# Image Captioning using Convolutional Neural Network and Long Short Term Memory Networks on COCO Dataset

**Jingyuan He**       **Shuyu Wang**       **Sohail Lokhandwala**

University of California San Diego
La Jolla, CA, 92092

*j7he@ucsd.edu*       *shw043@ucsd.edu*       *slokhandwala@ucsd.edu*

## Abstract

Captions for images from the COCO 2015 dataset are generated through an encoder-decoder structure. The encoder is a multi-layer convolutional neural network (CNN) or a ResNet-50 CNN pretrained on ImageNet that extracts the feature value from the images, and the decoder is a long short term memory (LSTM) unit architecture that is trained to generate a caption for each image by words based on the given vocabulary. Teacher forcing is used throughout the training stage and cross entropy loss of the generated caption encoded by the IDs of the words in the dictionary is computed. To evaluate the actual generated captions, individual BLEU-1 and BLEU-4 scores are computed to compute similarity of predictions to corresponding ground truth labels. The best custom CNN encoder-decoder model achieved a 42.47% BLEU-1 score and a 2.01% BLEU-4 score while the best overall model used ResNet-50 and achieved a performance of 58.04% on BLEU-1 and 3.54% on BLEU-4.

## 1    Introduction

In order to extract and translate the main information contained in an image to text, we implemented an encoder-decoder model built on convolutional neural networks and long short term memory networks. First, we need to process the images such that they are all of the standard shapes and could be operated and interpreted by machines. Thus, We convoluted the images to match the size of the LSTM and to filter features within the images such that it could be trained to predict the start word of the image caption. An alternative encoder architecture is the pretrained ResNet-50 from ImageNet. By feeding in the true captions into the LSTM at each time step, the model should learn to generate captions for an encoded image word by word. To compare the captions our model generated, we compute the loss over the captions and their BLEU scores against the set of original captions we have. The captions generated deterministically and stochastically with different temperatures are analyzed.

## 2    Related work

In the project, we are using the COCO dataset: Common Object in Context. It is frequently used for recognition in context as well as object segmentation because there are multiple human imputed and evaluated captions for each of the images in the dataset.

The first model that have a comprehensive view over instance object segmentation is proposed by Ronghang Hu1, Piotr Dollar, Kaiming He, Trevor Darrell, and Ross Girshick in 2018, with a Mask R-CNN model that segments the 3000 visual concepts from the Visual Genome dataset and masks

anotations from the 80 classes of COCO dataset. They adopt a transfer learning approach on the R-CNN model, which decomposes the segmentation problem in subprograms: object detection and mask prediction [1].

The multi-label feature of the COCO dataset also stimulates multi-label image classification work. Yucheng Li, Yale Song, and Jiebo Luo built proposed a new loss function and a label decision module to improving the pairwise ranking based on multi-label image classification, tested it on the VOC2007, NUS-WIDE and the MS-COCO dataset [2].

## 3        Baseline Models

In the task, we use convolutional neural networks to encode the images input and then decode it with the long short term memory architecture to produce the start word. From this start word, we continue to generate the captions word by word through the LSTM architecture.

### 3.1        Convolutional neural networks

Convolutional neural networks are useful for classifying features from images. There are three types of layers that build the CNN model which are convolutional, pooling, and fully-connected layers. The convolution layers can help extract features from the inputs with the filters. The pooling layers decrease the size of the features map and introduce translational invariance. In this project, we used two types of pooling layers: max and average pooling that get the maximum value from the feature map and the average value of the element in the area of the filter. Both of the layers can help to preserve the important features. The last layers are fully connected layers, They work with the inputs as flattened inputs, and classify the detected features into class labels.

