# Image Captioning with CNN-RNN and Pretrained Encoder

**Xiyan Shao**
UCSD
x3shao@ucsd.edu

**Yunxiang Chi**
UCSD
yuchi@ucsd.edu

**Zelong Wang**
UCSD
zew013@ucsd.edu

## Abstract

Image captioning is an important task in computer vision. In this paper, we implemented CNN-RNN models for image captioning where an encoder (pretrained or custom) extracts features from images, and then the LSTM decoder generates the captions. For generation during inference time, we used deterministic greedy search and sampled search. Our best model achieved a BLEU-1 score of 66.03 and a BLEU-4 score of 7.37 on the test set of the MS-COCO dataset. Randomly sampled outputs show that the model was able to generate captions that were (partially) semantically correct and grammatically correct.

## 1   Introduction

Image captioning is the process of generating textual description of images, and it is widely applied in search engines, accessibility tools, and content creation softwares. However, generating captions at scale with high accuracy has been a hard problem for traditional solutions. Modern learning-based solutions approach the task by bridging computer vision and natural languages models, and particularly the encoder-decoder architecture rendered excellent results that are close to human levels.

In this project, we trained CNN-RNN models on the COCO 2015 dataset for image captioning. For experiments, we tried different architecture of the encoder including a custom deep CNN encoder and a pretrained ResNet with frozen weights. We also altered the network topology and hyperparameters in an attempt to increase the accuracy. The models were evaluated using BLEU1 and BLEU4 scores, and we generated example captions for case studies.

## 2   Background/Related Work

We heavily referenced Dr. Cottrell's CSE151B slides (lecture 5 and 6) and CS224n lecture notes (5 and 6) by Chris Manning. We looked at and incorporated the pretrained model of He et al.'s *Deep Residual Learning for Image Recognition* (ResNet50). Also, we greatly relied on PyTorch's library page for `LSTM`, `Conv2d`, etc, including the examples.

# 3 Models

## 3.1 Descriptions

### 3.1.1 CNN Encoder

| Layer | Input Channels | Output Channels | Stride Size | Kernel Size | Activation Function | Padding Size |
|---|---|---|---|---|---|---|
| conv1 | 3 | 64 | 4 | 11 | ReLU | 0 |
| maxpool1 | 64 | 64 | 2 | 3 | - | 0 |
| conv2 | 64 | 128 | 1 | 5 | ReLU | 2 |
| maxpool2 | 128 | 128 | 2 | 3 | - | 0 |
| conv3 | 128 | 256 | 1 | 3 | ReLU | 1 |
| conv4 | 256 | 256 | 1 | 3 | ReLU | 1 |
| conv5 | 256 | 128 | 1 | 3 | ReLU | 1 |
| maxpool3 | 128 | 128 | 2 | 3 | - | 0 |
| fc1 | 128 | 1024 | - | - | ReLU | - |
| fc2 | 1024 | 1024 | - | - | ReLU | - |
| fc3 | 1024 | 300 | - | - | - | - |

Note: all of the convolution layers are batch normalized.

### 3.1.2 LSTM Decoder

The LSTM decoder consists of 2 layers. It accepts an embedding of size 300, and the hidden units are of size 512. There is a fully connected head over the LSTM layers that maps the outputs to the vocabulary space.

### 3.1.3 ResNet Decoder

The ResNet-50 network decoder is pretrained for an image classification task and has the weights frozen. It has 50 convolution layers with varying output channels, and there are extra residual connection to ease the training of the network and increase generalizability.

## 3.2 Task 1: Architecture Changes

We called the architecture-changed model as CNN2. In CNN2, we increase two convolutional layers(one with 256 output channel, 5*5 kernel, and 2 padding, the other with 256 output channel, 3*3 kernel, and 1 padding), 2 maxpool layers(both with 3*3 kernel and stride 2), and dropout 25 percent hidden units after first two fully-connected layers. We add 2 convolutional layers and maxpool layers because it will help us to extract more feature from convolution and summarize it from maxpool. In this way, the most expected feature can be extracted and learned by network at a better chance. Then, adding dropout for the first two fully-connected layers is to reduce the negative effect from overfitting.

## 3.3 Task 2: Hyperparameter Changes

### 3.3.1 Add Regularization (Weight Decay)

We added a $L2$ Regularization of $0.0001$ to the Adam optimizer in order to avoid overfitting. Regularization might help our model to generalize better so that it performs better with test data.

