
Image Captioning with Multiple Decoder Architectures: A Comparative Study of RNN Architectures

Owen Strength
Auburn University
ods0005@auburn.edu

Jacob Simmons
Auburn University
jss0112@auburn.edu

Abstract

Image captioning is a challenging task that requires combining computer vision and natural language processing. In this project, we conduct a comprehensive comparison of different recurrent neural network (RNN) architectures for caption generation, utilizing a pretrained ResNet50 model as the image encoder. We evaluate four distinct decoder architectures: basic RNN, GRU, and LSTM, all trained on the MSCOCO dataset. Our analysis focuses on the trade-offs between these architectures in terms of performance, complexity, and computational requirements. Results demonstrate the relative advantages of each approach, with particular emphasis on the benefits of long-term memory retention in the LSTM-based model. This comparative study provides insights into the optimal choice of recurrent architecture for image captioning tasks under different constraints.

1 Introduction

Generating meaningful captions for images is a fundamental problem in AI that requires understanding both visual and textual modalities. Applications include accessibility tools, content moderation, and automated reporting. The sequential nature of language makes recurrent architectures a natural choice for caption generation, as they process and generate text word by word while maintaining contextual information. While transformer-based models have shown impressive results in recent years, understanding the optimal recurrent architecture for this task remains valuable due to their efficiency and interpretability.

Our objective is to identify the best performing recurrent architecture for image captioning. Using ResNet50 as a fixed feature extractor, we evaluate four architectures with increasing complexity:

- **Basic RNN:** Serving as our baseline model
- **GRU:** Offering potentially better efficiency with moderate performance
- **LSTM:** Addressing long-term dependencies common in caption generation

Through systematic evaluation using BLEU, CIDEr, and METEOR metrics, we aim to determine whether more sophisticated architectures like LSTM provide meaningful improvements over simpler variants for image captioning. While the GRU architecture offers an interesting balance of complexity and capability, our primary focus is on identifying which model achieves the highest caption quality regardless of computational considerations.

2 Background

2.1 Image Feature Extraction

ResNet50 [3]: A deep convolutional neural network pre-trained on ImageNet [2], known for its ability to extract hierarchical visual features through residual connections. We utilize this as our encoder, freezing its weights to leverage transfer learning while adding a trainable embedding layer to adapt the features for caption generation.

2.2 Recurrent Architectures

2.2.1 Basic RNN

The simplest form of recurrent network processes sequences using a single hidden state, updating it at each time step t according to:

$$h_t = \tanh(W_{xh}x_t + W_{hh}h_{t-1} + b_h) \quad (1)$$

where h_t is the hidden state at time t , x_t is the input, W_{xh} and W_{hh} are weight matrices, and b_h is the bias term. The output is then computed as:

$$y_t = W_{hy}h_t + b_y \quad (2)$$

While conceptually simple, this architecture suffers from the vanishing gradient problem during backpropagation through time (BPTT). The gradient signal diminishes exponentially as it propagates backward through time steps, making it difficult to learn long-term dependencies. This occurs because the gradient depends on repeated multiplication of the weight matrix W_{hh} , leading to either vanishing or exploding gradients depending on the eigenvalues of this matrix.

2.2.2 GRU (Gated Recurrent Unit)

GRU addresses the vanishing gradient problem by introducing two gates while maintaining a relatively simple architecture. The key equations are:

$$z_t = \sigma(W_z[h_{t-1}, x_t] + b_z) \quad (3)$$

$$r_t = \sigma(W_r[h_{t-1}, x_t] + b_r) \quad (4)$$

$$\tilde{h}_t = \tanh(W_h[r_t \odot h_{t-1}, x_t] + b_h) \quad (5)$$

$$h_t = (1 - z_t) \odot h_{t-1} + z_t \odot \tilde{h}_t \quad (6)$$

where σ is the sigmoid function and \odot represents element-wise multiplication. The update gate z_t determines how much of the previous state to retain, while the reset gate r_t controls how much of the previous state to forget when computing the candidate state. This architecture is simpler than LSTM because:

- It combines the forget and input gates into a single update gate
- It merges the cell state and hidden state
- It requires fewer parameters and computations

