# SuperCap: Multi-resolution Superpixel-based Image Captioning

Henry Senior
Queen Mary University of London
London, UK

h.senior@qmul.ac.uk

Gregory Slabaugh
Queen Mary University of London
London, UK

g.slabaugh@qmul.ac.uk

Luca Rossi
The Hong Kong Polytechnic University
Hong Kong

luca.rossi@polyu.edu.hk

Shanxin Yuan
Queen Mary University of London
London, UK

shanxin.yuan@qmul.ac.uk

## Abstract

*It has been a longstanding goal within image captioning to move beyond a dependence on object detection. We investigate using superpixels coupled with Vision Language Models (VLMs) to bridge the gap between detector-based captioning architectures and those that solely pretrain on large datasets. Our novel superpixel approach ensures that the model receives object-like features whilst the use of VLMs provides our model with open set object understanding. Furthermore, we extend our architecture to make use of multi-resolution inputs, allowing our model to view images in different levels of detail, and use an attention mechanism to determine which parts are most relevant to the caption. We demonstrate our model's performance with multiple VLMs and through a range of ablations detailing the impact of different architectural choices. Our full model achieves a competitive CIDEr score of* 136.9 *on the COCO Karpathy split.*

## 1. Introduction

Image captioning, the challenging task of developing a model that describes a visual scene in natural language, has received considerable attention in recent years from the computer vision community. It is a task with great scientific impact (*e.g.*, ensuring models can truly describe the scene) and societal impact (*e.g.*, captioning models provide benefits in terms of accessibility through generating alt-text for websites). While early works focused on a standard Encoder-Decoder architecture [16] using a Convolutional Neural Network (CNN) to produce a feature vector that is decoded into a caption via a Recurrent Neural Network (RNN), later works moved to adopt object detectors.

There are typically two approaches to image captioning:

detector-based (Bottom Up Top Down (BUTD) [4] and its descendants [7, 43, 44]) or 'detector-free' models, the latter of which are mostly massively pretrained large Vision-Language Models (VLMs) built on top of the Transformer architecture [34] and pretrained on large datasets such as LAION-5B [31]. These models produce strong global features but risk missing out on more nuanced finer grain detail, something at which detector-based methods naturally excel. In this work, we investigate taking the best of both worlds, answering the question, 'how can we leverage VLMs to produce high quality captions whilst utilising finer object details?'. We achieve this through the use of superpixels: by segmenting the image into $M$ superpixels, we leverage the zero shot classification capabilities of VLMs to replace the object detector. To ensure finer-details are captured, we incorporate a multi-resolution approach, where the image is divided into a set of superpixels of different resolutions. We also make use of a whole image feature to maintain the global context.

Early deep learning approaches to image captioning focused on using a CNN to represent the entire image using a unique latent feature. This feature was then decoded into a caption using an RNN. Although this approach was initially successful, it produced captioning models that would focus on the global context of an image and would therefore miss specific objects [16]. As object detectors improved, they began to be used in image captioning models to address the weaknesses of earlier approaches. Following on from Anderson *et al.*'s BUTD [4] work, a large number of image captioning architectures based on object detectors have been produced [7, 43, 44]. Whilst these models all provide strong performance and novel ways of utilising the objects in different ways, they all share a common flaw - they are fundamentally limited by their object detectors. Specifically, 1) the closed set nature of object detectors, which

results in only objects in the training data of the detector being identified, and 2) the expensive computation required by object detectors (FasterRCNN [29] with a ResNet-50 [11] takes 447 GFLOPs). As object detectors are only as good as their training data, in scenarios where the detector misclassifies an object (or is unable to classify it), the caption will be incorrect.

The closed-set nature of detector-based image captioning, alongside the development of large pretrained VLMs has led researchers to develop detector-free captioning architectures [36]. These approaches largely address the problems of detector-based architectures by scaling the Transformer and the data used to train the model. Given the impressive Transformer scaling laws [15], the development in computing infrastructure, and the increasing accessibility of ever larger datasets, this approach has proven to be the dominant approach in recent years.

Methods aiming to unify detector-based and detector-free are limited [10]. Their architecture mirrors that of the detector-free models [36], a ViT combined with a transformer language model. However, they include a novel 'concept token network' that predicts concepts from the hidden layers of the ViT, essentially providing object detection via proxy. This approach leads to a complex model that is difficult to train. In comparison, our model is easily trainable end-to-end following standard procedures for an image captioning model trained on COCO [6, 24].

In this paper we present SuperCap, a novel detector-free image captioning architecture that joins together the benefits of detector-based and detector-free approaches to image captioning. The key contributions of our paper are as follows:

1. SuperCap is the first image captioning model to be based on superpixels instead of patches or object bounding boxes. Our model performs competitively on standard benchmarks whilst having an architecture that scales appropriately to future developments within the computer vision community.
2. To the best of our knowledge, we are the first image captioning architecture to make explicit use of multiple resolutions, thus ensuring that our model captures the local context (like detector-based models) and the global context (like pretrained detector-free models). We empirically show the advantage of our multi-resolution approach through extensive ablation studies.

The remainder of this paper is organised as follows. Section 2 reviews the related work, while Sections 3 introduces SuperCap. In Section 4 we discuss the implementation details of our model as well as of the training and evaluation protocol. We then conduct an extensive experimental evaluation in Section 5 and we conclude the paper in Section 6.

## 2. Related Work

In this section we detail the related work and background to provide context for the literature that works sits alongside. We start by giving an overview of image captioning, followed by detector-free approaches to this problem. The section concludes with a brief overview of superpixels.

