the pre-trained CNN model's requirement.

## 3.3. Feature Extraction and Modelling

This study experimented using five different pre-trained CNN models to extract image features. These five different pre-trained CNN models are namely Dense169, EfficientNetB0, IncepResNet V2, InceptionV3, and Xception. As CNN is not for a classification task, the last softmax layer is removed. By performing feature extraction, it allows the model to extract the important information about the objects contained in the image and the relationship between the objects. The extracted features are then stored in .npy files.

The DenseNet architecture [33] is a variant of Convolutional Neural Networks (CNN) where each dense layer is connected to one another and created a shortcut. This makes the model easy to train and parameter efficient as DenseNet allows features from different layers to be reused, thereby increasing the variation in the input of the next layer and the performance of the model.

We also experimented using EfficientNetB0, a state-of-the-art model in feature extraction [34]. Previous model used EfficientNetB0 as feature extractor in image captioning task for local tourism-related images and the model can generate detailed and sensible caption [35].

InceptionV3 is another variant of Inception that was developed and is one of the state-of-the-art pre-trained models [36]. Xception (Extreme version of Inception) is also another variant of Inception that has better performance compared to InceptionV3 [37]. Xception was developed using the same parameter as InceptionV3 with the aim of parameter model efficiency. Previous research on Xception shows that the model can provide the best performance and able to generate image description for eye disease [38]. InceptionResNet, as the name implies, is a model that combines a deep CNN layer with an Inception structure for computation efficiency and Residual Network (ResNet) to get the benefits of residual's optimization [39]. In this paper, we adopted IncepResNetV2, a reduced IncepResNet as one of the feature extractors.

This study follows the previous image captioning research in setting the model hyper-parameters values without significant architectural model modification [21]. We assign different values for each hyperparameter. The range of the hyperparameter values used in Transformer models can be seen in Table 2. All models run in 40 epochs due to computational limitations. The model uses Adam optimizer with β1=0.9 and β2=0.98 and sparse categorical as a loss function. The learning rate used during the training of the Transformer model varies according to the formula used in the

original paper by increasing and decreasing the learning rate value. The latent dimension in the multi-head attention module is set to 512 and the inner dimension in the feed-forward network module is set to 2048.

**Table 2** The range of hyperparameter values used in Transformer models

| Hyperparameter | Range |
|---|---|
| Attention heads | 4-16 |
| Batch Size | 32-128 |
| Dropout | 0.1-0.2 |

## 3.4 Evaluation

Different evaluation metrics are used to measure the performance of image captioning model. The models are evaluated using the standard evaluation metrics such as BLEU-n, CIDEr, ROUGE-L, and METEOR. In evaluating each text, BLEU (Bilingual Evaluation Understudy) considers the n-grams and ignores syntactical correctness [40]. BLEU score values range from 0.0 to 1.0. The higher the values indicate the best score between the reference and the candidate. BLEU score can be calculated using the formula in Eq. (1).

$$BLEU = \min\left(1, \frac{output-length}{reference-length}\right)\left(\prod_{i=1}^{4} precision_i\right)^{\frac{1}{4}} \quad (1)$$

There are various ROUGE metrics that can be used to evaluate different types of texts. ROUGE-L calculates the Longest Commong Subsequence (LCS) between generated text and reference [41]. The formula used to calculate ROUGE-L can be seen in the Eq. (2), Eq. (3), and Eq. (4).

$$recall = \frac{LCS(gram_n)}{count(gram_n)} \quad (2)$$

$$precision = \frac{LCS(gram_n)}{count(gram_n)} \quad (3)$$

$$F1\ score = 2 \times \frac{precision*recall}{precison+recall} \quad (4)$$

METEOR (Metric for Evaluation for Translation with Explicit Ordering) can capture semantic correlation between candidates and references [42]. The METEOR matrix calculates the accuracy, recall, and f-score of each word, stemmed word, and synonyms. The METEOR formula can be written as in Eq. (5), Eq. (6), Eq. (7), Eq. (8), and Eq. (9). Where $P$ is precision, $R$ is recall, and $F_{mean}$ is used to calculate the F score. $m$ is the number of unigrams in the candidate that are also found in the references. $w_t$ is the number of unigrams in the candidate, $w_r$ is the number of unigrams in the reference. p is a penalty with c as the number of chunks and $u_m$ as the unigrams that have been mapped. Meanwhile, M is the formula to calculate the METEOR score.