2.2.3 LSTM (Long Short-Term Memory)

LSTM provides the most comprehensive control over information flow through three gates and a separate memory cell. The mathematical formulation is:

$$f_t = \sigma(W_f[h_{t-1}, x_t] + b_f) \text{ (Forget Gate)} \quad (7)$$

$$i_t = \sigma(W_i[h_{t-1}, x_t] + b_i) \text{ (Input Gate)} \quad (8)$$

$$o_t = \sigma(W_o[h_{t-1}, x_t] + b_o) \text{ (Output Gate)} \quad (9)$$

$$\tilde{c}_t = \tanh(W_c[h_{t-1}, x_t] + b_c) \text{ (Candidate Cell State)} \quad (10)$$

$$c_t = f_t \odot c_{t-1} + i_t \odot \tilde{c}_t \text{ (Cell State)} \quad (11)$$

$$h_t = o_t \odot \tanh(c_t) \text{ (Hidden State)} \quad (12)$$

The LSTM's complexity provides several advantages:

- Separate cell state (c_t) provides a direct pathway for gradient flow, mitigating vanishing gradients
- Independent forget gate allows explicit control over memory retention
- Output gate enables the network to maintain information in the cell state without affecting the hidden state
- Three separate gates provide fine-grained control over information flow

2.3 Dataset

MSCOCO Dataset: The Microsoft Common Objects in Context (MSCOCO) [5] dataset is a large-scale dataset designed for multiple computer vision tasks including image captioning. It contains over 330,000 images with 5 human-annotated captions per image, providing a rich and diverse range of scenes and descriptions. We utilize the 2017 training split which consists of 118,287 images for training and 5,000 images for validation.

The dataset's key characteristics make it particularly suitable for image captioning:

- High-quality annotations with multiple captions per image, capturing different human perspectives
- Complex scenes containing multiple objects in their natural context
- Diverse vocabulary with over 30,000 unique words in the captions
- Wide range of visual scenarios including indoor/outdoor scenes, human activities, and object interactions

For our experiments, we use the Karpathy splits, a widely adopted configuration in the image captioning community that provides 113,287 training images, 5,000 validation images, and 5,000 test images. Each image is preprocessed using standard transformations including resizing to 256x256 pixels, random cropping to 224x224, horizontal flipping for augmentation, and normalization using ImageNet statistics.

3 Method

3.1 Architecture Overview

Our system consists of two main components:

Encoder (ResNet50):

- Pretrained weights frozen to preserve ImageNet knowledge
- Final fully connected layer replaced with a trainable embedding layer (size: 256)
- Outputs feature vectors suitable for caption generation

Decoder Variants: All decoders share common characteristics:

- Embedding layer for word tokens (size: 256)
- Hidden size of 512 units
- Single layer architecture
- Output layer mapping to vocabulary size

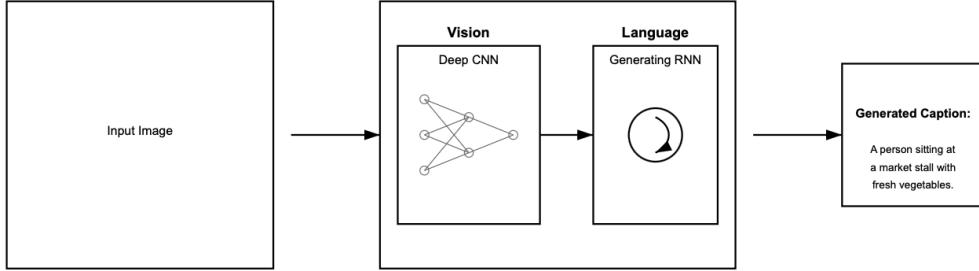


Figure 1: Overview of the image captioning architecture. The model consists of a visual encoder (Deep CNN) that processes the input image and a language decoder (RNN) that generates the output caption. The CNN extracts visual features which are then fed into the RNN to generate natural language descriptions word by word.

3.2 Implementation Details

Each decoder architecture builds upon the previous, increasing in complexity and theoretical capacity. The basic RNN serves as our baseline, using a simple recurrent layer with tanh activation to process sequential data. Building on this, the GRU introduces update and reset gates to better control information flow and improve gradient propagation through the network. The LSTM further enhances the architecture by implementing separate memory cells with input, forget, and output gates, allowing for more sophisticated memory retention and feature processing.