### 2.1. Image Captioning

Given an input image $\mathcal{I}$, an image captioning model is tasked with producing a caption $\mathcal{C} = \{w_i\}_{i=1}^N$ comprised of $N$ words in a predetermined dictionary $w_i \in \mathcal{D}$. Most image captioning methods begin by defining a region extractor $\varphi$ which produces a set of $M$ regions $\mathcal{R} = \{\mathbf{r}_i\}_{i=1}^M$ detected within the image, *i.e.*,

$$\varphi : \mathcal{I} \mapsto \mathcal{R} . \tag{1}$$

Typically regions correspond to the bounding boxes of objects [4, 40, 43] or patches [10]. These regions are then processed by a region feature extractor $\Psi$ which produces a set of latent features $\bar{\mathcal{R}} = \{\mathbf{r}_i \in \mathbb{R}^N\}_{i=1}^M$, *i.e.*,

$$\Psi : \mathcal{R} \mapsto \bar{\mathcal{R}} . \tag{2}$$

In the case of bounding boxes, $\Psi$ is typically a ResNet [11] model [4, 40, 43] or in the case of patches, $\Psi$ is comprised of a Vision Transformer [9].
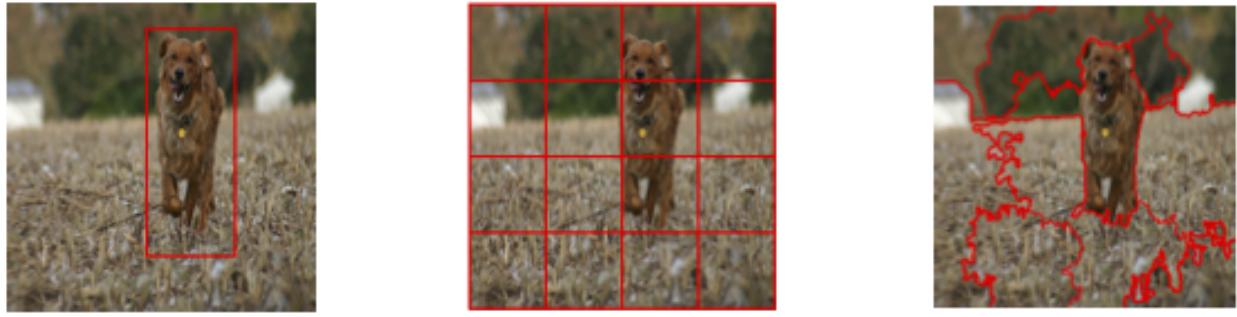
Finally, the main captioning model $\Omega$ uses the latent features $\bar{\mathcal{R}}$ alongside the previously predicted word $w_{t-1}$ to predict the next word in $\mathcal{C}$, *i.e.*,

$$\Omega : (\bar{\mathcal{R}}, w_{t-1}) \mapsto w_t . \tag{3}$$

Traditionally [4] this has been implemented by a dual LSTM but more contemporary approaches [10, 36] make use of a Transformer.

Typically, $\varphi$, $\Psi$, and $\Omega$ are neural networks. We note $\varphi$ and $\Psi$ are not usually included in the end-to-end training of the captioning network, instead it is trained to complete a standalone task such as object detection.

Most models use the $\varphi$ and $\Omega$ defined in [4]. That is, a Visual Genome [20] based FasterRCNN [29] object detector with a ResNet backbone [11] representing the region operators ($\varphi$ and $\Psi$) and a language model consisting of a dual LSTM with attention providing the captioning model ($\Omega$). The majority of methods (especially those in GNN-based image captioning [32]) focus on augmenting the region encoder that produces rich latent features that capture the necessary information to produce high quality captions. Yao *et al*. [40] implemented a GNN-based architecture that used both semantic and spatial graphs. Further approaches built on the foundation laid by Yao *et al*. [40], with Zhong *et al*. [44] making use of a GCN [18] to perform semantic graph decomposition on the semantic graph. They address the problem of which nodes and edges of the graph to include in the final caption.

(a) **Bounding Boxes** One common way in which regions are extracted is through the use of bounding boxes. This is used by Anderson *et al.* [4] and subsequent works [7, 40, 44].

(b) **Patches** A more contemporary approach to region extraction is through the use of patches [10]. However, it results in objects being divided between features, as shown above.

(c) **Superpixels** Our approach utilises superpixels to create regions that encompass the whole image but ensure that objects are not unnecessarily divided between features.

Figure 1. **Visual Representation of Different Region Extraction Techniques.** Here we show the three different ways in which regions can be extracted from images, bounding boxes, patches, and superpixels.

## 2.2. Detector-Free Image Captioning

One way in which research has addressed the issues introduced by object detectors is Vision-Language Pretraining (VLP), the practice of using massive amounts of data (typically from multiple datasets) to give the model a rich image-language embedding suitable for a large number of downstream tasks. Pretrained models are then finetuned on specific datasets, honing their performance at a particular task. This approach has been the main approach to detector-free image captioning [21, 36, 37] as large Vision-Language models learn to identify objects as an emergent behaviour learnt from the vast amounts of data.

More recently, Fang *et al.* [10] proposed a technique based on ViT [9] that boasts strong performance without the need for VLP. Their approach uses a novel Concept Token Network, a Transformer-based [34] network that takes the intermediary tokens of ViT [9] as input, and predicts the top-$k$ most likely concept tokens from a corpus generated from COCO[6, 24]. These concept tokens are then concatenated with the final output feature of ViT to create a final latent feature vector that is decoded into a caption by a BERT-based [8] language model. Fang *et al.*'s [10] Concept Token Network is essentially used to perform object detection via proxy and leads to a network architecture that requires an involved training strategy. Their reliance on large network architectures such as ViT and BERT results in an overall model with 190 million parameters. Meanwhile, the architecture proposed in this paper sits at only 93 million trainable parameters in its main configuration.

Another approach to detector-free image captioning, which uses non-traditional region definitions (segmentation masks), is Segment and Caption Anything [13]. This model combines SAM [19] and an LLM with a thin trainable network allowing for parameter efficient training of a region captioning model. However, to achieve strong performance, it requires large amounts of weak pretraining, something our model avoids via the use of VLMs.

For a more detailed survey on image captioning, we direct the reader to Stefanini *et al.* [33] and Senior *et al.* [32], the latter of which has a particular focus on GNN-based approaches.

