# CNN LSTM Encoder Decoder Model for Caption Generation on COCO Dataset

**Jason Stanley**
Department of Electrical and Computer Engineering
University of California, San Diego
La Jolla, CA 92093
jtstanle@ucsd.edu


**Elliot Lee**
Department of Computer Science and Engineering
University of California, San Diego
La Jolla, CA 92093
eal001@ucsd.edu


**Michael Ma**
Department of Electrical and Computer Engineering
University of California, San Diego
La Jolla, CA 92093
jum002@ucsd.edu


**Annabella Macaluso**
Department of Electrical and Computer Engineering
University of California, San Diego
La Jolla, CA 92093
amacalus@ucsd.edu

## Abstract

In our paper we propose different methods and findings for visual-learning on the Microsoft COCO dataset using Convolutional Neural Networks (CNNs) and Long Short-Term Memory (LSTM) to create a Encoder Decoder architecture. We also compare our model with a popular alternative, the ResNet-50 model (Residual Network 50) by Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun.

After trying many types of CNN-LSTM networks, with changes to CNN architecture, LSTM architecture, hyper parameters, and using pre-trained models, that using a strong image classifier like Resnet or using a large LSTM network both resulted in good image captioning performance. In task 1 we had our best network be a Custom CNN-LSTM with an larger hidden layer, which achieved Bleu-1 and Bleu-4 scores of 47.36 and 1.79 respectively on the holdout test set. For task 2, we saw our resnet-LSTM network achieved Bleu-1 and Bleu-4 scores of 47.81 and 1.70 respectively on the holdout test set.

## 1   Introduction

The topic of learning-based vision provides new opportunities and challenges, especially within the area of learning visual representations such as objection detection. What do we mean by objection detection? To provide some intuition, when providing a person a photo, it'd be natural for them to

classify different objects within the image that they recognize. For instance, they may be provided a picture of a birthday party and could identify common objects such as a cake, or people or birthday presents. Or they could be provided an image of a table with different blocks of varying shapes and based on prior knowledge throughout their life could classify their shapes using semantics solely based on visual cues. This type of innate human ability to easily classify objects would present massive benefits in varying fields such as robotics, quality control for medicine and agriculture, security, video production and more. But how can we imbue computers with this same innate human ability?

Through the usage of Convolutional Neural Networks and Long Short-Term Memory, we explore the process of combining semantic language and visual learning in order to perform object detection and captioning on a given image. We generated a CNN LSTM architecture where the CNN takes an input of an image and abstracts the essential details of the image, through feature extraction. These features are passed into an LSTM model to perform sequence prediction and output a caption that describes objects in an image in relation to each other in English. We trained this model using the MS COCO dataset and experimented with different layer sizes and parameters to see how we could create a model that would generalize to the MS COCO dataset well. Finally, we tested using a popular network, ResNet-50 to compare out results with another known option.

Through these steps we hope to provide more insight on the process of object detection and caption generation using CNN LSTM models.

## 2   Background/Related Work

For training our model we used the extensive Common Objects in Context dataset (COCO). This dataset is a large-scale object detection, segmentation and captioning dataset that has millions of object instances, 80 object categories and 5 captions per image. To read more about the Microsoft COCO dataset please see this paper describing the COCO dataset; Microsoft COCO: Common Objects in Context

Similar papers that have used this this dataset inconjunction with their work include

- MDETR - Modulated Detection for End-to-End Multi-Modal Understanding. In this paper by Aishwarya Kamath, Mannat Singh, Yann LeCun, Gabriel Synnaeve and Ishan Misra they propose MDETR, an end-to-end modulated detector that performs object detection using a rawtext query, like a caption or a question. To pre-train their model they used multiple datasets including the MS COCO dataset.
- YOLOv4: Optimal Speed and Accuracy of Object Detection. In this paper by Alexey Bochkovskiy, Chien-Yao Wang and Hong-Yuan Mark Liao, they propose a new CNN that can perform accurate objection detection real-time on a conventional GPU and does not require a large number of GPUs which other modern nerual networks require. They validate the accuracy of the detector on the MS COCO (test-dev 2017) dataset.

In addition, the following lecture slides from Professor Gary Cottrell's Deep Learning class at the University of California, San Diego and supplemental materials were highly utilized for this project.