CIDEr (Consensus-based Image Description Evaluation) calculates n-gram using term frequency-inverse document frequencies (TF-IDF) to calculate the cosine similarity between reference and candidate captions [43]. CIDEr turns each word

into their root or stemmed form and only considers words that are important and hold significant meanings.

$$P = \frac{m}{w_t} \tag{5}$$

$$R = \frac{m}{w_r} \tag{6}$$

$$F_{mean} = \frac{10PR}{R+9P} \tag{7}$$

$$p = 0.5\left(\frac{c}{u_m}\right)^3 \tag{8}$$

$$M = F_{mean}(1-p) \tag{9}$$

## 4.0 RESULT AND DISCUSSION

In this study, we set up an experiment using several variants of CNN as encoder. These CNN has been previously trained using ImageNet dataset to perform feature extraction [44]. The selected pre-trained CNN variants are namely DenseNet169, EfficientNetB0, IncepResNetV2, InceptionV3, and Xception. Furthermore, we assigned different values for batch size, dropout, and attention heads on Transformer model to obtain the best model. Table 3 shows the CNN variants and the assigned hyperparameter values for each model.

In training the model, a total of 600 images and 3000 captions were used and then divided into two datasets: training set and testing set. We split the training and the testing set by following a commonly used ratio of 8:2 [31]. All datasets are images taken indoors by selecting 10 different objects retrieved from MSCOCO (bed, chair, couch, hanging lamp, oven, sink, tv, washing machine, windows, and potted plant).

Table 4 shows the BLEU, METEOR, CIDER and ROUGE-L scores for each model. It can be seen in the table that Model 5 and Model 9 using the EfficientNetB0 as feature extractor have higher scores compared to other models. Model 5 obtained a BLEU-1 score of 0.711, BLEU-2 of 0.552, BLEU-3 of 0.434, BLEU-4 of 0.344, ROUGE-L of 0.535, METEOR of 0.264, and CIDEr of 0.492, while Model 6 reaches a BLEU-1 score of 0.712, BLEU-2

of 0.561, BLEU-3 of 0.437, BLEU-4 of 0.341, ROUGE-L of 0.522, METEOR of 0.258, and CIDEr of 0.464. Furthermore, we also performed hyper-parameter tuning on the Transformer model. By assigning different hyper-parameter values and feature extractor, our experiment shows that higher batch size values and feature extractor EfficientNetB0 play a role in increasing the evaluation metrics scores. Model 5 has a higher score on four different metrics, namely BLEU-4, ROUGE-L, METEOR, and CIDER, whereas Model 9 has higher scores on BLEU metric: BLEU-1, BLEU-2, and BLEU-3.

At this point, we only compared each feature extractor model with the highest score, Model 1, Model 3, Model 5, Model 9, Model 10, and Model 12. From Table 5, these models work well in generating captions that are appropriate to the given image. Models can recognize the computer device by generating terms such as "computer/computer", "menyala/turned on monitor", "laptop", and "papan ketik/keyboard". However, Model 1 generated two different location in one sentence "di atas meja/on top of the table" and "sisi kanan/on the right" to describe the location information of the computer device. Model 3 and Model 5 generated the correct object and its color "putih/white", while Model 9 generated a different color "hitam/black". Model 10 identified the object and the color correctly whereas Model 12 generated "papan ketik berwarna hitam/black keyboard" that still corresponds to the given image.

**Table 3** Hyperparameter settings

| Model # | Feature Extractor | Batch Size | Dropout | Number of Attention Heads |
|---------|-------------------|------------|---------|---------------------------|
| 1 | DenseNet169 | 128 | 0.2 | 4 |
| 2 | DenseNet169 | 32 | 0.2 | 8 |
| 3 | InceptionResNetV2 | 128 | 0.2 | 4 |
| 4 | InceptionResNetV2 | 64 | 0.1 | 16 |
| 5 | EfficientNetB0 | 128 | 0.2 | 4 |
| 6 | EfficientNetB0 | 32 | 0.2 | 4 |
| 7 | EfficientNetB0 | 64 | 0.1 | 8 |
| 8 | EfficientNetB0 | 32 | 0.1 | 16 |
| 9 | EfficientNetB0 | 128 | 0.2 | 16 |
| 10 | InceptionV3 | 128 | 0.2 | 4 |
| 11 | InceptionV3 | 32 | 0.1 | 8 |
| 12 | Xception | 128 | 0.2 | 4 |
| 13 | Xception | 64 | 0.1 | 8 |