## 2.3. Superpixels

Superpixel segmentation is a well established technique within computer vision for grouping pixels together into coherent and homogeneous boundaries. Unlike bounding boxes, superpixels are non-rigid and therefore easily handle edges of objects that are not straight. There are many approaches to producing superpixels for an image, including clustering, watershed, density, and graph-based ones, to name a few. One widely adopted approach is the Simple Linear Iterative Clustering (SLIC) proposed by Achanta *et al.* [1]. As a cluster-based approach, SLIC works by converting the pixels into a 5-dimensional space comprised of the $(L, a, b)$ pixel colour values of the CIELAB colour space and their $(x, y)$ location in the image. Once the pixels for the image are converted into this space, SLIC then clusters the values into $K$ groups, which then go on to become the $K$ superpixels. Although the SLIC superpixel algorithm appears simple at first glance, it provides an elegant and robust technique to provide the superpixel segmentation of a given image.

## 3. Method

In this section we introduce our architecture, both for the superpixel preprocessing and the captioning model itself. We

start by outlining how our model extracts superpixel regions from a given image, following on with how we convert superpixels to region features which are used for captioning. Finally, we conclude by detailing how the architecture is extended to support multi-resolution.

### 3.1. Superpixel Regions Extraction ($\varphi$)

The region extraction forms the first component of our architecture. Given an image, $\mathcal{I}$, it must produce a set of distinct regions $\mathcal{R}$ that can be used as input to a captioning architecture *i.e.* $\mathcal{R} = \varphi(\mathcal{I})$. Whilst patches [10] or object detectors [4] have been traditionally used for this task, we instead take the approach of using superpixels (shown in Figure 1). We choose to make use of superpixel regions as they produce contiguous groups of pixels that make up an object. This ensures that our regions are more homogenous and reduces the risk of objects being broken up unnecessarily.

Specifically, we use the SLIC [1] superpixel segmentation technique for its computational efficiency [2]. Using SLIC, our model will produce $M_k$ superpixel regions, $\mathcal{R}_k = \{r_i\}_{i=1}^{M_k}$. To ensure that our model is able to capture different levels of image details, we extend the superpixel region extraction to include multiple resolutions. This results in $\mathcal{R} = \{\mathcal{R}_1, ..., \mathcal{R}_k\}$ regions being produced for $k$ resolutions.

### 3.2. Regions Encoder ($\Psi$)

Once a set of regions have been defined, they must be converted into latent features, $\bar{\mathcal{R}}$, that can be used in the captioning model *i.e.* $\bar{\mathcal{R}} = \Psi(\mathcal{R})$. To achieve this, we look towards VLMs such as CLIP [28] and BLIP [21]. These large models produce features that combine semantic elements of both visual elements and language elements and are therefore a good fit for image captioning. We produce latent features for superpixels by generating a bounding box around each superpixel and feeding those extracted pixels into the VLM. This approach allows features to be generated for coherent regions of the image rather than smaller broken up regions, as is the case for patch based methods. For each resolution, $k$, the latent features are defined as

$$\bar{\mathcal{R}}_k = \{\mathbf{r}_i \in \mathbb{R}^N\}_{i=1}^M, \qquad \mathbf{r}_i = VLM(r_i) \quad (4)$$

Additionally, we also introduce a global image feature $\mathbf{r}_0 \in \mathbb{R}^N$ to capture the wider context of the image.

### 3.3. Captioning Model ($\Omega$)

Once the superpixel region features ($\bar{\mathcal{R}}$) have been produced, the captioning model must refine them to better account for the wider context of the image. Our model implements this through an Encoder-Decoder transformer [34] model. In order to support multiple resolutions, we stack $k$

Encoders, one per resolution, and concatenate their output before inputting the attended features into the decoder.

The encoder consists of a stack of a classic Transformer comprised of a multi-head self-attention blocks followed by a feed-forward network. The superpixel features for a given resolution $\bar{\mathcal{R}}_k$ are passed through the Transformer to calculate their corresponding attended features $\bar{\mathcal{R}}_k'$. This is achieved by first converting the input features $\mathcal{R}$ into corresponding queries, keys, and values by reprojecting them according to

$$Q = \bar{\mathcal{R}}_k W_Q, \quad K = \bar{\mathcal{R}}_k W_K, \quad V = \bar{\mathcal{R}}_k W_V, \quad (5)$$

where $W_Q, W_K, W_V$ are all learnt weight matrices. The attention weights ($\alpha$) are then calculated with

$$\alpha = \frac{QK^T}{\sqrt{d_k}}, \quad (6)$$

where the scaling value of $d_k$ is the dimension of projected superpixel features. These attention weights are then used to scale the input features following

$$\text{head}(\bar{\mathcal{R}}_k) = \text{softmax}(\alpha)V. \quad (7)$$

For multi-head self-attention, the results of multiple heads are concatenated together (denoted with $[\cdot\|\cdot]$) and reprojected with a learnt weight matrix $W_O$

$$\text{Multihead}(Q, K, V) = [head_1\|...\|head_n]W_O \quad (8)$$

Finally, the latent features are produced with a feed-forward network

$$\bar{\mathcal{R}}_k' = \sigma(W_1\bar{\mathcal{R}}_k)W_2, \quad (9)$$

where $\sigma$ is the ReLU non-linearity and $W_1$, $W_2$ are both learnable weight matrices. The encoder consists of a stack of multi-head self-attention blocks followed by a feed-forward network. The superpixel features $\mathcal{R} = \{\mathbf{r}_i\}_{i=0}^M$ are passed through the self-attention mechanism to calculate their corresponding attention weights.

Once done for each resolution, the attended features are concatenated together to produce the final set of image features, $\mathbf{R}$ that will be used by the decoder to produce the caption.

$$\mathbf{R} = [\bar{\mathcal{R}}_1'\|..\|\bar{\mathcal{R}}_k'] \quad (10)$$

where $[\cdot\|\cdot]$ denotes concatenation.