- Lecture Recurrent Nets Part 1 Page 52 - 55
- Lecture Generative Modeling with RNNs Page 9 - 15
- Lecture Introduction to Convolutional Networks Page 5 - 7
- Lecture Convolutional Networks, Part II Page 3 - 10
- CS231n Convolutional Neural Networks for Visual Recognition

The following PyTorch libraries were important in implementing our model architectures. For more information on how they function please view the following documentation pages

- PyTorch LSTM
- PyTorch BatchNorm2D
- PyTorch Embedding

- PyTorch Linear
- PyTorch CrossEntropyLoss

## 3  Models

This section contains an overview of the CNN, LSTM and ResNet-50 architecture and decisions behind why they're structured like so.

### 3.1  CNN, LSTM and ResNet-50 Architecture Description

The CNN architecture from Task 1 takes an image as input with 3 channels and a scaled size of 256 by 256. The CNN consists of convolutional, pooling and fully connected layers at the end to get a specified embedding size. These layers allows the network to abstract the features from the image to get the core idea of the elements within it. We go into more detail of each layer in the CNN below in this table:

| CNN Architecture | | | | | | |
|---|---|---|---|---|---|---|
| Layer | Input Channels | Output Channels | Stride Size | Kernel Size | Activation Function | Padding Size |
| Convolutional 1 | 3 | 64 | 4 | 11x11 | BN + ReLU | 0 |
| MaxPool 1 | 64 | 64 | 2 | 3x3 | None | 0 |
| Convolutional 2 | 64 | 128 | 1 | 5x5 | BN + ReLU | 2 |
| MaxPool 2 | 128 | 128 | 2 | 3x3 | None | 0 |
| Convolutional 3 | 128 | 256 | 1 | 3x3 | BN + ReLU | 1 |
| Convolutional 4 | 256 | 256 | 1 | 3x3 | BN + ReLU | 1 |
| Convolutional 5 | 256 | 128 | 1 | 3x3 | BN + ReLU | 1 |
| MaxPool 4 | 128 | 128 | 2 | 3x3 | None | 0 |
| Adaptive AvgPool | 128 | 1 | 1 | 1x1 | None | 0 |
| Fully Connected 1 | 128 | 1024 | 1 | N/A | ReLU | N/A |
| Fully Connected 2 | 1024 | 1024 | 1 | N/A | ReLU | N/A |
| Fully Connected 3 | 1024 | 300 (EmbeddedSize) | 1 | N/A | N/A | N/A |

The CNN makes up the Encoder portion of the model, as it encodes the features of the image to then be decoded by the LSTM architecture. The following decisions we made about the LSTM architecture are listed in this table

| LSTM Architecture | | | |
|---|---|---|---|
| Layer | Input Channels | Output Channels | Recurrence |
| Embedding | 1 | 300 | sentence length |
| LSTM 1 | 300 | 512 | sentence-length + 1 |
| LSTM 2 | 512 | 512 | sentence-length + 1 |
| Fully Connected | 512 | 14663 | sentence-length + 1 |

Note: sentence-length ranges from 1 to max-length 20

The LSTM is a variety of reccurent neural networks (RNNs) that solves the vanishing and exploding gradient problem. But it still contains the essence of an RNN where it takes sequential data and uses patterns to determine what the next likely output should be. In our case, this next likely output should be a cohesive, human-readable sentence that describes the objects that were in the encoded image. The output we expect from this model should predict a ¡start¿ key word to indicate the start of the sentence when the image is passed into the LSTM network, and should predict a sequence of words that relate to the image passed into the LSTM encoder, and once completed this sentence should end with the keyword ¡end¿. At the end there should be a human-readable sentence.

Another type of convolutional network used was the ResNet-50 architecture. As the name suggests, ResNet-50 contains 50 layers with 48 convolutional layers, one MaxPool layer and one Average Pool layer. Something novel about ResNet in general is that they tackle the vanishing gradient problem where as a model backpropagates the gradient gets so small it "vanishes". Thus, to combat