**Table 4** BLEU, ROUGE-L, METEOR, and CIDEr evaluation scores

| No | Feature Extractor | BLEU-1 | BLEU-2 | BLEU-3 | BLEU-4 | ROUGE-L | METEOR | CIDER |
|----|-------------------|--------|--------|--------|--------|---------|--------|-------|
| 1 | DenseNet169 | 0.688 | 0.527 | 0.401 | 0.306 | 0.498 | 0.251 | 0.415 |
| 2 | DenseNet169 | 0.671 | 0.526 | 0.409 | 0.323 | 0.498 | 0.253 | 0.356 |
| 3 | IncepResNetV2 | 0.665 | 0.519 | 0.399 | 0.310 | 0.499 | 0.248 | 0.408 |
| 4 | IncepResNetV2 | 0.638 | 0.475 | 0.353 | 0.265 | 0.462 | 0.234 | 0.331 |
| 5 | EfficientNetB0 | 0.711 | 0.552 | 0.434 | 0.344 | 0.535 | 0.264 | 0.492 |
| 6 | EfficientNetB0 | 0.659 | 0.5 | 0.381 | 0.296 | 0.479 | 0.244 | 0.333 |
| 7 | EfficientNetB0 | 0.665 | 0.501 | 0.381 | 0.292 | 0.499 | 0.244 | 0.411 |
| 8 | EfficientNetB0 | 0.627 | 0.473 | 0.354 | 0.266 | 0.479 | 0.246 | 0.311 |
| 9 | EfficientNetB0 | 0.712 | 0.561 | 0.437 | 0.341 | 0.522 | 0.258 | 0.464 |
| 10 | InceptionV3 | 0.659 | 0.505 | 0.381 | 0.291 | 0.490 | 0.244 | 0.340 |
| 11 | InceptionV3 | 0.575 | 0.414 | 0.303 | 0.223 | 0.424 | 0.209 | 0.269 |
| 12 | Xception | 0.655 | 0.502 | 0.384 | 0.293 | 0.485 | 0.237 | 0.410 |
| 13 | Xception | 0.64 | 0.477 | 0.356 | 0.272 | 0.461 | 0.233 | 0.311 |

**Table 5** Prediction Results

| Image | Caption |
|---|---|
|  | **Model 1**: *Di atas meja kerja terdapat sebuah komputer berwarna putih dalam keadaan menyala dengan sebuah komputer yang diletakkan pada sisi kanan.* (On the worktable there is a turned-on white computer with a computer placed on the right side.) **Model 3**: *Di atas meja kerja berwarna putih terdapat sebuah komputer dengan layar komputer berwarna putih dalam keadaan menyala.* (On the white desk there is a computer with the computer screen on.) **Model 5**: *Terdapat sebuah laptop berwarna putih dalam keadaan menyala dengan sebuah laptop yang diletakkan di atas meja.* (There is a turned-on white laptop with the laptop placed on top of the table.) **Model 9**: *Terdapat sebuah laptop berwarna hitam dalam keadaan menyala dengan sebuah laptop berwarna hitam di atas meja kerja.* (There is a turned-on white laptop with the laptop placed on the worktable.) **Model 10**: *Di sisi kanan terdapat sebuah laptop berwarna putih dalam keadaan tertutup di atas meja kerja panjang berwarna hitam.* (On the right side there is a closed white laptop on a long black desk.) **Model 12**: *Terdapat sebuah papan ketik berwarna hitam yang diletakkan di atas meja.* (There is a black keyboard placed on the table.) |

Model 5 and Model 9 use EfficientNetB0 with different attention of 4 and 16. Both models have the highest scores compared to other models. Model 5 has the highest scores on four different evaluation metrics (BLEU-4, METEOR, CIDER, and ROUGE-L), while Model 9 has the highest BLEU scores (BLEU-1, BLEU-2, and BLEU-3). We compared the quality of both models by analyzing a few of their generated captions shown in Table 6. In image#1, Model 5 generated a well-formed caption for a close-up image. It can recognize the laptop color and the state of the computer, whereas Model 9 generated a different laptop color. In image#2, Model 5 works well in recognizing the black laptop with the screen on, on top of the worktable, while Model 9 generated an incorrect table characteristic. In image#3, although Model 5 failed to recognize the tv's characteristics and color, it can identify the characteristics of the desk "wooden desk" and locate the position of the desk and the television. Model 9 can recognize the sofa but failed to generate the correct color and characteristics.