### 3.3.2 Add Image Augmentation

After adding the $L2$ Regularization, we also add an option of doing image augmentation. We use both *.transforms.RandomCrop* and *.transforms.ColorJitter* to transform the images. First, we apply resize to all training images and make all training images of size 316. Then we apply a RandomCrop of image size 256. After that, we apply ColorJitter with brightness factor, contrast factor, saturation factor, and hue factor chosen uniformly from the range $[\max(0, 1 - 0.3), 1 + 0.3]$

When the epoch is even, our model is trained by transformed training images. When the epoch is odd, our model is trained by images from the original dataset. And with this augmented dataset and early stopping, we can train the model with more epochs. However, due to time concerns, we trained it with 14 epochs.

We decide to use RandomCrop because in some cases, the caption is ignoring information near the edges of the image while in other cases, all elements within an image are described in the caption. But in most cases, elements at the center of the image should be correctly described to match the target caption. Therefore, our RandomCrop randomly captures 0.64 percent of the image. And we wish a combination of cropped and uncropped training datasets would make the model better at recognizing both cases.

We want to use ColorJitter because sometimes color information can be misleading. For example, the model might think of a man surfing in a muddy river as riding a horse on the road simply because the brown color connects more to horses or streets than water.

# 4 Results

## 4.1 Results Table

| Model | BLEU-1 | BLEU-4 |
|---|---|---|
| CNN Default | 55.52 | 3.85 |
| CNN Architecture Change | 51.70 | 3.10 |
| Resnet Default | 64.66 | 6.17 |
| Resnet Regularization | 66.03 | 7.37 |
| Resnet Image Augmentation | 64.55 | 6.26 |

All models' BLEU score are above expectation, but the model CNN Architecture Change didn't work quite well compared to others, and we'll talk about that in the discussion section.
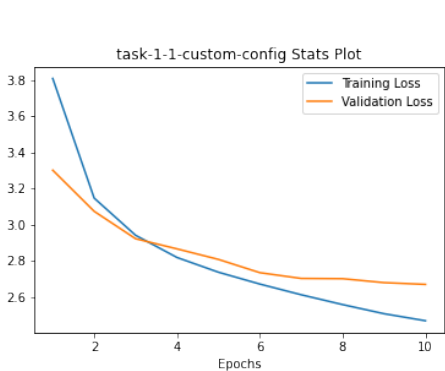
## 4.2 Results Plot



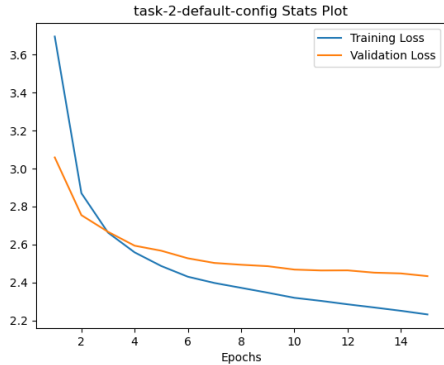Figure 1: Training and Validation Loss of CNN Default Architecture

Figure 2: Training and Validation Loss of Resnet with Regularization

The best model for task 1 is CNN with default setting, its test loss is 2.67, and the best model for task 2 is Resnet with Regularization with test loss 2.44.

Separately, both of them show the similar pattern: both losses reduce quickly when training on the first three epochs and reduce gradually for the rest epochs. When comparing them together, we found the difference that the training loss quick-reducing period for resnet is longer than that of CNN, which

its slope become smaller until around epoch 5. The validation loss show similar pattern for both, but it beginning point(the loss at the beginning) for Resnet is 0.2 less than that of CNN.

## 5 Captions



Image 422583 (Task 1, Good)

- **Reference Caption:** ['Man catching air, while skateboarding in a public park.', 'A man flying through the air while riding a skateboard.', 'A man doing tricks on a skateboard outdoors in a city.', 'A skateboarder practicing a jump under a bridge', 'A man dressed for cold weather doing a jump on a skeateboard.']

- **Deterministic**: a man is riding a skateboard down a street

- **Sampled** ($t = 0.4$): a man is riding a skateboard down a street

- **Sampled** ($t = 0.001$): a man is riding a skateboard down a street

- **Sampled** ($t = 5$): blue/green suburbs gassy leaned outreached cuddled wilderness body city taking yamaha have wiimote want conditions parasurfing window kiwi



Image 253470 (Task 1, Good)

- **Reference Caption:** ['The gold tie looks especially bright against a dark shirt.', 'A man wearing glass and a tie looking to the side.', 'a man wearing glasses and an orange tie', 'A yellow tie stands in contrast to the black and white image of a man wearing glasses.', 'A man is standing and looking attentively to the side.']