this they include additional skip connections which allows the network to skip layers during training that aren't relevant and allows the network to be deeper. Our architecture Using the resent50 model removed the final layer, instead replacing it with a fully connected layer that had an output of dimension 300. This decision was made to better fit the pre-trained model with our existing LSTM architecture.
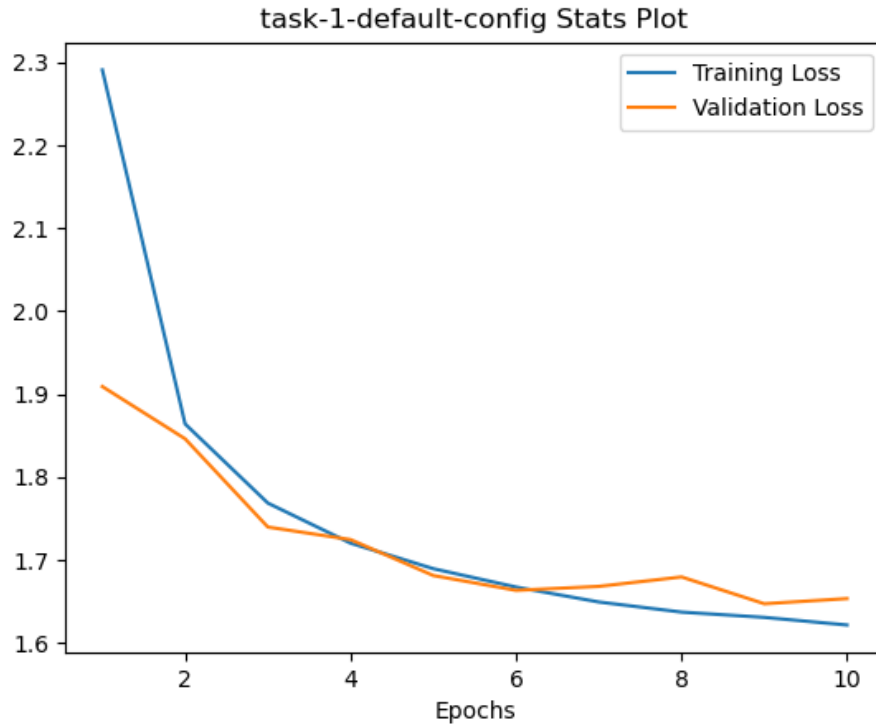


Figure 1: Default Training/Validation Loss

## 3.2 Task 1 Architecture Changes

For our model we changed the embedding and hidden size to confirm the behaviour of our model.

We reduced the embedding size from the original 300 to 150, and expected a decrease in accuracy from our model on our validation set because with a reduction in embedding size when we pass this into the decoder there's effectively less information for the LSTM to go off of, thus we expect less accurate sentences to be produced. Our thoughts were correct as we achieved BLEU1 score of 40.837 and a BLEU4 score of 1.539 on the validation set.

To view the BLEU results please refer to our Results section. The reduced embedding size results are labelled as "CNN-LSTM Embedding Size 150" and the increased Hidden Size results are labelled as "CNN-LSTM Hidden Size 1024".
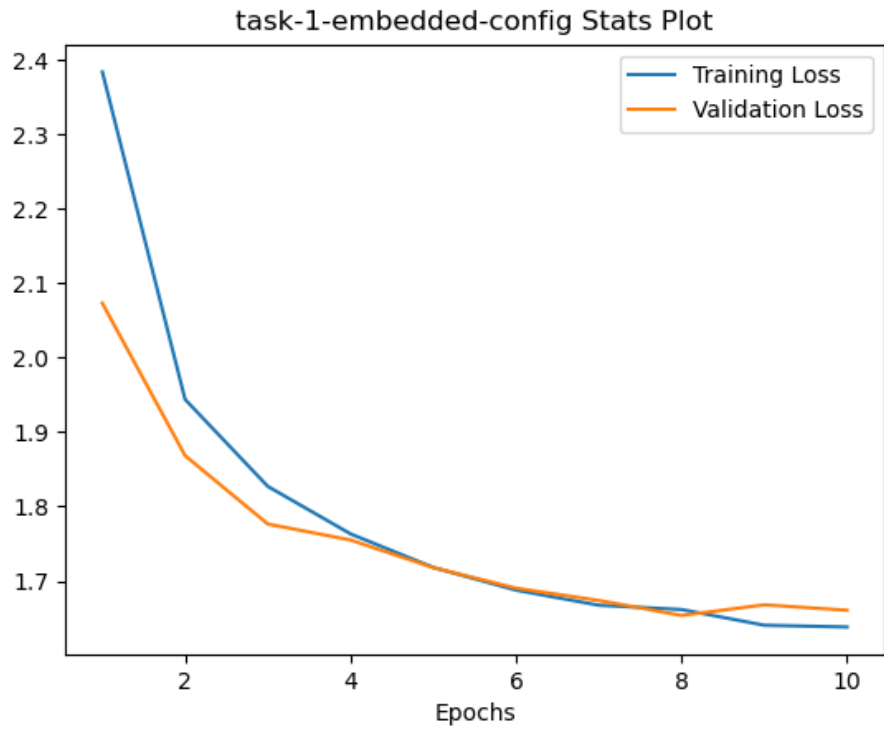
4