Table 1: Convolutional Network Baseline Architecture

| Layer \ Parameter | Input Channel | Output Channel | Stride Size | Kernel Size | Activation Function | Padding Size |
|---|---|---|---|---|---|---|
| conv1 | 3 | 64 | 4 | 11 | ReLU | 0 |
| maxpool1 | 64 | 64 | 2 | 3 | ReLU | 0 |
| conv2 | 64 | 128 | 1 | 5 | ReLU | 2 |
| maxpool2 | 128 | 128 | 2 | 3 | ReLU | 0 |
| conv3 | 128 | 256 | 1 | 3 | ReLU | 1 |
| conv4 | 256 | 256 | 1 | 3 | ReLU | 1 |
| conv5 | 256 | 128 | 1 | 3 | ReLU | 1 |
| maxpool2 | 128 | 128 | 2 | 3 | ReLU | 0 |
| adaptive_ avgpool | 128 | 128 | 0 | 6 | None | 0 |
| fc1 | 128 | 1024 | 0 | 0 | ReLU | 0 |
| fc2 | 1024 | 1024 | 0 | 0 | ReLU | 0 |
| fc3 | 1024 | 300 | 0 | 0 | None | 0 |

## 3.2 Long short term memory

LSTM is a recurrent neural network unit that feeds the output from the last step into the current input. It has four gates: input, input modulation, output, and forget. With these gates, the model is able to keep important information and discard the insignificant data during training. In the case of caption generation, we are predicting the next word based on the current output.

During training, validating, and testing, we did teacher forcing on the LSTM model. Teacher forcing feeds the actual label into the LSTM model after we compute the loss of our predicted word with this label word. This prevents inaccurate predictions by the model from affecting the next word prediction and largely skewing loss values.

Additionally, inputs to the LSTM units are embedded into representations of words or images to allow the LSTM to begin its prediction on the image and use the following predicted words to predict the next time step of words.

Our default configuration file uses 2 LSTM unit layers of 512 units each, with an input embedding size of 300.

Table 2: LSTM Baseline Structure

| Hyperparameter | Value |
|---|---|
| Hidden Units | 512 |
| Embedding Units | 300 |
| Number of Layers | 2 |
| Batch-First | True |

### 3.2.3 Deterministic and Stochastic Caption Generation

There are two approaches to selecting the next word in a sequence when generating captions: deterministically and stochastically. Deterministically generating captions involves selecting the word with the maximum value output. Stochastically generating samples involves creating a probability distribution of the outputs and sampling from this distribution. To provide more flexibility in creating this sampling distribution, we use the technique of using a weighted softmax on the outputs.

The formula for this is provided below.

$$y^j = exp(o^j/\tau) / \sum_n exp(o^n/\tau)$$

Equation 1: Weighted Softmax

where $y$ is the generated distribution, $o$ is the output from the previous layer, $n$ is the size of the vocabulary and $\tau$ is the "temperature" term.

Our default configuration file uses a stochastic sampling method with a default temperature of 0.4.

## 3.3 ResNet-50 model

ResNet-50 is a 50-layer neural network developed by Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. They made use of residual networks to enable effective optimization after evaluating residual nets up to 152 layers deep. The result improves the object detection on the

COCO dataset by 28%. Its performance also helped the team win the ImageNet detection, localization, and COCO segmentation [3].

We are using the ResNet-50 model pre-trained on ImageNet, while freezing all the layers except for the last one and changing the last layer to be a fully-connected layer (nn.Linear) from PyTorch.

Our default configuration file for a ResNet encoder uses this pre-trained version of ResNet-50 from ImageNet.

### 3.4     Default Hyperparameters

The set of default hyperparameters that we used when initially training both the Custom CNN and ResNet-50 are given below.

Table 3: Default Hyperparameters

| Hyperparameter | Value |
|---|---|
| Epochs | 10 |
| Learning Rate | 0.0005 |
| Optimizer | Adam [4] |
| Weight Decay | 0.0001 |
| Scheduler | StepLR |
| Step Size | 3 |
| Early Stopping | True |
| Patience Epochs | 3 |

## 4     Encoder-decoder model architecture modifications

Several architectural changes were made to the CNN to improve the performance of the model: we aim to get a higher BLEU score and lower loss. The intuition and hypothesis of the performance of these changes are discussed in the sections below.

### 4.1     Encoder CNN modified model

#### 4.1.1     Dropout

Dropout is a commonly used regularization method that helps the model generalize. By adding a dropout layer at the last fully connected layer of the CNN model, we expect to downgrade over-fitting and to see a performance on the test set that is close to the one we have on the training data.

#### 4.1.2     Embedding size

Embedding size constraints the number of orthogonally unique vocabulary representations and features we can represent the images by. We increased the embedding size from 300 to 512, so we expect the network to have encoded more information about the image in the embedding and generate better captions. That is, the model should be trained to recognize more specific objects

and the substantial differences between the similar objects. In this case, there should be an improvement on both BLEU-1 and BLEU-4 scores.