**Table 6** Comparison between the machine generated caption using EfficientNetB0 with attention 4 and EfficientNetB0 with attention 16

| No | Image | Caption |
|---|---|---|
| 1 |  | **Model 5:** *Terdapat sebuah laptop berwarna putih dalam keadaan menyala di atas tempat tidur berwarna putih.* (There is a white laptop with the screen on, on a white bed.) **Model 9:** *Di sisi kiri terdapat sebuah laptop berwarna hitam yang diletakkan di atas tempat tidur berwarna putih.* (On the left side there is a black laptop placed on a white bed.) |
| 2 |  | **Model 5:** *Di atas meja kerja terdapat sebuah laptop berwarna hitam dalam keadaan menyala.* (On the desk there is a black laptop with the screen on.) **Model 9:** *Terdapat sebuah meja kayu berukuran besar dengan sebuah laptop yang diletakkan di sisi kiri meja kerja.* (There is a large wooden table with a laptop placed on the left side of the worktable.) |
| 3 |  | **Model 5:** *Di sisi kanan terdapat sebuah meja dari kayu dengan sebuah televisi tabung berwarna hitam.* (On the right side there is a wooden desk and a black tube television.) **Model 9:** *Di sisi kanan terdapat sebuah sofa panjang berwarna coklat tua dengan bantal.* (On the right side there is a long dark brown sofa with pillows.) |

We tested our Indonesian image captioning using EfficientNetB0 with a batch size of 128, a dropout of 0.2, and an attention head of 4 that has the highest scores on four different metrics on images that are taken in Indonesian indoor environment. These images are collected from Google. As can be seen in Table 7, the model can still generate decent Indonesian caption on images that are culturally different from MSCOCO image dataset. Model is still able to identify the silver sink in image#1, a table with a blue tablecloth in image#2, and a wooden table in image#3.

**Table 7** Caption generation results on indoor images in Indonesia

| No | Image | Caption |
|----|-------|---------|
| 1 |  | *Di sisi kanan terdapat sebuah wastafel tanpa bingkai berwarna silver.* (On the right side there is a silver frameless sink.) |
| 2 |  | *Di sisi kanan terdapat sebuah meja kayu dengan taplak berwarna biru dan beberapa barang di atasnya.* (On the right side there is a wooden table with a blue tablecloth and several items on it.) |
| 3 |  | *Di sisi kiri terdapat meja dari kayu berwarna cokelat dengan beberapa barang di atasnya.* (On the left side there is a brown wooden table with some items on it.) |

This study compares our EfficientNetB0-Transformer model to several studies on Indonesian image captioning that have previously been conducted. Previous Indonesian image captioning are mostly using the MSCOCO and Flickr datasets that have been translated into Indonesian using machine translation such as Google translate or professional English to Indonesian machine translation. As can be seen in Table 7, the model using EfficientNetB0 and Transformer with the attention heads set to 4, dropout of 0.1, and batch size of 128 and only using a small dataset, the model is still comparable to the previous Indonesian image captioning studies.

We also compare our model to our previous work [30]. Compared to our previous model (Table 8 model#4), in this paper we use a combination of EfficientNetB0 and vanilla Transformer instead of IncepResNet V2. By using 600 images and increasing the captions for each image to 5 captions, the model has a far higher BLEU scores for all n-gram.

## 5.0 CONCLUSION

Several conclusions in this study can be drawn as follows: 1) we proposed an Indonesian image captioning to get an indoor visual understanding, 2) performed an experiment using several pre-trained CNN variants, 3) created a few rules in creating the Indonesian dataset, 4) performed a hyper-parameter tuning to get the best result.