Figure 2: Embedding Size Reduction 150 Training/Validation Loss

For a second architecture change we increased the hidden layer size for the LSTM from the original 512 to 1024, and expected an increase in accuracy from our model because we have a larger network and more memory from past sentences. We wanted to make this change because we were interested if it would affect the performance but also if it would affect the training time, as there are now much more nodes in the LSTM, but we saw minimal change in training time. Our thoughts were correct as we achieved a BLEU1 score of 47.36 and a BLEU4 score of 1.79 on the validation set.
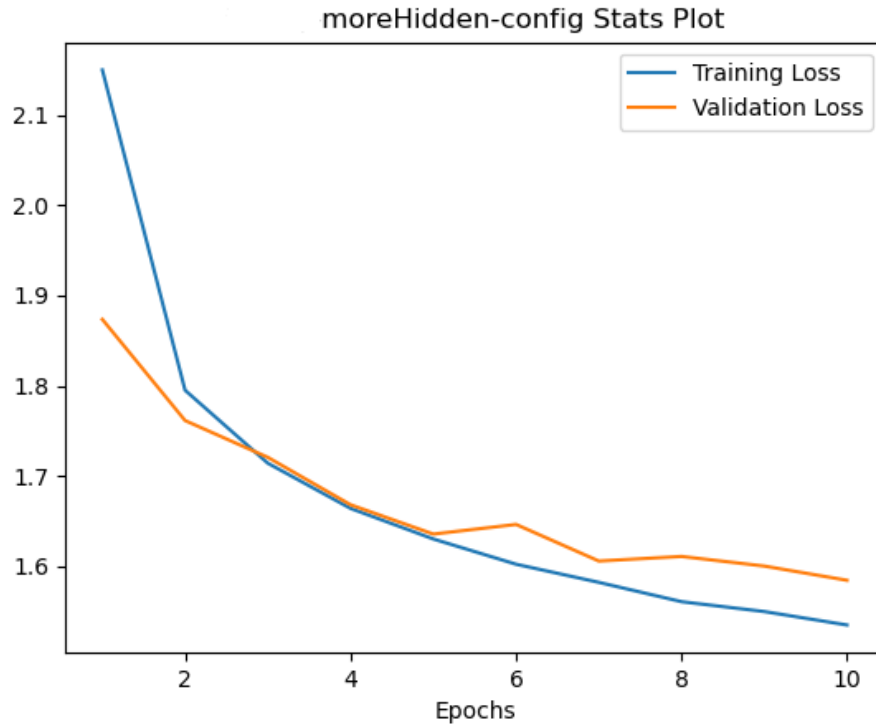
Figure 3: LSTM Hidden Size Increase 1024 Training/Validation Loss