The decoder component of SuperCap then uses $\mathbf{R}$ alongside the previously generated word $w_{t-1}$ to autoregressively predict the next most likely word $w_t$.

## 4. Implementation

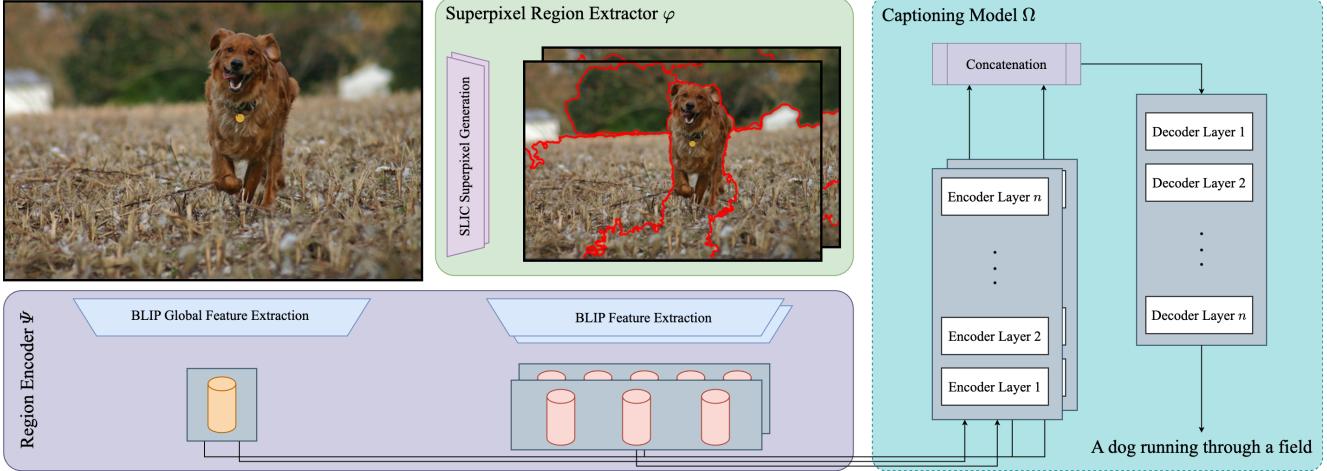In the interest of reproducibility, in section we detail our implementation choices as well as the training and evaluation

Figure 2. **Architecture Diagram** Our SuperCap model, like other image captioning models, has three main components: 1) the region extractor ($\varphi$), 2) the region encoder ($\Psi$), and 3) the captioning model ($\Omega$). The initial encoder takes in an image and uses BLIP to extract a global image feature. It simultaneously extracts superpixel regions for the image at different resolutions, which are in turn converted to BLIP region features. The secondary encoder consists of a stack of encoder transformers, one per resolution, which receive both the global feature and superpixel features for their resolution. Latent features produced by these encoders are then concatenated together and fed into a decoder transformer which autoregressively produces the final image caption.

protocol used. Our code is publicly available [1] [2]

## 4.1. Model

Our region extractor ($\varphi$) consists of a SLIC [1] superpixel segmentation process that computes $M$ superpixels from an image $\mathcal{I}$. Each superpixel is then converted into an $N$-dimensional feature vector by our region encoder ($\Psi$), which is implemented via either BLIP [21] ($N = 768$) or CLIP [28] ($N = 512$). We demonstrate the effectiveness of both models in Section 5. Once the superpixel features $\mathcal{R} = \{\mathbf{r}_i \in \mathbb{R}^N\}_{i=1}^M$ are computed, an additional global feature $\mathbf{r}_0 \in \mathbb{R}^N$ is extracted using the same feature extractor. Our captioning model ($\Omega$) is implemented via a six layer, eight head encoder-decoder Transformer. It makes use of sinusoidal positional encoding and predicts the next most likely token across a one-hot encoded dictionary of $10,000$ words. The dictionary is generated following the standard practice in image captioning [16] of using words that appear more than five times in the corpus of captions.

## 4.2. Dataset

The Microsoft Common Objects in Context (COCO) dataset [6, 24] is a widely used dataset for image captioning. While the original release incorporated an evaluation server with a hidden test set, the community has unanimously adopted what has become known as the Karpathy

Table 1. **Comparison of Superpixel Resolutions**. We compare the impact of dividing the image into different numbers of superpixels using CLIP as the feature encoder. Results are for the Karpathy Test Split [16] of the COCO dataset [6, 24], where B-4, M, R, C, S denote BLEU@4, METEOR, ROUGE-L, CIDEr and SPICE scores respectively. Best results in **bold**.

| # Superpixels | B-4 ↑ | M ↑ | R ↑ | C ↑ | S ↑ |
|---|---|---|---|---|---|
| 1 | 35.0 | 27.2 | 56.4 | 119.1 | 20.6 |
| 2 | 21.5 | 20.1 | 46.0 | 71.4 | 13.4 |
| 3 | 25.5 | 22.4 | 49.4 | 88.0 | 15.9 |
| 4 | 34.7 | 27.2 | 56.4 | 118.4 | 20.7 |
| 5 | 37.5 | 28.4 | 58.1 | 127.1 | 22.1 |
| 6 | 36.0 | 27.7 | 57.0 | 122.5 | 21.3 |
| 7 | 36.5 | 27.9 | 57.3 | 122.7 | 21.4 |
| 8 | 36.3 | 27.9 | 57.3 | 123.2 | 21.5 |
| 9 | 37.1 | 28.1 | 57.6 | 125.4 | 21.8 |
| 10 | 37.9 | 28.6 | **58.2** | **128.4** | 22.4 |
| 15 | **38.0** | **28.7** | **58.2** | **128.4** | **22.5** |
| 20 | 37.8 | 28.5 | 57.9 | 127.0 | 22.3 |
| 25 | 37.2 | 28.4 | 57.8 | 126.58 | 22.3 |
| 50 | 36.0 | 27.9 | 56.9 | 121.7 | 21.5 |
| 75 | 35.4 | 27.6 | 56.3 | 119.6 | 21.5 |
| 100 | 34.5 | 27.1 | 55.7 | 115.9 | 20.9 |