From our experiment, EfficientNetB0 using a dropout of 0.1, a batch size of 128, and attention heads of 4 has the highest result on four different metrics BLEU-4 of 0.344, ROUGE-L of 0.535, METEOR of 0.264, and CIDEr of 0.492. This model is also tested on images taken in Indonesian indoor environments and although the test images are culturally different from MSCOCO images dataset, model can still generate decent captions that correspond to the given images.

**Table 8** Image captioning evaluation comparison to previous Indonesian image captioning studies

| No | Dataset | Model | Total Images | Captions per image | BLEU score (n-gram) | | | |
|----|---------|-------|-------|------|------|------|------|------|
| | | | | | 1 | 2 | 3 | 4 |
| 1 | Flickr - FEEH - ID [14] | CNN-LSTM | 8099 | 5 | 50.0 | 31.4 | 23.9 | 13.1 |
| 2 | Flickr30k-ID [45] | CNN-GRU | 31783 | 5 | 36.7 | 17.8 | 6.7 | 2.0 |
| 3 | MSCOCO & Flickr30k [46] | ResNet101-LSTM with adaptive attention | 180k | 5 | 67.8 | 51.2 | 37.5 | 27.4 |
| 4 | Modified MSCOCO [30] | IncepResNet-Transformer | 771 | 3 | 52.8 | 35.4 | 22.8 | 14.6 |
| **5** | **Modified MSCOCO** | **EfficientNetB0-Transformer** | **600** | **5** | **71.1** | **55.2** | **43.4** | **34.4** |

## References

[1] H. Lu, R. Yang, Z. Deng, Y. Zhang, G. Gao, and R. Lan. 2021. "Chinese Image Captioning via Fuzzy Attention-based DenseNet-BiLSTM," *ACM Transactions on Multimedia Computing, Communications, and Applications*, 17(1s): 1–18. doi: 10.1145/3422668.

[2] Y. Ming, N. Hu, C. Fan, F. Feng, J. Zhou, and H. Yu. 2022. "Visuals to Text: A Comprehensive Review on Automatic Image Captioning," *IEEE/CAA Journal of Automatica Sinica*, 9(8): 1339–1365. doi: 10.1109/jas.2022.105734.

[3] G. Xu, S. Niu, M. Tan, Y. Luo, Q. Du, and Q. Wu. 2021. "Towards Accurate Text-based Image Captioning with Content Diversity Exploration,". doi: 10.48550/ARXIV.2105.03236.

[4] S. Herdade, A. Kappeler, K. Boakye, and J. Soares. 2019. "Image captioning: Transforming objects into words," *Advances in neural information processing systems*, 32: 1–11.

[5] C. Yan, B. Gong, Y. Wei, and Y. Gao. 2021. "Deep Multi-View Enhancement Hashing for Image Retrieval," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 43(4): 1445–1451. doi: 10.1109/tpami.2020.2975798.

[6] X. Yang, H. Zhang and J. Cai. 2019. "Learning to Collocate Neural Modules for Image Captioning," *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*: 4249-4259. doi: 10.1109/ICCV.2019.00435.

[7] C. Janiesch, P. Zschech, and K. Heinrich. 2021. "Machine learning and deep learning," *Electronic Markets*, 31(3):685–695. doi: 10.1007/s12525-021-00475-2.

[8] Y. LeCun, Y. Bengio, and G. Hinton. 2015. "Deep learning," *Nature*, 521(7553): 436–444. doi: 10.1038/nature14539.

[9] A. Sherstinsky. 2020. "Fundamentals of Recurrent Neural Network (RNN) and Long Short-Term Memory (LSTM) network," *Physica D:*

*Nonlinear Phenomena*, 404: 132306. doi: 10.1016/j.physd.2019.132306.

[10] G. Li, L. Zhu, P. Liu, and Y. Yang. 2019. "Entangled Transformer for Image Captioning," *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*. doi: 10.1109/iccv.2019.00902.

[11] P. Zeng, H. Zhang, J. Song, and L. Gao. 2022. "S2 Transformer for Image Captioning," *Proceedings of the International Joint Conferences on Artificial Intelligence, 5*.