## 3.3 Task 2 Hyperparamter Changes

For this, we decided to change out the optimizer, and instead of using an Adam Optimizer, used the Stochastic Gradient Decent (SGD) optimizer, and increased the learning rate by 100 times. increasing from $5 \times 10^{-}4$ to $5 \times 10^{-}2$ We decided to make these changes because we noticed that the training slowed down between the first 1-2 epochs and beyond then, and wanted to look at the differences between different optimizers. We saw an increase in BLEU-1 score performance, but BLEU-4 performance dropped off, and we had much higher loss compared to the Adam optimizer. As we can see in the loss graph where the minimum is aprox 2.2 compared to 1.6 of the default configuration. This, along with the much worse captioning from the samples we viewed, led us to not select this as our best model, despite it having the best BLEU-1 score. And we noticed that BLEU4 tended to be a better metric of a good model for captioning.
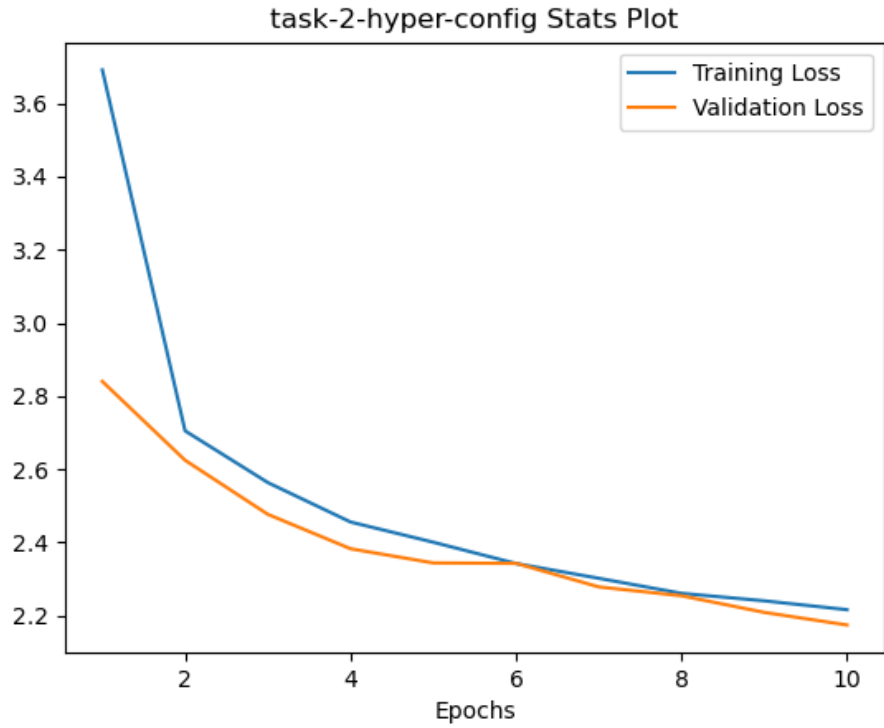
Figure 4: LSTM Hidden Size Increase 1024 Training/Validation Loss

# 4 Results

As described in the previous section our models had different performances depending on what hyperparameters and architecture changes were made. The following table collects the results BLEU-1 and BLEU-4 scores for each of the models we ran from task 1 and 2.

| BLEU Scores | | |
|---|---|---|
| Model | BLEU-1 | BLEU-4 |
| CNN-LSTM Original | 46.97 | 1.63 |
| CNN-LSTM Embedding Size 150 | 40.837 | 1.539 |
| CNN-LSTM Hidden Size 1024 | 47.36 | 1.79 |
| Resnet-50 | 47.81 | 1.70 |
| Resnet-50 new Hyper-params | 49.24 | 1.63 |

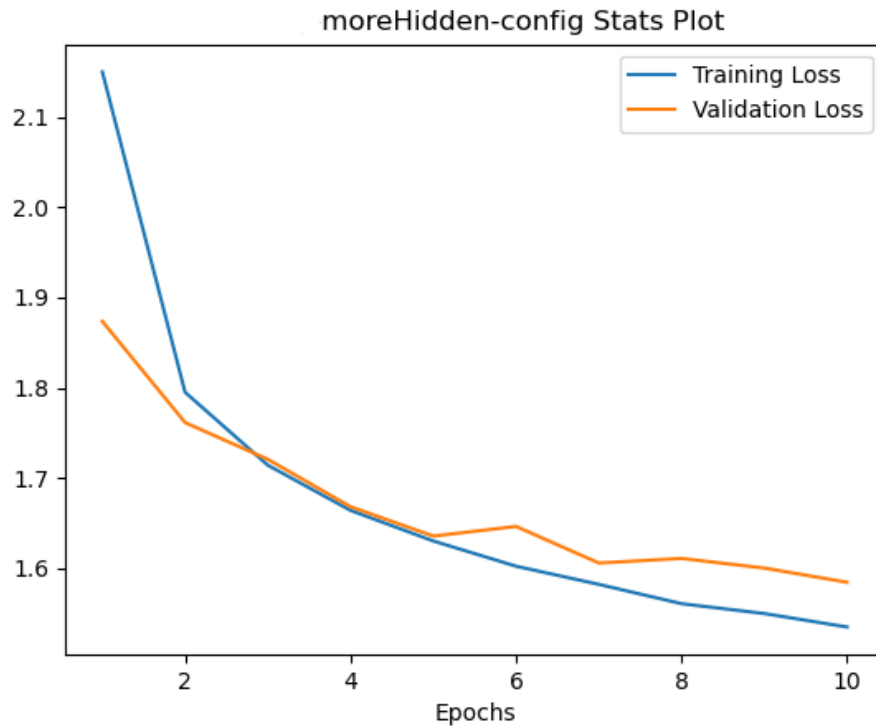## 4.1 Training and Validation Loss Plots