### 4.1.3 Hypothesis

According to each of these claims, we hypothesize that the model has a higher BLEU1 and BLEU4 score on generating captions for the testing data than the baseline model, and has a better performance on providing the correct caption for similar objects respectively, with a 512 embedding size and dropout units on the last fully connected layer.

## 4.2 Decoder LSTM modified model

### 4.2.1 Layers

Adding the number of layers of the LSTM model can help us to increase the depth of the neural network and thereby allow the network to learn more complex features about the relations between words. We increased the layers from 2 to 3.

### 4.2.2 Hidden units

Adding more hidden units to each of the layers in the LSTM can help to increase the width of the neural network by making more varieties of relationships among the inputs. We experiment to halve the hidden units to be 256, which should downgrade the performance slightly.

### 4.2.3 Hypothesis

According to each of these claims, we hypothesize that the model has a similar BLEU1 and BLEU4 score on generating captions for the testing data as the baseline model, with a 3 layers LSTM and 256 number of hidden units.

## 5 ResNet-50 model hyperparameter modification

Several architectural changes were made to the ResNet-50 model to improve the performance of the model: we aim to get a higher BLEU score. The intuition and hypothesis of the performance of these changes are discussed in the sections below.

## 5.1 Weight decay

Weight decay improves generalization of the model and prevents overfitting. Thus, we want to increase the decay from $10^{-4}$ to $10^{-2}$. This smaller decay in the learning rate should allow for faster learning in earlier epochs.

## 5.2 Step size

The model will update the learning rate by multiplying in the decay every (step size) epochs. The default step size was set to 3. To update the learning rate more precisely with the modified decay, we wanted to decay more often. Thus, we changed the step size parameter to 2 and expected a more frequent decay on the learning rate.

## 5.3 Learning Rate

Since we decay the learning rate more frequently with a larger amount, we should increase the initial learning rate as well such that it doesn't converge too quickly in the wrong direction. Therefore, we increased the learning rate from $5\times10^{-4}$ to $5\times10^{-3}$.

## 5.4 Number of epoch

As we are using decay with step size, we want to train for enough epoch such that these

modifications took effect. The default number of epochs — 10 — is a small number such that the model may not decay the learning rate to a small enough value to properly converge. Therefore, we increased the number of epochs to 20.

### 5.5    Hypothesis

According to each of these claims, we hypothesize that the model has a higher BLEU1 and BLEU4 score on generating captions for the testing data than the baseline model, and has a better performance on providing the correct caption for similar objects respectively, with a $10^{-2}$ decay, a step size of 2, a $5\times10^{-3}$ learning rate, and 20 number of epoch during training.

## 6    Results

### 6.1    BLEU performance

BLEU-Score, the Bilingual Evaluation Understudy, is a standardized method to evaluate the performance of machine translation. The individual BLEU-1 score measures the match of single words whereas the individual BLEU-4 score measures the match of 4 consecutive words. The performance of our experimented models on BLEU-1 and BLEU-4 scores are presented below.

Table  4: BLEU Score Performance Table

| Model | BLEU-1 (%) | BLEU-4 (%) |
|---|---|---|
| Custom Encoder-Decoder: default (3.1-3.2) | 41.70 | 1.47 |
| Custom Encoder-Decoder: 3 layer LSTM (4.2) | 47.66 | 1.46 |
| Custom Encoder-Decoder: dropout and embed (4.1) | 42.67 | 2.01 |
| ResNet-50: default (3.3, 3.2) | 58.04 | 3.54 |
| ResNet-50: LR decay changes (5.1-5.4) | 39.42 | 1.48 |