[12] Y. Zhang, X. Shi, S. Mi, and X. Yang. 2021. "Image captioning with transformer and knowledge graph," *Pattern Recognition Letters*, 143: 43–49. doi: 10.1016/j.patrec.2020.12.020.

[13] M. R. S. Mahadi, A. Arifianto, and K. N. Ramadhani. 2020. "Adaptive Attention Generation for Indonesian Image Captioning," *2020 8th International Conference on Information and Communication Technology (ICoICT)*. doi: 10.1109/icoict49345.2020.9166244.

[14] E. Mulyanto, E. I. Setiawan, E. M. Yuniarno, and M. H. Purnomo. 2019. "Automatic Indonesian Image Caption Generation using CNN-LSTM Model and FEEH-ID Dataset," *2019 IEEE International Conference on Computational Intelligence and Virtual Environments for Measurement Systems and Applications (CIVEMSA)*. doi: 10.1109/civemsa45640.2019.9071632.

[15] A. A. Nugraha, A. Arifianto, and Suyanto. 2019. "Generating Image Description on Indonesian Language using Convolutional Neural Network and Gated Recurrent Unit," *2019 7th International Conference on Information and Communication Technology (ICoICT)*. doi: 10.1109/icoict.2019.8835370.

[16] J. Devlin, S. Gupta, R. Girshick, M. Mitchell, and C. L. Zitnick. 2015. "Exploring nearest neighbor approaches for image captioning." *arXiv preprint arXiv:1505.04467*.

[17] O. Vinyals, A. Toshev, S. Bengio, and D. Erhan. 2015. "Show and tell: A neural image caption generator," *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. doi: 10.1109/cvpr.2015.7298935.

[18] MD. Z. Hossain, F. Sohel, M. F. Shiratuddin, and H. Laga. 2019. "A Comprehensive Survey of Deep Learning for Image Captioning," *ACM Computing Surveys*, 51(6): 1–36. doi: 10.1145/3295748.

[19] S. Pa Pa Aung, W. Pa Pa, and T. Lay Nwe. 2020. "Automatic Myanmar image captioning using CNN and LSTM-based language model." *Proceedings of the 1st Joint Workshop on Spoken Language Technologies for Under-resourced languages (SLTU) and Collaboration and Computing for Under-Resourced Languages (CCURL)*. 2020. "Automatic Myanmar Image Captioning using CNN and LSTM-Based Language Model," 139-143.

[20] S. Sukhbaatar, A. Szlam, J. Weston, and R. Fergus. 2015. "End-to-end memory networks." *Advances in neural information processing systems* 28.

[21] A. Vaswani et al. 2017. "Attention is all you need." *Advances in neural information processing systems* 30.

[22] Y. Pan, T. Yao, Y. Li, and T. Mei. 2020. "X-Linear Attention Networks for Image Captioning," *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. doi: 10.1109/cvpr42600.2020.01098.

[23] T. Xian, Z. Li, C. Zhang, and H. Ma. 2022. "Dual Global Enhanced Transformer for image captioning," *Neural Networks*, 148: 129–141. doi: 10.1016/j.neunet.2022.01.011.

[24] S. Dubey, F. Olimov, M. A. Rafique, J. Kim, and M. Jeon. 2023. "Label-attention transformer with geometrically coherent objects for image captioning." *Information Sciences* 623: 812-831.

[25] W. Liu, S. Chen, L. Guo, X. Zhu, and J. Liu. 2021. "Cptr: Full transformer network for image captioning." *arXiv preprint arXiv:2101.10804*.

[26] Z. Ren, S. Gou, Z. Guo, S. Mao, and R. Li. 2022. "A Mask-Guided Transformer Network with Topic Token for Remote Sensing Image Captioning," *Remote Sensing*, 14(12): 2939. doi: 10.3390/rs14122939.

[27] X. Chen et al. 2015. "Microsoft coco captions: Data collection and evaluation server." *arXiv preprint arXiv:1504.00325*.

[28] B. A. Plummer, L. Wang, C. M. Cervantes, J. C. Caicedo, J. Hockenmaier, and S. Lazebnik. 2016. "Flickr30k Entities: Collecting Region-to-Phrase Correspondences for Richer Image-to-Sentence

Models," *International Journal of Computer Vision*, 123(1): 74–93. doi: 10.1007/s11263-016-0965-7.