Figure 5: Custom CNN+ LSTM with 1024 hidden nodes Losses

For Task 1, We observed that it took a long time to train, and even though we trained for over two hours, it only completed 10 epochs and would probably benefit from even more training. This is to be expected though from the large amount of complexity in our network, there are a large amount of weights and many layers that must be passed through and back propagated over to train. We also noticed that pulling random images to generate captions for rarely worked well, and that certain topics the network learned to caption well quickly (such as tennis and baseball), probably due to them being more common in the training data.
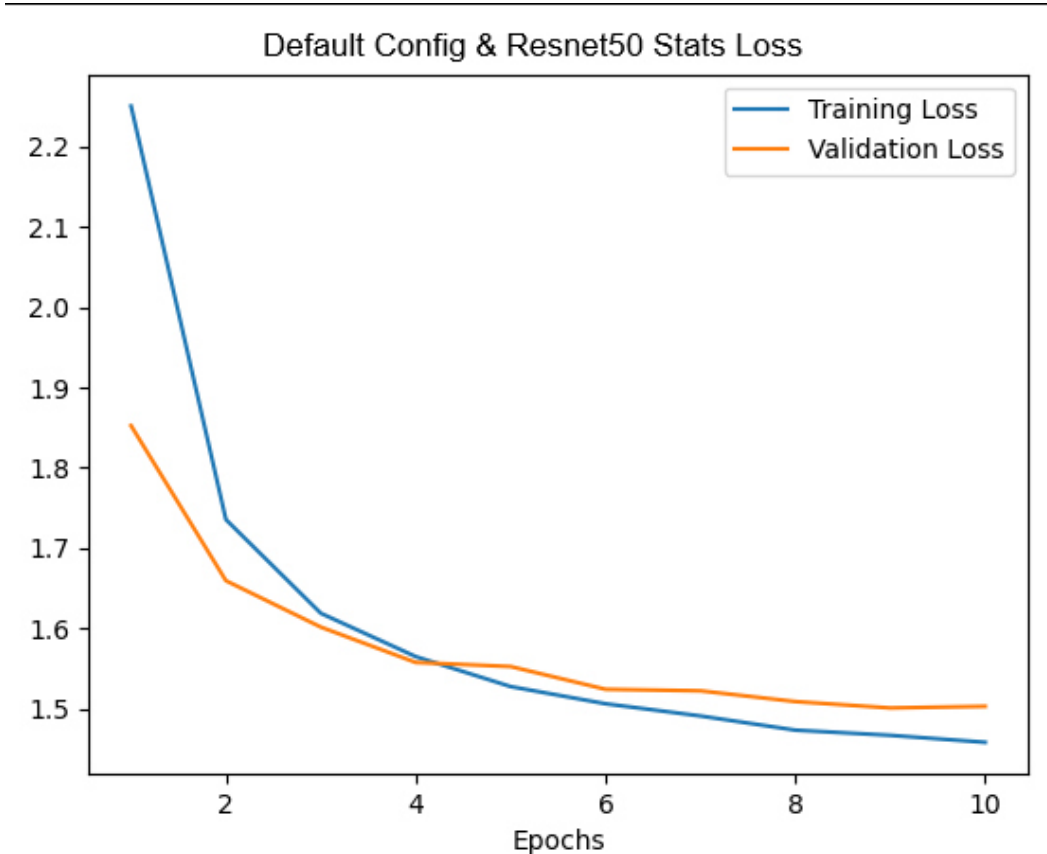
Figure 6: Resnet50 + LSTM Default Configuration Losses

For task 2 We observed that we were able to achieve a slightly lower loss compared to task 1, and that the performance in terms of caption generation was noticably better after switching to Resnet50. This improvement is expected, as we are starting with a more powerful, pre-trained convNet, which means more time can be focused on training just the LSTM portion of the Encoder-decoder network, This network generalized better to more objects / topics, but still struggled to always caption a random image well, and tended to work better simple images compared to busy scenes.