### 6.2    Best models

#### 6.2.1    Best Custom CNN Encoder-Decoder model

The best Custom CNN Encoder-Decoder model was the Custom Encoder-Decoder: dropout and embed model proposed in section 4.1.
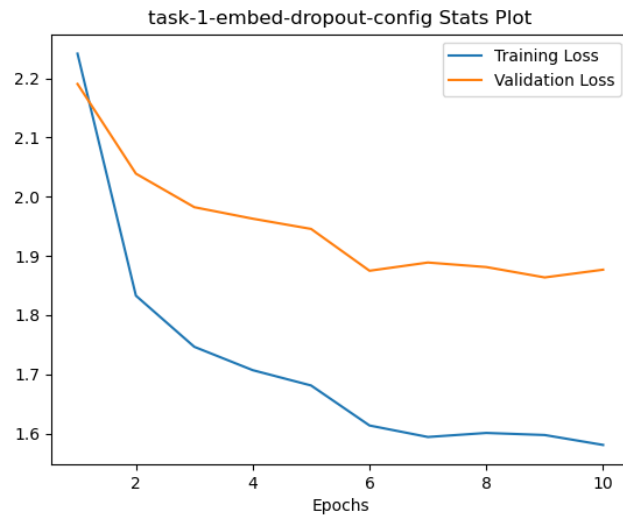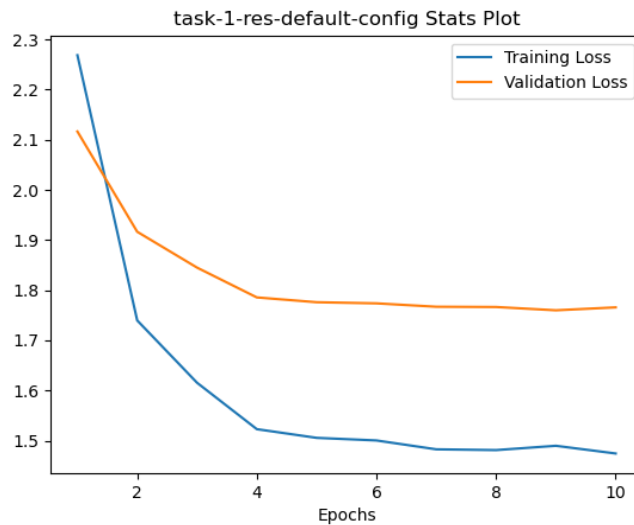
Figure 1: Training and Validation Loss of Best Custom Encoder-Decoder Model

The test loss for this model was 1.8798. It had an average BLEU-1 score of 42.67% and BLEU-4 score of 2.01% across the test set.

The validation and training loss in this model both still appear to be trending downwards. As such, the model would likely have achieved better performance if trained for more epochs. It is also possible that the learning rate decayed too fast and weight updates were trivial or near zero. As such, the model would have effectively stopped learning.

### 6.2.2 Best ResNet-50 model

The best ResNet-50 model was the first one trained using default hyperparameters and architecture.



Figure 2: Training and Validation Loss of Best ResNet-50 Model

The test loss for this model was 1.7617. It had an average BLEU-1 score of 58.04% and BLEU-4

score of 3.54% across the test set.

Unlike the loss plot of the Best Custom CNN model (Fig. 1), the loss values of this model appear to have plateaued in the provided graph (Fig. 2). As such, the previous model would have likely benefited from training more epochs whereas this model was trained for what appears to be a sufficient number of epochs given the other hyperparameters.

### 6.3 Discussion

Both of the custom CNN Encoder-Decoder models had a nice performance. The Encoder-Decoder model with modified LSTM structure has the highest BLEU-1 score. Considering the slightly lower BLEU-4 score, our hypothesis is correct. We can also say that the model actually has a better performance, since the decrease in BLEU-4 score is trivially small. The other model with modified CNN structure has the highest BLEU-4 score and a slightly improved BLEU-1 score. Such performance confirmed the hypotheses we made in section 4.1.3 and 4.2.3.

The modified ResNet-50 model has a bad performance on its BLEU scores. This contradicts our hypothesis that the modified model will have a better performance compared to the default ResNet-50 model. We suspect that the reason for this bad performance is that we are decaying the weight too fast and too often. We should aim for a lower weight decay such that we adjust the weights slower and more consistently to prevent reducing the learning rate down to a trivial (nearly zero) value too soon.

Across the custom CNN model and the ResNet-50 models, we can observe that the ResNet-50 default has the best performance, which is expected because it is a carefully developed model. The reason why the modified ResNet-50 model had the worst performance might be because the weight decay changes explained above. As all the layers in this deep network except for the last fully connected layer are not trainable, the increased weight decay might not be compatible with this architecture.