[29] F. M. Shah, M. Humaira, M. A. R. K. Jim, A. Saha Ami, and S. Paul. 2021. "Bornon: Bengali Image Captioning with Transformer-Based Deep Learning Approach," *SN Computer Science*, 3(1). doi: 10.1007/s42979-021-00975-0.

[30] D. H. Fudholi and R. A. N. Nayoan. 2022. "The Role of Transformer-based Image Captioning for Indoor Environment Visual Understanding," *International Journal of Computing and Digital Systems*, 12(3): 479–488. doi: 10.12785/ijcds/120138.

[31] A. Gholamy, V. Kreinovich, and O. Kosheleva. 2018. "Why 70/30 or 80/20 relation between training and testing sets: A pedagogical explanation.".

[32] H. Rashid, A. S. M. R. Al-Mamun, M. S. R. Robin, M. Ahasan, and S. M. T. Reza. 2016. "Bilingual wearable assistive technology for visually impaired persons," *2016 International Conference on Medical Engineering, Health Informatics and Technology (MediTec)*. doi: 10.1109/meditec.2016.7835386.

[33] G. Huang, Z. Liu, L. Van Der Maaten, and K. Q. Weinberger. 2017. "Densely Connected Convolutional Networks," *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. doi: 10.1109/cvpr.2017.243.

[34] M. Tan and Q. v. Le. 2019. "Efficientnet: Rethinking model scaling for convolutional neural networks." *International conference on machine learning*. PMLR.

[35] D. H. Fudholi *et al.* 2021. "Image Captioning with Attention for Smart Local Tourism using EfficientNet," *IOP Conference Series: Materials Science and Engineering*, 1077(1): 012038. doi: 10.1088/1757-899x/1077/1/012038.

[36] V. Maeda-Gutiérrez *et al.* 2020. "Comparison of Convolutional Neural Network Architectures for Classification of Tomato Plant Diseases," *Applied Sciences*, 10(4):1245. doi: 10.3390/app10041245.

[37] F. Chollet. 2017. "Xception: Deep Learning with Depthwise Separable Convolutions," *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. doi: 10.1109/cvpr.2017.195.

[38] S. Vellakani and I. Pushbam. 2020. "An enhanced OCT image captioning system to assist ophthalmologists in detecting and classifying eye diseases," *Journal of X-Ray Science and Technology*, 28(5): 975–988. doi: 10.3233/xst-200697.

[39] C. Szegedy, S. Ioffe, V. Vanhoucke, and A. Alemi. 2017. "Inception-v4, inception-resnet and the impact of residual connections on learning." *Proceedings of the AAAI conference on artificial intelligence,* 31(1).

[40] K. Papineni, S. Roukos, T. Ward, and W.-J. Zhu. 2001. "BLEU," *Proceedings of the 40th Annual Meeting on Association for Computational Linguistics - ACL '02*. doi: 10.3115/1073083.1073135.

[41] C. Y. Lin. 2004. "Rouge: A package for automatic evaluation of summaries." *Text summarization branches out*.

[42] A. Lavie and A. Agarwal. 2007. "Meteor: an automatic metric for MT evaluation with high levels of correlation with human judgments." *In Proceedings of the Second Workshop on Statistical Machine Translation (StatMT '07)*. Association for Computational Linguistics, USA, 228–231.

[43] R. Vedantam, C. L. Zitnick, and D. Parikh. 2015. "CIDEr: Consensus-based image description evaluation," *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. doi: 10.1109/cvpr.2015.7299087.

[44] M. Huh, P. Agrawal, and A. A. Efros. 2016. "What makes ImageNet good for transfer learning?." *arXiv preprint arXiv:1608.08614*.

[45] A. A. Nugraha, A. Arifianto, and Suyanto. 2019. "Generating Image Description on Indonesian Language using Convolutional Neural Network and Gated Recurrent Unit," *2019 7th International Conference on Information and Communication Technology (ICoICT)*. doi: 10.1109/icoict.2019.8835370.

[46] M. R. S. Mahadi, A. Arifianto, and K. N. Ramadhani. 2020. "Adaptive Attention Generation for Indonesian Image Captioning," *2020 8th International Conference on Information and Communication Technology (ICoICT)*. doi: 10.1109/icoict49345.2020.9166244.