Split [16]. This is a train/val/test split of: $113,287$ / $5000$ / $5000$. Each image is accompanied by five human generated captions, some with differing focuses or verbosity. Overall, it provides a wide range of real life imagery alongside high

---

[1]Region Extraction: https://anonymous.4open.science/r/superpixel-features-1754/

[2]Captioning code: https://anonymous.4open.science/r/DEFIGNN-B91F/README.md

Table 2. **Comparison of Superpixel Resolutions with Global Feature**. We compare the impact of dividing the image into different numbers of superpixels using CLIP as the feature encoder alongside the inclusion of a global image feature. Results are for the Karpathy Test Split [16] of the COCO dataset [6, 24], where B-4, M, R, C, S denote BLEU@4, METEOR, ROUGE-L, CIDEr and SPICE scores respectively. Best results in **bold**.

| # Superpixels | B-4 ↑ | M ↑ | R ↑ | C ↑ | S ↑ |
|---|---|---|---|---|---|
| 1 | 35.0 | 27.2 | 56.4 | 119.1 | 20.6 |
| 10 | **38.1** | 28.6 | 58.4 | **129.9** | 22.4 |
| 15 | 38.0 | **28.7** | 58.4 | 129.2 | **22.5** |
| 25 | 38.0 | **28.7** | **58.5** | 129.8 | 22.4 |
| 50 | 37.9 | 28.6 | 58.3 | 129.1 | 22.2 |
| 75 | 36.9 | 28.3 | 57.9 | 126.0 | 21.9 |
| 100 | 36.3 | 27.8 | 57.1 | 124.4 | 21.5 |

quality captions.

### 4.3. Training

We train our model for 30 epochs using the Adam [17] optimiser with the same learning rate scheduling strategy as [7, 34] during the cross-entropy loss phrase. Following standard practices in image captioning [10, 30, 40, 43, 44], we then use the CIDEr-based reinforcement learning strategy to prevent our model suffering from the exposure bias problem. The finetuning is completed over an additional 30 epochs, optimised by Adam [17] with a fixed learning rate of $5e - 6$.

### 4.4. Evaluation

In line with other works in the field [10, 30, 40, 43, 44], we use a number of evaluation metrics [3, 5, 23, 27, 35].

BLEU@$N$ [27] is a precision-based n-gram matching metric originally adopted from machine translation. It has long been used in image captioning and is one of the key evaluation metrics. Additionally from natural language processing, ROUGE [23] is a set of metrics used to evaluate the quality of automatic summarisation and machine translation systems. It takes into account: n-gram recall between the candidate and reference set, a comparison of the longest common sub-sequence, a comparison of the weighted longest common sub-sequence, and finally, the skip-bigram co-occurrence statistic. The final natural language processing evaluation used in image captioning is METEOR [5], which uses the $F_{mean}$ to take into account both recall and precision. CIDEr [35] and SPICE [3] are both evaluation metrics designed specifically for image captioning and have therefore become important benchmarks since their introduction.

Finally, note that during inference, we use a beam width of 5 and a max decoding length of 20 words.

Table 3. **Comparison of Different Initial Encodings**. We compare the impact of various initial encodings with all ablations including an appropriate global feature. Results are for the Karpathy Test Split [16] of the COCO dataset [6, 24], where B-4, M, R, C, S denote BLEU@4, METEOR, ROUGE-L, CIDEr and SPICE scores respectively. Best results in **bold**.

| Encoding | B-4 ↑ | M ↑ | R ↑ | C ↑ | S ↑ |
|---|---|---|---|---|---|
| 15 SPs & CLIP | 38.0 | 28.7 | 58.4 | 129.2 | 22.5 |
| 15 SPs & BLIP | **40.1** | **29.6** | **59.6** | **135.7** | **23.3** |
| 16 Patches & BLIP | 39.2 | 29.2 | 59.2 | 133.2 | 23.0 |

## 5. Results

Given the setting described in the previous Section, we perform a number of experiments aimed at evaluating the impact of our architectural choices on model performance as well as comparing SuperCap with state-of-the-art (SOTA) alternatives. We first compare the impact of various superpixel resolutions with the model running in a single resolution capacity. After seeing the strong performance of a single superpixel feature (essentially a global feature), we evaluate the impact of including a global superpixel feature into the architecture. Next we compare the impact of different region extractors and region encoders, comparing CLIP and BLIP superpixel features with BLIP patch features. Finally, we compare four different multi-resolution approaches before assembling our final model architecture to perform comparison against SOTA models.

### 5.1. Impact of Different Superpixel Resolutions

Using CLIP to encode the superpixel region features, we compare the the performance of our single resolution model when changing the number of superpixels that the image is divided into (shown in Table 1). We observe that whilst a single superpixel performs well (as it is essentially just a global context feature), there is a sweet spot in performance around 10 to 15 superpixels per image. Increasing the number of superpixels results in a decreasing performance. We hypothesise that this is because the superpixels are too small to contain enough of the image for the VLM to produce a meaningful feature.