## 7 Captions

The captions generated by the best CNN model and the best ResNet-50 model are presented in the sections below. For each model, 3 images with the best generated captions and 3 images with the worst generated captions are presented. For each of these images, we provided the caption generated with temperatures of 0.4 by default, the captions generated deterministically, the ones generated with a very low temperature of 0.001, and the captions resulted from a very high temperature of 5.

### 7.1 Best encoder-decoder model: dropout and embed

For the dropout and embedding modified Custom CNN encoder-decoder model, which was our best performing custom CNN model, the following are some examples of both good and bad predictions that the network made using deterministic and stochastic sampling techniques.

Table 5: Best Encoder-Decoder Model's Best and Worst Predictions

| Best 3 Predictions: | Worst 3 Predictions: |
|---|---|
|  |  |
| Figure 3: COCO_val2014_000000393056.jpg | Figure 4: COCO_val2014_000000479762.jpg |

**Best 3 Predictions:**



Figure 3: COCO_val2014_000000393056.jpg

**Ground Truth Captions:**
a man wearing a wet suit riding the wave
A person riding a wave on top of a surfboard.
A surfer rides a small wave in the ocean.
a man riding a wave on a blue ocean
A young man ridding a small wave on a surfboard.

**Temp = 0.4 (default):**
a man riding a wave on a surfboard.

**Deterministic:**
a man riding a wave on a surfboard.

**Temp = 0.001:**
a man riding a wave on a surfboard.

**Temp = 5:**
exposure find status jester panorama extra wheelchairs fixtures issue early tabl stop hitch eaten showed tay rotting platforms sock determined

**Worst 3 Predictions:**



Figure 4: COCO_val2014_000000479762.jpg

**Ground Truth Captions:**
Three adult zebra are walking through the dirt path.
A trio of dwarf zebras wander their pen.
A pack of zebras running and playing in the open ,
There are baby zebras standing together outside .
A family of zebras walking together in the dirt.

**Temp = 0.4 (default):**
a man riding skiing a horse on a wave

**Deterministic:**
a man is riding a skateboard on a ramp.

**Temp = 0.001:**
a man is riding a skateboard on a ramp.

**Temp = 5:**
customized magenta scape crotch scissor outs sow fix choosing mid-air mid-swing fitted drapes performer joint style equipped march salads beet

Figure 5: COCO_val2014_000000223289.jpg

Figure 6: COCO_val2014_000000385103.jpg

**Ground Truth Captions:**
Close up of a plate with food on it.
A piece of fish on a sandwich next to a lemon.
This is a fish sandwich on a bun with a slice of lemon on the side of the plate.
A sandwich that is on a plate with the top piece of bread off.
A plate that has bread and chicken on it.

**Temp = 0.4 (default):**
a plate of food with a sandwich and it.

**Deterministic:**
a plate of food with a sandwich and a fork.

**Temp = 0.001:**
a plate of food with a sandwich and a fork.

**Temp = 5:**
11:30 salami index clothes pantry walled boxer hyena cubicles piping scissors laptops overhand colorized sharpies necklaces sparkler drifting blackberries potato

**Ground Truth Captions:**
A room filled with different types of items all around .

All white kitchen with brown counter tops and red fire extinguisher .
A kitchen with cream colored walls and brown counters.
A modern kitchen has an abundance of counter space.
THIS IS A PICTURE OF A LARGE CLEAN KITCHEN

**Temp = 0.4 (default):**
a man sitting holding a teddy teddy dog

**Deterministic:**
a man is holding a cell phone in his hand.

**Temp = 0.001:**
a man is holding a banana in his mouth.

**Temp = 5:**
components lighter trey dill paved body main flip furred gathered to muggy escapes signs coin expose souvenir bouquets rode stony

Figure 7: COCO_val2014_000000278095.jpg

**Ground Truth Captions:**
A man power sliding on a long board
A young man sitting on his skateboard touching the ground
A man riding a skateboard down a street.
Young man with skateboard appearing like he just fell down.
A man doing a trick on a skateboard in the middle of the street.