- **Deterministic**: a man in a suit and tie is standing in a room

- **Sampled** ($t = 0.4$): a man in a suit and tie holding a tie

- **Sampled** ($t = 0.001$): a man in a suit and tie is standing in a room

- **Sampled** ($t = 5$): mules stoop hoagie salvador hes cockatiel paint hotdog slender wreaths kisses breasts ken glassed enterprise salutes prepare rules

Image 479094 (Task 1, Good)

- **Reference Caption:** ['A beach with a lot of people parasailing in the water', 'Many people on the beach are flying kites.', 'People are flying kites on the beach on a hazy day.', 'As some kites lie on the beach others are in full flight.', 'A group of people on the beach watching parasailers. ']
- **Deterministic**: a large airplane is parked on the runway
- **Sampled** ($t = 0.4$): a group of people standing on a beach by a beach
- **Sampled** ($t = 0.001$): a large airplane is parked on the runway
- **Sampled** ($t = 5$): kitche ref feeds pointed stony night-time broccoli shopped loan diced pained devastated outcropping camerea create grass mound texing

Image 108862 (Task 1, Bad)

- **Reference Caption:** ['An orange and an apple on a table. ', 'A close up view of an orange and an apple on a wooden table. ', 'An apple and an orange sitting on a table.', 'A pair of fruits sit on a table top.', 'a close up of an orange and an apple ']
- **Deterministic**: a plate of food with a sandwich and a fork
- **Sampled** ($t = 0.4$): a plate of food with a sandwich and a knife
- **Sampled** ($t = 0.001$): a plate of food with a sandwich and a fork
- **Sampled** ($t = 5$): central mustard universal beetlejuice cobbler wipeout cubs stand road blowup sidelines fax mid-century pickle schoolbus c motorcyle washington

Image 521976 (Task 1, Bad)

- **Reference Caption:** ['a man is surfing on a wave at the beach', 'Young men are waiting to ride in a wave simulator.', 'A group of male surfers practicing in a river.', 'A group of boys surfing in water ', 'A man hitting a wave on a surf board in a man made wave pool.']
- **Deterministic**: a man riding a horse on a dirt road
- **Sampled** ($t = 0.4$): a man riding a horse down a road with a dog
- **Sampled** ($t = 0.001$): a man riding a horse on a dirt road
- **Sampled** ($t = 5$): focusing developed hello grassland brim urges chariot cuddles luxury ha ganache misc comforter faced google extinguisher propelled 23rd

Image 360701 (Task 1, Bad)

- **Reference Caption:** ['A blue and yellow airplane flying through a blue sky.', 'a blue and yellow plane flying in the sky with a smoke trail behind it', 'A blue and yellow airplane performing a stunt.', 'Airplane with smoke coming out flying through blue skies.', 'An airplane flying through the sky does sky writing.']
- **Deterministic**: a person on a snowboard in the air
- **Sampled** ($t = 0.4$): a man flying through the air while a man is on the beach
- **Sampled** ($t = 0.001$): a person on a snowboard in the air
- **Sampled** ($t = 5$): af surfboarding representatives threesome peaches weaved while kneepads monroe intensely incoming pipeline slug hairy

Image 273362 (Task 1, Good)

- **Reference Caption:** ['A guy doing something with a Frisbee out in a field.', 'A guy lunging forward to catch a white Frisbee.', 'A man catching a Frisbee during a game of Frisbee Golf.', 'Man reaching to catch a frisbee in a park.', 'A man in a park about to catch a Frisbee. '],

- **Deterministic**: a man holding a frisbee in a field .

- **Sampled** ($t = 0.4$): a man holding a frisbee in a field .

- **Sampled** ($t = 0.001$): a man holding a frisbee in a field .

- **Sampled** ($t = 5$): bug teacups cycling generators battered rested shoes overturned dessert commodes sonic jellies fulled pockets solid bulldozer observes ave.



Image 462138 (Task 1, Good)

- **Reference Caption:** ['A man and woman sitting by a wooden table with champagne and cake.', 'Man and older woman smiling with champagne and cake', 'An old lady and a young man sitting in front of a table on which a cake is kept.Two wine glasses are also there.', 'On a beige couch, a young man and older woman admire a cake on a coffee table.', 'Family celebrating; sitting on couch with glasses of wine and food. ']

- **Deterministic**: a group of people standing around a table with a cake .

- **Sampled** ($t = 0.4$): a group of people sitting at a table with a cake

- **Sampled** ($t = 0.001$): a group of people standing around a table with a cake .