3.3 Model Architectures

Each decoder is implemented with the following specifications:

- **Basic RNN:**
 - Simple recurrent layer
 - Tanh activation function
 - Single hidden state
- **GRU:**
 - Update and reset gates

- Improved gradient flow
- Combined hidden and cell states

- **LSTM:**

- Separate memory cell
- Input, forget, and output gates
- Enhanced memory retention

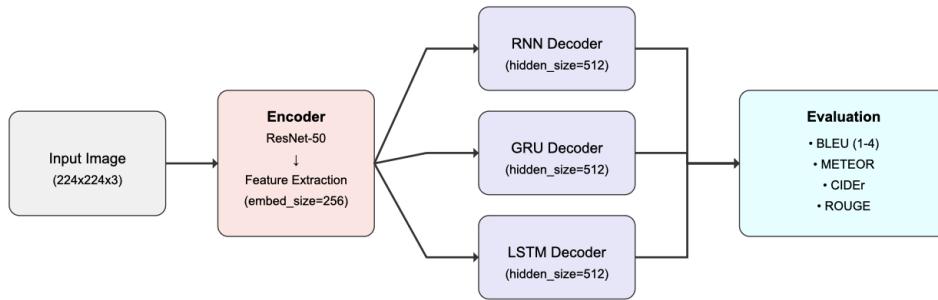


Figure 2: Overview of the separate decoder architectures

3.4 Training Protocol

Our training configuration was carefully selected to balance model performance and computational constraints. We employed the Adam optimizer with a learning rate of 1e-3, chosen for its adaptive learning rate properties and consistent performance in similar vision-language tasks. Cross-entropy loss was selected as our objective function due to its effectiveness in classification tasks and caption generation. We set the batch size to 128, which maximized GPU memory utilization while maintaining stable training dynamics. The models were trained for 5 epochs each, with each epoch taking approximately 50 minutes on either an NVIDIA A6000 GPU or a MacBook Pro M3 Pro with 36GB unified memory. This training duration proved sufficient for model convergence while remaining computationally feasible. For data augmentation, we implemented random horizontal flips and the standard ResNet preprocessing pipeline of resizing images to 256x256 followed by random crops to 224x224, maintaining consistency with the pretrained ResNet50’s training regime.

4 Hyperparameter Tuning

Our hyperparameter selection was primarily guided by the seminal “Show and Tell” paper by Google [7]. Following their successful architecture, we adopted their core hyperparameters: a learning rate of 1e-3, embedding dimension of 256, hidden size of 512, and a single layer for our recurrent networks. The vocabulary size was determined by our tokenizer’s vocabulary. While the original paper explored various configurations, we found these parameters provided a strong baseline for our implementation.

5 Evaluation and Results

We evaluate our models using standard metrics in image captioning:

BLEU (Bilingual Evaluation Understudy): Measures similarity between generated and reference captions by comparing overlapping word sequences of varying lengths (n-grams), from single words (BLEU-1) to four-word sequences (BLEU-4) [6]. This metric evaluates the precision of word matches, with BLEU-4 being the standard reported score as it better captures phrase-level accuracy.

CIDEr (Consensus-based Image Description Evaluation): Measures consensus in generated captions, weighing n-grams by TF-IDF scores [?]. This metric focuses on capturing the importance of specific words in the image context, making it particularly effective for evaluating the relevance of generated captions.

METEOR (Metric for Evaluation of Translation with Explicit ORdering): Measures similarity between captions by aligning words based on exact matches, stems, and synonyms, with additional consideration for word order [1]. This metric often correlates better with human judgments than pure n-gram based approaches.

ROUGE (Recall-Oriented Understudy for Gisting Evaluation): Measures the overlap between generated captions and reference captions by comparing both consecutive word matches and in-sequence matches [4]. This metric focuses on recall and is particularly sensitive to shared word sequences, making it effective for evaluating the completeness of generated captions.