**Temp = 0.4 (default):**
a man standing a skateboard riding a park lot.

**Deterministic:**
a man is riding a skateboard on a ramp.

**Temp = 0.001:**
a man is riding a skateboard on a ramp.

**Temp = 5:**
wineglasses cooks snowshoes albeit really groupe heavy expose score chid handcuffs modular drenched bib demands obstacles cranberry circuit customers cookers



Figure 8: COCO_val2014_000000123289.jpg

**Ground Truth Captions:**
A crowded city street filled with traffic at night.
A traffic light over many different passing cars.
A night time city view with vehicle lights and street lights .
A city street with lights at night with vehicles
A busy street with passing traffic at night

**Temp = 0.4 (default):**
a black that standing in the grass next a frisbee in

**Deterministic:**
a man is holding a tennis racket on a tennis court.

**Temp = 0.001:**
a man is holding a tennis racket on a tennis court.

**Temp = 5:**
outstanding carrots golf collectors mouthful jams delicatessen fit sand zombie traipsing inner greenwich eyeballs whizzes scoop blanketed salami strong machete

## 7.2 Best ResNet-50 model : default

For the default ResNet-50 model, which was our best performing ResNet-50 model, the following are some examples of both good and bad predictions that the network made using deterministic and stochastic sampling techniques.

Table 6: Best ResNet-50 Model's Best and Worst Predictions

| Best 3 Predictions: | Worst 3 Predictions: |
|---|---|
|  Figure 9: COCO_val2014_000000075285.jpg |  Figure 10: COCO_val2014_000000371879.jpg |

**Best 3 Predictions:**



Figure 9: COCO_val2014_000000075285.jpg

**Ground Truth Captions:**
A desk with a keyboard , mouse and computer monitor.
An empty desk chair pushed up to a small computer desk.
A keyboard , mouse and computer monitor on a desk with a laptop.
A table with a computer and desk chair along a wall under a window.
A computer desk that also has a laptop on it.

**Temp = 0.4 (default):**
a table with a computer and a laptop monitor on it.

**Deterministic:**
a desk with a laptop and a laptop on it

**Temp = 0.001:**
a desk with a laptop and a laptop on it

**Temp = 5:**
completed average certain warplane telephone stencils shown lamp beef round bar ponytail majestic pain grocery braces bathrooms session initial pre

**Worst 3 Predictions:**



Figure 10:
COCO_val2014_000000371879.jpg

**Ground Truth Captions:**
A collaboration  of people in different pictures doing things
A series of images of young men painting and holding kites.
Various children and adults are making their own kites.
A variety of colorful pictures with people doing various activities.
A few people working with colored fabrics in different ways.

**Temp = 0.4 (default):**
a woman carrying a teddy animal stuffed bears

**Deterministic:**
a woman is holding a teddy bear in a park.

**Temp = 0.001:**
a woman is holding a teddy bear in a park.

**Temp = 5:**
buggy ballons toy accents seal bunting sprawled roundabout war christ forging expressing mallet gorgeous trot spread like meant mein ranch

Figure 11:
COCO_val2014_000000082715.jpg

**Ground Truth Captions:**
A person that is surfing in the water.
A man on a surfboard surfing a wave in the ocean.
A man riding a wave on a surfboard.
a surfer in a white shirt surfing on a sunny day
A man on a surfboard riding an ocean wave.

**Temp = 0.4 (default):**
a man riding a wave on the ocean.

**Deterministic:**
a man riding a wave on a surfboard.

**Temp = 0.001:**
a man riding a wave on a surfboard.

**Temp = 5:**
bathtubs crumbling rotten flag acrobatically collapse cashews multi-colored dc fliers stretching tourists napkins band beverages and/or soaring shoulder obstructed brighten



Figure 12:
COCO_val2014_000000469088.jpg

**Ground Truth Captions:**
Two dogs that are standing up holding each other .
Two dogs on their hind legs playing with each other .
Two dogs standing on their back legs wrestling with one another.
Two little , well groomed dogs hugging each other energetically
A couple of dogs standing each other up being playful .