- **Sampled** ($t = 5$): photoshopped reporters garnished bride amusement tartar stems bid ms technique muli-colored machete them uses sly slippery numbers mysterious

6

Image 360701 (Task 1, Good)

- **Reference Caption:** ['A blue and yellow airplane flying through a blue sky.', 'a blue and yellow plane flying in the sky with a smoke trail behind it', 'A blue and yellow airplane performing a stunt.', 'Airplane with smoke coming out flying through blue skies.', 'An airplane flying through the sky does sky writing.']
- **Deterministic**: a plane flying in the air with a blue sky .
- **Sampled** ($t = 0.4$): a small jet jet flying through the air .
- **Sampled** ($t = 0.001$): a plane flying in the air with a blue sky .
- **Sampled** ($t = 5$): headed stret along alfredo bicyclists worker commercial discs tip flown docket bunting pecans inspects coming longboarder cloud-filled buidling



Image 108862 (Task 1, Bad)

- **Reference Caption:** ['An orange and an apple on a table. ', 'A close up view of an orange and an apple on a wooden table. ', 'An apple and an orange sitting on a table.', 'A pair of fruits sit on a table top.', 'a close up of an orange and an apple ']
- **Deterministic**: a small bowl with a bunch of apples on it .
- **Sampled** ($t = 0.4$): a yellow plate with an apple and a knife
- **Sampled** ($t = 0.001$): a small bowl with a bunch of apples on it .
- **Sampled** ($t = 5$): smiley manicure mason borne worst outside selling har-tru bp seals dug stretched zerbra rapid architecture kempa reacts alaskan

Image 043580 (Task 1, Bad)



Image 445887 (Task 1, Bad)

- **Reference Caption:** ['Many birds are perched on the weather vane of a building.', 'Several birds perched on top of a weather vane.', 'Many birds are sitting on the pole erected at the top of a temple.', 'A group of birds resting on roof and weather vane.', 'The birds are sitting atop the flag on the tower.']

- **Deterministic**: a clock on a pole with a street sign .

- **Sampled** ($t = 0.4$): a clock that is flying in the air .

- **Sampled** ($t = 0.001$): a clock on a pole with a street sign .

- **Sampled** ($t = 5$): stirs sheepdog lounge earbuds butts thing wasteland reins nunchuck books skateboad oxen refridgerator there buxom features pitts top

- **Reference Caption:** ['A pizza sitting on top of a BBQ grill next to a tree.', 'A pizza is cooking on a grill outside.', 'A large cooked pizza on a grill outside.', 'Pizza on the grill keeps the kitchen cool on a hot day.', 'a silver BBQ with a square pizza and grass and bushes']

- **Deterministic**: a pizza is being prepared in a pan .

- **Sampled** ($t = 0.4$): a pizza is sitting on a stove top and oven

- **Sampled** ($t = 0.001$): a pizza is being prepared in a pan .

- **Sampled** ($t = 5$): grinds cart truck basil except plantains dragged coming exposure expandable envrioment centre dealership chamber garland sku foldable attacking

## 6 Discussion

For the two models in task one, their BLEU1 and BLEU4 scores are both above the expectation level, which is 45 for BLEU1 and 1.9 for BLEU4. Compared within those two models, we got the best model of default-setting CNN that has better BLEU1 socre, BLEU4 score, and smaller Test Loss.

Using the deterministic approach means choosing the one with maximum value, but we found that deterministic approach didn't quite work well. We believe the reason is that take skateboard as an example, if the test case is exactly skateboard, then the same feature will make the model return the correct description. However, if the test case is a little bit different, like surfboard, or quite different, like airplane, similar extracted feature will also make the model return the skateboard related description. In another word, it doesn't provide a range of possibility given extracted feature, and that's why its performance is bad at testing.

Then, we tested models' performance with very low temperature(0.001) and very high temperature(5).

When using very low temperature, models' performance is a little worse than using 0.4 temperature since the description is similar for most of them but doesn't make sense for some of them. A very low temperature means it's basically deterministic, and deterministic approach will reduce the performance as we discussed above.

When using very high temperature, the output description are totally nonsense, not even a readable sentence. The reason is that when the temperature is very big, the distribution will not be a normal distribution, but almost a uniform distribution, which will cause the output just randomly pick work without consider anything. From our trials in task 2, we find that adding a weight decay of 0.0001 leads to the best result. Compared with the task 2 base model without regularization, weight decay controls the overfitting issue and allows our model to generalize better to the test set with higher BLEU scores.

After augmenting the training images by applying RandomCrop and ColorJitter, both BLEU scores decrease. We think the augmentation makes the model worse due to three reasons.

First, some valuable information that usually appears at the edges of an image might be cropped off. For example, the road or the water.

Second, the information at the center of an image often correlated with the information at the edges. For example, the water at the bottom of an image and the ship or surfboard in the center, or the road at the bottom and the car or horse in the center. And cropping the edges could lead to the missingness of this connection.

Third, we find that colors are often a very important description in the captions. However, after ColorJitter, the color might be distorted. Therefore, the model couldn't properly generate the color descriptions. Thus, the BLEU scores are lower.