### 5.2. Impact of Global Features

Using CLIP to encode the superpixel region features alongside a global image feature, we compare the the performance of our single resolution model when changing the number of superpixels that the image is divided into (shown in Table 2). By including the global feature, all superpixel resolutions see an improvement in performance when compared with their counterparts in Table 1. Note also that while the impact of adding the global feature appears to be uneven across the range of superpixel resolution (*i.e.*, it has

Table 4. **Comparison of Different Multi-Resolution Approaches**. We compare the impact of various multi-resolution approaches alongside the impact of the inclusion of a global image feature $r_0$ (denoted with FC in the table). Results are for the Karpathy Test Split [16] of the COCO dataset [6, 24], where B-4, M, R, C, S denote BLEU@4, METEOR, ROUGE-L, CIDEr and SPICE scores respectively. Best results in **bold**.

| Encoding | B-4 ↑ | M ↑ | R ↑ | C ↑ | S ↑ |
|---|---|---|---|---|---|
| Method 1 | 37.6 | 28.3 | 57.9 | 127.3 | 22.2 |
| Method 1 & FC | 37.9 | 28.7 | 58.4 | 129.5 | 22.5 |
| Method 2 | 38.5 | 28.8 | 58.7 | 129.9 | 22.7 |
| Method 2 & FC | 38.8 | 29.1 | 59.0 | 131.8 | **22.8** |
| Method 3 | 36.4 | 27.8 | 57.2 | 123.4 | 21.2 |
| Method 3 & FC | 37.5 | 28.3 | 58.1 | 127.2 | 21.8 |
| Method 4 | 37.4 | 28.2 | 57.8 | 125.9 | 22.0 |
| Method 4 & FC | 36.0 | 27.5 | 57.3 | 121.6 | 20.9 |
| Dual & FC | **39.2** | **29.2** | **59.1** | **132.4** | **22.8** |

a larger positive impact for higher resolutions), this only minimally affects the optimum.

## 5.3. Comparison of Different Initial Encoders

Here we explore the impact of different initial encoders by comparing both CLIP [28] and BLIP [21] embeddings of the same SLIC [1] superpixels. As shown in Table 3, for the same superpixels, BLIP is superior in performance across all the benchmarks. We then hold the embedding model constant and change the region representation from superpixels to the more classic patches. We create 16 patches by dividing the image into a $4 \times 4$ grid and produce BLIP features for each patch. Whilst traditional patching methods like those used in Vision Transformer [9] use a smaller grid with more regions, we demonstrate in Tables 1 and 2 that higher region counts correlate with lower performance. Therefore, in order to ensure a fair comparison we use a comparable amount of regions.

## 5.4. Different Multi-Resolution Approaches

Extending our architecture to support multiple resolutions can be achieved in a variety of different ways. **Method 1** is the most straightforward approach, and involves keeping the architecture the same and simply concatenating the different resolutions together. This essentially increased the number of tokens that our model receives. **Method 2**, which proved to be the most effective, involves stacking the secondary encoder, one per resolution. The outputs of these encoders are then concatenated and fed into a single decoder language model. **Method 3** and **Method 4** investigate using a mixture of experts (MoE) approach to the decoder, following previous captioning work [39]. **Method 3** passes the resolutions through a single encoder with each resolu-

tion going to its own decoder. The outputs of the decoders are then passed through a soft routing network to select the final token. **Method 4** follows the same MoE approach but uses one encoder per resolution.

We test each of the four methods with CLIP [28] superpixels generated at four resolutions, 25, 50, 75, 100 to ensure that the models receive inputs that cover a wide range of scale. Our ablation concludes by testing the best multi-resolution approach (**Method 2**) with dual inputs (referred to in Table 4 as 'dual') 10 and 25 superpixels.

## 5.5. Comparison with SOTA

In order to demonstrate the performance of our model, we compare it against similar detector-free image captioning models. To the best of our knowledge, as this is a new field, there only exists [10]. We therefore include some common detector-based image captioning models in order to place our work within a wider context.

When compared against the next best model [10], our model has considerably fewer parameters. As shown in Table 5, we match or exceed the performance during the important CIDEr optimisation stage. This optimisation stage was introduced to reduce the exposure bias problem caused by training language models with teacher forcing [30], which is used during the cross-entropy loss training. Given the significant reduction in parameter count, alongside the strong performance, it is clear that our architecture is very parameter efficient.

## 5.6. Qualitative Results

In order to demonstrate our model in action, we present a representative selection of five randomly selected Karpathy test set [16] images, accompanied by our generated model and their ground truth captions. Our model produces both accurate and descriptive captions and demonstrates an ability to identify individual objects (such as the chair and bed in the top left image). Even in cases where the caption does not perfectly match the ground truth, such as in the bottom right picture, the caption is nevertheless an accurate description of the image.

Like all image captioning models, SuperCap does also suffer from mild hallucination, as in the case of the top right image. Our model describes a 'train' being present on top of the cake, whereas in actuality it is a children's character. Another possible shortcoming present in Figure 3 is the production of identical captions for similar images, as in the case for the two central images.

## 6. Conclusion

In this paper we introduced SuperCap, a superpixel-based image captioning architecture with support for multi-resolution. Our architecture makes use of the zero shot classification capabilities of large pretrained models,

Table 5. **Comparison With SOTA**. The below table compares our SuperCap to the current SOTA models that perform detector-free image captioning without vision-language pretraining. Results are for the Karpathy Test Split [16] of the COCO dataset [6, 24]. End-to-end pretrained models are denoted with a *. Best scores for *'detector'* based models are indicated in *italics* and the best scores for **detector-free** models are indicated in **bold**.

| Model | # Trainable Params ↓ | BLEU-4 ↑ | METEOR ↑ | ROUGE-L ↑ | CIDEr ↑ | SPICE ↑ |
|---|---|---|---|---|---|---|
| Detector Models | | | | | | |
| RFNet [14] | - | 36.5 | 27.7 | 57.3 | 121.9 | 21.2 |
| BUTD [4] | - | 36.3 | 27.7 | 56.9 | 120.1 | 21.4 |
| SGAE [38] | - | 38.4 | 28.4 | 58.6 | 127.8 | 22.1 |
| AoANet [12] | - | 38.9 | 29.2 | 58.8 | 129.8 | 22.4 |
| M2 [7] | - | 39.1 | 29.2 | 58.6 | 131.2 | 22.6 |
| X-LAN [26] | - | 39.5 | 29.5 | 59.2 | 132.0 | 23.4 |
| RSTNet [42] | - | 40.1 | 29.8 | *59.5* | 135.6 | 23.3 |
| OSCAR$_b$* [22] | *198M* | 40.5 | 29.7 | - | 137.6 | 22.8 |
| VinVL$_b$* [41] | 260M | *40.9* | *30.9* | - | *140.4* | *25.1* |
| Detector Free Models | | | | | | |
| CLIP-CAP [25] | - | 33.5 | 27.5 | - | 113.1 | 21.05 |
| BLIP* [21] | - | **40.4** | - | - | 136.7 | - |
| ViTCap [10] | 190M | 40.1 | 29.4 | 59.4 | 133.1 | 23.0 |
| SuperCap CLIP (Ours) | **45M** | 39.2 | 29.2 | 59.1 | 132.4 | 22.8 |
| SuperCap BLIP (Ours) | 93M | 40.0 | **29.8** | **59.7** | **136.9** | **23.6** |