**Temp = 0.4 (default):**
a dog is is sitting next a bench

**Deterministic:**
a dog is sitting on a bench next to a dog.

**Temp = 0.001:**
a dog is sitting on a bench with a dog.

**Temp = 5:**
bonnet adrift cross-county fill meets chowder roadside issue logo-emblazoned petrol twilight papayas ripened ashore fluids pecking design girafee dismantled deaker



Figure 13:
COCO_val2014_000000063950.jpg



Figure 14:
COCO_val2014_000000126606.jpg

| Ground Truth Captions: | Ground Truth Captions: |
|---|---|
| A black and white cat is sitting in the sink. A black and white cat laying in a bathroom sink. A cat laying in a white sink next to a toilet. Black and White cat lays inside of the sink The cat is relaxing in the bathroom sink. | Group of motorcycle riders looking over traffic on the street Several people wearing jackets with foreign writing and motorcycle helmets. A bang of bikers sitting on the side of a road. A small group of motorcyclists stand next to a motorcycle. The police officers are standing on the side of the road. |
| **Temp = 0.4 (default):** a cat is sitting in a bathroom sink. | **Temp = 0.4 (default):** a man in a skateboard down a man in it |
| **Deterministic:** a cat is sitting on a toilet in a bathroom | **Deterministic:** a man riding a bike down a street. |
| **Temp = 0.001:** a cat is sitting on a toilet in a bathroom. | **Temp = 0.001:** a man riding a motorcycle down a street. |
| **Temp = 5:** bountiful count chests jeweled wooded coop includes toothpick bagels wristwatch spending current clamp rue him frolicking await tech paraphernalia oddly | **Temp = 5:** supports latch flavored pom-poms geared penn any affection build roaring limbs looks marking ext coming part margherita bots reservoir brook |

# 8    Discussions

## 8.1    Temperature

According to the above table, the captions generated under a very low temperature, 0.001 in our case, is almost deterministic. This occurs because the sampling probability of everything but the maximum value is driven to nearly zero when the temperature softmax is applied. This can be seen in Equation 1, where as the temperature approaches low values, the distribution favors only the maximum value. As a result, a probabilistic sample of this distribution results in the highest value in the original distribution nearly every time.

A high temperature of 5 has a low performance, even though it generates rare words in the overall vocabulary. Once again, referring to Equation 1, we can see that as the temperature approaches extremely large values, the resulting distribution becomes nearly uniform and each word in the vocabulary is nearly equally likely to be predicted, making the generated caption nonsensical.

The captions generated with a default temperature of 0.4 tends to give a most precise description of the images. The worst captions generated with this temperature seem to result from overfitting, in which the model is outputting a well-structured yet non-related caption, which is probably what the model learnt from some other images with comparably common objects.

## 8.2    Best encoder-decoder model: dropout and embed

We chose the Encoder-Decoder model with a modified CNN of 512 dropout size and 512 embedding units to be the best custom CNN model. This model has a slightly higher BLEU-1 score (42.67%) and a significantly higher BLEU-4 score (2.01%) compared to the baseline Encoder-Decoder model.

The model is selected to be the best because it has a well BLEU-1 score, the highest BLEU-4 score, and the lowest test loss in our custom CNN Encoder-Decoder model experiments, which means that while correctly recognizing the objects, the model is able to formulate phrases in a

readable manner rather than merely spitting out the names of the objects within the images.

The test loss is very close to the validation loss of the last epoch, confirming that we are not overfitting on the training data. The increase in embedding size could be the reason for the improvement in BLEU-1 score. As the below table indicates, the model is able to distinguish between different categories of food.

Table 7: Best Encoder-Decoder Model's Performance on Similar Objects

| Image of Sandwich | Image of Pizza |
|---|---|
| Figure 15: COCO_val2014_000000223289.jpg<br><br>**Ground Truth Captions:**<br>Close up of a plate with food on it.<br>A piece of fish on a sandwich next to a lemon.<br>This is a fish sandwich on a bun with a slice of lemon on the side of the plate.<br>A sandwich that is on a plate with the top piece of bread off.<br>A plate that has bread and chicken on it.<br><br>**Temp = 0.4 (default):**<br>a plate of food with a sandwich and it.<br>**Deterministic:**<br>a plate of food with a sandwich and a fork. | Figure 16: COCO_val2014_000000174123.jpg<br><br>**Ground Truth Captions:**<br>A huge plate of yummy food with fork to eat .<br>A portion of a pizza is sitting on a tray and someone is holding a fork and a knife.<br>a person holding a knife and fork over a pizza<br>A white plate of pizza on a table.<br>A pizza topped with cheese , tomato sauce , and mushrooms being sliced on a plate.<br><br>**Temp = 0.4 (default):**<br>a pizza of food with a on a table |

Despite the fact that among the three images, the default temperature tends to give the most descriptive caption, the deterministically generated captions seems to perform better on Figure 15: it recognizes the plate, the sandwich, and the tip of the fork on the edge of the image whereas the captions generated with the default temperature of 0.4 only recognizes the plate and the sandwich.