### 4.2 Task 1 vs Task 2

Comparing the two networks, Task 2 using Resnet-50 was a better solution. It gave better captions to more images, and had higher BLEU scores on the test set. We can attribute this to the 50 layers of Resnet and think it was performing better because it was better able to pull out the features of an image and was less sensitive to the backgrounds.

## 5   Captions

For these captions, we provide a Refrence Caption from the dataset as an example of what a "correct" answer would be. The temperature used for these refers to the Stochastic softmax done to add some randomness into the sentences produced. From the captions we can see that using a medium, .4 in this case, temperature performed well, a very high (5) temperature gave completely random words, and a very small temperature (.001) performed identically to a deterministic softmax. Deterministic simply means that we take the word that is most likely instead of randomly selecting one of the usually few likely words.

## 5.1 BEST Task 1: CNN-LSTM with Hidden size = 1024



Figure 7:
Reference captions: a child with a laptop on a table.
Temp=.4: a man is looking at a laptop on a table.
Deterministic: a man is sitting on a bench with a laptop.
High temp (5): tucking aloft bore spiderman onthe phones righteous ceramics jointly investigate down wendy director facilities gross motorcyles gingerbread jello yaks talbot
Low temp (.001): a man is sitting on a bench with a laptop.

Figure 8:
Reference caption: a woman is holding a tennis racquet preparing to serve the ball.
Temp=.4: a woman holding a tennis racquet on a tennis court.
Deterministic: a man is standing in a field with a frisbee.
High temp (5): yong numbers loosing windshields truffles most vice odd missile measurements ingredients fritos toy patients picket weight escorted collection skier stone-paved
Low temp (.001): a man is standing in a field with a frisbee.

Figure 9:
Reference caption: a person riding a skate board on a skate park.
Temp=.4: a person riding a skateboard down a street.
Deterministic: a man is standing in front of a large building.
High temp (5): toa offers disembark player paradise powdered threes hurt summit moldings fluies hardcover them eggbeater vase negatives visitor medals fog pastrami
Low temp (.001): a man is standing in front of a large building .

## 5.2 Bad Captions Task 1: CNN-LSTM with Hidden size = 1024



Figure 10:
Reference caption: car driving down a road behind a lot of sheep.
Temp=.4: a man is standing in a field of food.
Deterministic: a man is standing in front of a large building.
High temp (5): blinded broadcast concentrates threes doll carrot glob room ramen cleaning english livery doughnut motions graces peperoni checked dugout scopes snatching
Low temp (.001): a man is standing in front of a large building.

Figure 11:
Reference caption: happy boy showing off a packaged electric toothbrush.
Temp=.4: a woman is holding a banana and a knife.
Deterministic: a man is sitting on a bench with a dog.
High temp (5): crucifix costume national needless neutral clementines he bell elevating shower panned action valance goat deluxe cutters warranty instructions hots setting
Low temp (.001): a man is sitting on a bench with a dog.

Figure 12:
Reference caption: a man with a leather jacket sitting on a motorcycle in a street.
Temp=.4: a young boy standing on a skateboard on a bench.
Deterministic: a man is standing in front of a large building.
High temp (5): roosts whispy rained outs tattooed recliner grease guy recliner fog bowties lapse conductor anchor stand-up drums coax tommy wharf statues
Low temp (.001): a man is standing in front of a large building.

## 5.3 BEST Task 2: RESNET50-LSTM



Figure 13:
Reference captions: this is a photo of a large crowd of young adults sitting on a long row of park benches
Temp=.4: a woman is sitting on a bench in a black and white photo
Deterministic: a person is standing on a skateboard in the rain
High temp (0.9): the school lady is holding up the water
Low temp (0.3): a man and a woman standing next to a fire hydrant

16

Figure 14:
Reference captions:
Temp=.4: a person is holding a frisbee in the air
Deterministic: a woman is standing in front of a fire hydrant
High temp (0.9): woman holding a skateboard and a child on the concrete
Low temp (0.3): a woman is standing next to a fire hydrant

Figure 15:
Reference captions:
Temp=.4: a man is standing next to a fire hydrant
Deterministic: a person is holding a frisbee in the air
High temp (0.9): a woman stands on a mug in front of a fire hydrant
Low temp (0.3): a woman is standing in front of a fire hydrant