Ours: *a bedroom with a bed and a chair*
GT: *A view of a bed, chair, and window treatments.*

Ours: *a man riding a wave on a surfboard in the water*
GT: *A man riding a wave on top of a surfboard.*

Ours: *a cake with a train on top of a table*
GT: *A cake with icing and cartoon cake toppers.*

Ours: *a group of people on the beach with surfboards*
GT: *Four kids sitting on surfboards with man in background.*

Ours: *a man riding a wave on a surfboard in the water*
GT: *Two people surfing a beautiful sky blue wave.*

Ours: *a person holding a pizza in front of a pan*
GT: *A pizza is being cut into on a sheet.*

Figure 3. **Test Set Performance.** Images were randomly selected from the test set and captioned with the best performing model settings using beam width 5. Best viewed in colour.

meaning that it is easily scalable to account for advancements in these large models. We hope that our work demonstrates to the community that irregular groups of homogeneous pixels can be used instead of the traditional object bounding boxes or patches. Future work will investigate alternative ways to handle multiple resolutions.

# References

[1] Radhakrishna Achanta, Appu Shaji, Kevin Smith, Aurelien Lucchi, Pascal Fua, and Sabine Süsstrunk. Slic superpixels, 2010. 3, 4, 5, 7

[2] Radhakrishna Achanta, Appu Shaji, Kevin Smith, Aurelien Lucchi, Pascal Fua, and Sabine Süsstrunk. Slic superpixels compared to state-of-the-art superpixel methods. *IEEE transactions on pattern analysis and machine intelligence*, 34(11):2274–2282, 2012. 4

[3] Peter Anderson, Basura Fernando, Mark Johnson, and Stephen Gould. Spice: Semantic propositional image caption evaluation. In *European conference on computer vision*, pages 382–398, 2016. 6

[4] 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 *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6077–6086, 2018. 1, 2, 3, 4, 8

[5] Satanjeev Banerjee and Alon Lavie. Meteor: An automatic metric for mt evaluation with improved correlation with human judgments. In *Proceedings of the acl workshop on intrinsic and extrinsic evaluation measures for machine translation and/or summarization*, pages 65–72, 2005. 6

[6] Xinlei Chen, Hao Fang, Tsung-Yi Lin, Ramakrishna Vedantam, Saurabh Gupta, Piotr Doll'ar, and C Lawrence Zitnick. Microsoft coco captions: Data collection and evaluation server. *arXiv preprint arXiv:1504.00325*, 2015. 2, 3, 5, 6, 7, 8

[7] Marcella Cornia, Matteo Stefanini, Lorenzo Baraldi, and Rita Cucchiara. Meshed-memory transformer for image captioning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10578–10587, 2020. 1, 3, 6, 8

[8] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018. 3

[9] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020. 2, 3, 7

[10] Zhiyuan Fang, Jianfeng Wang, Xiaowei Hu, Lin Liang, Zhe Gan, Lijuan Wang, Yezhou Yang, and Zicheng Liu. Injecting semantic concepts into end-to-end image captioning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 18009–18019, 2022. 2, 3, 4, 6, 7, 8

[11] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 2

[12] Lun Huang, Wenmin Wang, Jie Chen, and Xiao-Yong Wei. Attention on attention for image captioning. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 4634–4643, 2019. 8

[13] Xiaoke Huang, Jianfeng Wang, Yansong Tang, Zheng Zhang, Han Hu, Jiwen Lu, Lijuan Wang, and Zicheng Liu. Segment and caption anything. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 13405–13417, 2024. 3

[14] Wenhao Jiang, Lin Ma, Yu-Gang Jiang, Wei Liu, and Tong Zhang. Recurrent fusion network for image captioning. In *Proceedings of the European conference on computer vision (ECCV)*, pages 499–515, 2018. 8

[15] Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361*, 2020. 2

[16] Andrej Karpathy and Li Fei-Fei. Deep visual-semantic alignments for generating image descriptions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3128–3137, 2015. 1, 5, 6, 7, 8

[17] Diederik P Kingma and Jimmy Lei Ba. Adam: Amethod for stochastic optimization. In *Proc. 3rd Int. Conf. Learn. Representations*, pages 1–15, 2014. 6

[18] Thomas N Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. *arXiv preprint arXiv:1609.02907*, 2016. 2

[19] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 4015–4026, 2023. 3

[20] Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A Shamma, et al. Visual genome: Connecting language and vision using crowdsourced dense image annotations. *International journal of computer vision*, 123:32–73, 2017. 2

[21] Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In *International Conference on Machine Learning*, pages 12888–12900. PMLR, 2022. 3, 4, 5, 7, 8

[22] Xiujun Li, Xi Yin, Chunyuan Li, Pengchuan Zhang, Xiaowei Hu, Lei Zhang, Lijuan Wang, Houdong Hu, Li Dong, Furu Wei, et al. Oscar: Object-semantics aligned pre-training for vision-language tasks. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXX 16*, pages 121–137. Springer, 2020. 8

[23] Chin-Yew Lin. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*, pages 74–81, 2004. 6

[24] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Doll'ar, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European conference on computer vision*, pages 740–755, 2014. 2, 3, 5, 6, 7, 8