## 8.3 Best ResNet-50 model : default

According to what we learnt from the residual network paper, ResNet-50 improves the performance of object detection on COCO dataset by 28% [3]. The deep layers of the network with different convolutional layers and different activation functions should be carefully selected. Therefore, it is expected that the ResNet-50 default model achieves significantly higher BLEU-1 and BLEU-2 scores compared to the custom CNN models we implemented.

## 8.4 Modified encoder-decoder model: LSTM

The 3-layer LSTM model gives the highest BLEU-1 score, which is 47.66% yet a slightly lower BLEU-2 score of 1.46% compared to the baseline model. The high BLEU-1 score indicates that it is able to recognize more objects in the images. The decrease in BLEU-4 score probably indicates that while it is capable of recognizing subjects, the model cannot put them into structured sentences.

### 8.5    Modified ResNet-50 model: weight decay related modification

The modified ResNet-50 was the worst model when evaluated in terms of BLEU scores and its test losses, which ended at a high value of 3.15312.
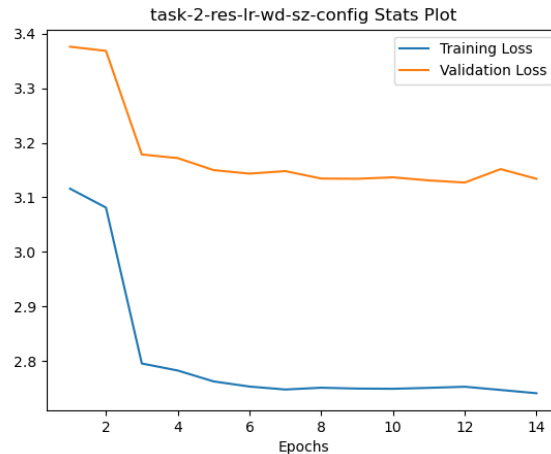


Figure 17: Training and Validation Loss of Best Custom Encoder-Decoder Model

Inspecting the best captions generated by this model (the ones with the highest BLEU scores), the model failed to recognize any objects but only predict very common words like "man" and "a". As discussed above in section 6.3, the reason why this model performs badly is likely that the decay value was too large and the step size caused this decay to occur far too often. This resulted in the model attaining a near zero learning rate early on in its learning, preventing the model from learning important relationships between images and words as well as between individual words.

### 8.6    Best Hyperparameters

Using the default hyperparameters as discussed in section 3 seemed to generally provide highly accurate and functional models. While we did not have the opportunity to individually tune any of these hyperparameters in a controlled setting (where we did not modify other variables), the default values allowed for models with well generalized models that produced largely accurate predictions. As seen in section 8.5, trying to tune a number of these parameters together resulted in degraded model performance.

## 9    Team contribution

Sohail Lokhandwala was primarily responsible for debugging the model during run time, implementing helper functions in the experiment.py file, and debugging some starter code.

Jingyuan He implemented a number of helper functions in experiment.py, and implemented the ResNet-50 changes.

Shuyu Wang implemented and designed architecture changes.

All team members were responsible for collaboratively debugging and determining hyperparameter changes.

## Acknowledgments

## References

[1] Ronghang Hu, Piotr Dollár, Kaiming He, & Trevor Darrell, Ross Girshick. Learning to Segment Every Thing. Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2018, pp. 4233-4241.

[2] Yuncheng Li, Yale Song, Jiebo Luo. Improving Pairwise Ranking for Multi-label Image Classification. Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2017, pp. 3617-3625.

[3] K. He, X. Zhang, S. Ren and J. Sun, "Deep Residual Learning for Image Recognition," 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2016, pp. 770-778, doi: 10.1109/CVPR.2016.90.

[4] Kingma D, Ba J. Adam: A method for stochastic optimization. arXiv preprint arXiv:14126980. 2014