## 5.4 Bad Captions Task 2: RESNET50-LSTM



Figure 16:
Reference captions:
Temp=.4: a person is sitting on a skateboard in a parking lot
Deterministic: a person is standing on a bench with a skateboard
High temp (0.9): a kid colored toy school a beagle on a wall
Low temp (0.3): a woman is holding a frisbee in the air

Figure 17:
Reference captions:
Temp=.4: a woman is sitting on a bench while holding a cell phone
Deterministic: a woman is standing on a skateboard in the rain
High temp (0.9): a couple of fans on a skateboard looking at a toy
Low temp (0.3): a woman is standing by a fire hydrant

Figure 18:
Reference captions:
Temp=.4: a person in a suit is holding a skateboard
Deterministic: a person is standing on a skateboard in the rain
High temp (0.9): a woman holds a phone to the snowy swiss in the rain
Low temp (0.3): a man is sitting on a bench with a skateboard

## 6  Discussion

One aspect of our results was how on multiple instances of training, the network latched onto certain common phrases that occurred often. All of the generated captions begin with the word "a" and often the subject of the generated caption is "a man", "a woman" or "a person" presumably because these are common objects within images in the COCO dataset. Another insight we had into these results were regarding the calculation of loss. Because the sentence is padded, the quickest way to reduce loss is to correctly match the end and pad tags with each teacher-forced iteration. This means that the sentence biases towards simpler and shorter sentences that use more common subjects and actions. Early on in our network we notice that the network generates captions like "a a in a on a a" and similarly worded intellectual insights. This is probably due to the likelihood of any singular word being "a", "in" or "on" is very likely. Since these captions are generated early on, this is probably before the LSTM learns the hidden weights between each time step well enough to understand some form of grammar.

The Resnet50 Network was the best network for generating captions. This judgement is based off of the average BLEU Scores generated for each network. In our experience, the Resnet50 network performed better. This may be because it has already been pre-trained on a larger set of data. It is also a deeper network than own convectional neural network, with 48 convolutional layers. The added depth, allows it to generalize more abstract features from our image and thus may have

allowed the network to perform better.

The deterministic approach does not work well because it simply selects the option output from our softmax layer that has the maximum likelihood of occurring. This means that the network is going to choose the single most likely option based on some input features. Then, the captions will not be unique due many of the input images sharing similar features. Then, many of the captions will be exactly the same for a variety of images. Thus, that's why when the temperature approaches 0, the captions generated would be very close to a deterministic generation as the distribution is completely uniform. However, by having a higher temperature such as 0.4, we'd get a different caption each time. Thus, the network won't latch onto a single outcome that may be completely wrong and the temperature combats this issue in the deterministic method. However, too high of a temperature value injects too much randomness into the decision-making process and thus that's why we observed more randomness in the captions generated, which were not very human-readable.

# 7 Team contributions

## 7.1 Jason Stanley

On this PA I worked on implementing the CNN, connecting the CNN to the LSTM and calculating the cross entropy loss for the LSTM outputs. As a group we all contributed to bug fixes along the way and split up the work for training models and contributing to the write up. Another component I helped with was implementing the pretrained resnet and freezing the past layers, which Elliot lead and I helped with. Overall this group was great to work with and everyone contributed a lot. We worked well together and through in person meetings and discord we were able to get work done in a good amount of time.

## 7.2 Michael Ma

I worked on the LSTM class with initializing layers and the forward function. Throughout the assignment I worked on testing and debugging the LSTM forward function, and implemented caption generation with sampling. I also helped programming the training function and the generate caption function. Finally, I wrote the README file for the Github repository.

## 7.3 Elliot Lee

I helped with the initial implementation of the LSTM encoder/decoder architecture, attempted bug fixing with our training function, and helped with the implementation and testing of the resnet50-lstm architecture changes. Additionally, I trained this model, generated different stochastic captions and recorded statistics for this network. Also I helped write sections of the report.

## 7.4 Annabella Macaluso

For this PA I contributed towards building the CNN and debugging the implementation of the CNN-LSTM. I went to office hours to debug our code and fixed bugs with computing the loss. Ultimately we had to train longer and we were able to receive promising results. I trained models on how architecture changes affected our model and assisted in the writing and formatting of this paper submission.

**References**

[1] Zhang, Aston and Lipton, Zachary C. and Li, Mu and Smola, Alexander J. (2021) Dive into Deep Learning.