[25] Ron Mokady, Amir Hertz, and Amit H Bermano. Clipcap: Clip prefix for image captioning. *arXiv preprint arXiv:2111.09734*, 2021. 8

[26] Yingwei Pan, Ting Yao, Yehao Li, and Tao Mei. X-linear attention networks for image captioning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10971–10980, 2020. 8

[27] Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, pages 311–318, 2002. 6

[28] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021. 4, 5, 7

[29] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. *Advances in neural information processing systems*, 28, 2015. 2

[30] Steven J Rennie, Etienne Marcheret, Youssef Mroueh, Jerret Ross, and Vaibhava Goel. Self-critical sequence training for image captioning. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7008–7024, 2017. 6, 7

[31] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. Laion-5b: An open large-scale dataset for training next generation image-text models. *Advances in neural information processing systems*, 35:25278–25294, 2022. 1

[32] Henry Senior, Gregory Slabaugh, Shanxin Yuan, and Luca Rossi. Graph neural networks in vision-language image understanding: A survey. *arXiv preprint arXiv:2303.03761*, 2023. 2, 3

[33] Matteo Stefanini, Marcella Cornia, Lorenzo Baraldi, Silvia Cascianelli, Giuseppe Fiameni, and Rita Cucchiara. From show to tell: A survey on deep learning-based image captioning. *IEEE transactions on pattern analysis and machine intelligence*, 45(1):539–559, 2022. 3

[34] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017. 1, 3, 4, 6

[35] Ramakrishna Vedantam, C Lawrence Zitnick, and Devi Parikh. Cider: Consensus-based image description evaluation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4566–4575, 2015. 6

[36] Jianfeng Wang, Zhengyuan Yang, Xiaowei Hu, Linjie Li, Kevin Lin, Zhe Gan, Zicheng Liu, Ce Liu, and Lijuan Wang. Git: A generative image-to-text transformer for vision and language. *arXiv preprint arXiv:2205.14100*, 2022. 2, 3

[37] Haiyang Xu, Ming Yan, Chenliang Li, Bin Bi, Songfang Huang, Wenming Xiao, and Fei Huang. E2e-vlp: end-to-end vision-language pre-training enhanced by visual learning. *arXiv preprint arXiv:2106.01804*, 2021. 3

[38] Xu Yang, Kaihua Tang, Hanwang Zhang, and Jianfei Cai. Auto-encoding scene graphs for image captioning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10685–10694, 2019. 8

[39] Xu Yang, Jiawei Peng, Zihua Wang, Haiyang Xu, Qinghao Ye, Chenliang Li, Songfang Huang, Fei Huang, Zhangzikang Li, and Yu Zhang. Transforming visual scene graphs to image captions. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 12427–12440. Association for Computational Linguistics, 2023. 7

[40] Ting Yao, Yingwei Pan, Yehao Li, and Tao Mei. Exploring visual relationship for image captioning. In *Proceedings of the European conference on computer vision (ECCV)*, pages 684–699, 2018. 2, 3, 6

[41] Pengchuan Zhang, Xiujun Li, Xiaowei Hu, Jianwei Yang, Lei Zhang, Lijuan Wang, Yejin Choi, and Jianfeng Gao. Vinvl: Revisiting visual representations in vision-language models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5579–5588, 2021. 8

[42] Xuying Zhang, Xiaoshuai Sun, Yunpeng Luo, Jiayi Ji, Yiyi Zhou, Yongjian Wu, Feiyue Huang, and Rongrong Ji. Rstnet: Captioning with adaptive attention on visual and non-visual words. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 15465–15474, 2021. 8

[43] Zhong-Qiu Zhao, Peng Zheng, Shou-tao Xu, and Xindong Wu. Object detection with deep learning: A review. *IEEE transactions on neural networks and learning systems*, 30 (11):3212–3232, 2019. 1, 2, 6

[44] Yiwu Zhong, Liwei Wang, Jianshu Chen, Dong Yu, and Yin Li. Comprehensive image captioning via scene graph decomposition. In *European Conference on Computer Vision*, pages 211–229, 2020. 1, 2, 3, 6

# SuperCap: Multi-resolution Superpixel-based Image Captioning

## Supplementary Material

### 7. Further Qualitative Results

Ours: *a giraffe walking in the grass in a field*
GT: *a giraffe walks on the tundra tree-lined park*

Ours: *two giraffes and a giraffe laying in a field*
GT: *Three giraffes are sitting on the ground*

Ours: *a young boy sitting on a bench in a park*
GT: *A young man sitting on a park bench next to a playground*

Ours: *a green bus driving down a street*
GT: *a green bus is driving down the street*

Ours: *a wooden bench sitting on top of a park*
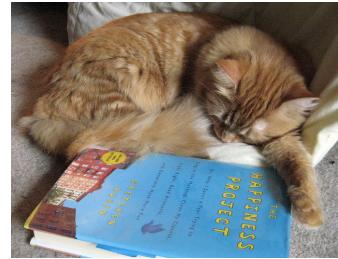GT: *A closeup of an old and mossy outdoor garden bench*

Ours: *a group of birds eating from a nest*
GT: *A mother bird feeds her little baby chicks*

Ours: *a black bird perched on top of a tree branch*
GT: *A couple of black birds are standing on a tree branch*

Ours: *a white bird flying over a body of water*
GT: *A white and gray bird soaring over the blue ocean*

Ours: *a cat laying on a bed next to a book*
GT: *A cat sleeping on a pillow next to a book*

Ours: *a stop sign on the side of a street*
GT: *A stop sign is viewed from a low angle*

Ours: *a cat and a dog sitting in the grass*
GT: *A dog standing next to a cat in a dirt field*

Ours: *a cat sitting on top of a chair with a tie*
GT: *A kitten peeking out from behind a tablecloth on a chair*

Figure 4. **Qualitative Performance.** Example captions from the Karpathy test set. Best viewed in colour.

2