Table 1: BLEU, CIDEr, METEOR, and ROUGE Scores for Different Decoder Architectures

Metric	GRU	LSTM	Naive RNN
BLEU-1	63.5	66.3	64.7
BLEU-2	43.3	45.4	42.7
BLEU-3	30.0	31.7	29.1
BLEU-4	20.8	22.0	20.1
CIDEr	66.1	68.7	60.8
METEOR	43.1	43.6	40.7
ROUGE	50.0	50.9	49.2

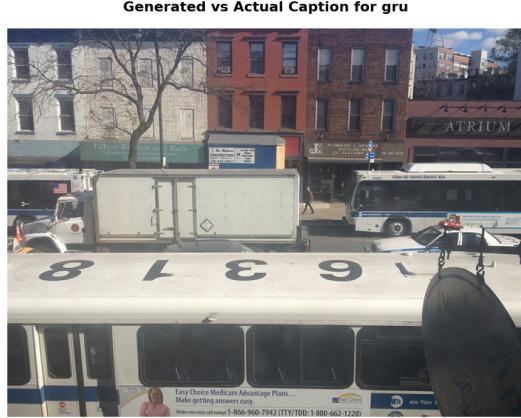


Figure 3: RNN Generated Caption

The performance of different decoder architectures for image captioning using a ResNet encoder is summarized in Table 1. The results indicate that the LSTM decoder achieves the best scores across all metrics, including BLEU, CIDEr, METEOR, and ROUGE. This highlights the effectiveness of LSTM in generating high-quality captions.

However, it is noteworthy that the GRU decoder performs very close to the LSTM in all metrics. For instance, the BLEU-4 score for the GRU is 20.8 compared to 22.0 for the LSTM, and the CIDEr score for the GRU is 66.1 versus 68.7 for the LSTM. The METEOR and ROUGE scores also show minimal differences between the two architectures.

Given that GRU is computationally simpler and requires fewer parameters compared to LSTM, its performance makes it an attractive alternative in scenarios where computational efficiency is a priority.



Actual: there are city buses that are lined on the street in front of the building
Predicted: a large passenger bus parked next to a building.

Figure 4: GRU Generated Caption



Actual: a train traveling through a rural country side covered in grass.
Predicted: a train traveling through a rural countryside.

Figure 5: LSTM Generated Caption

The Naive RNN decoder, while functional, lags behind both the GRU and LSTM in every metric, highlighting its limitations in this context.

6 Conclusion

In this work, we explored different recurrent neural network architectures for image captioning, implementing and comparing Naive RNN, LSTM, and GRU models. We also attempted to implement an attention mechanism with LSTM [8], though this model did not converge during training. The attention mechanism was designed to allow the model to focus on different parts of the image features while generating each word, potentially enabling more precise and contextually relevant captions. While the theoretical benefits of attention are well-documented in literature, getting such models to converge often requires careful hyperparameter tuning and longer training times than were available for this project.

Our experimental results showed that the LSTM model achieved the best performance metrics, followed closely by the GRU model, while the Naive RNN performed notably worse. This aligns with theoretical expectations, as both LSTM and GRU are designed to handle the vanishing gradient problem that plagues simple RNNs. This allows them to better capture long-term dependencies in sequential data. The LSTM’s superior performance can be attributed to its more sophisticated gating mechanism, which allows for better control over information flow through the network.

However, when considering the trade-off between performance and model complexity, the GRU is a particularly attractive option. The GRU model offers an excellent balance of efficiency and effectiveness with fewer parameters and a simpler architecture than LSTM, but nearly matching performance. This makes it potentially the best choice for practical applications where computational resources may be limited.

Future work could focus on successfully implementing the attention mechanism, which would require more extensive hyperparameter optimization and training time. Additionally, exploring more recent architectures like Transformers could provide interesting comparisons to these traditional recurrent approaches. The integration of pre-trained vision models for feature extraction could also potentially improve the quality of the generated captions across all architectures.

References

- [1] S. Banerjee and A. Lavie. METEOR: An automatic metric for MT evaluation with improved correlation with human judgments. In J. Goldstein, A. Lavie, C.-Y. Lin, and C. Voss, editors, *Proceedings of the ACL Workshop on Intrinsic and Extrinsic Evaluation Measures for Machine Translation and/or Summarization*, pages 65–72, Ann Arbor, Michigan, June 2005. Association for Computational Linguistics.
- [2] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pages 248–255, 2009.
- [3] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. *CoRR*, abs/1512.03385, 2015.
- [4] C.-Y. Lin. ROUGE: A package for automatic evaluation of summaries. In *Text Summarization Branches Out*, pages 74–81, Barcelona, Spain, July 2004. Association for Computational Linguistics.
- [5] T. Lin, M. Maire, S. J. Belongie, L. D. Bourdev, R. B. Girshick, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick. Microsoft COCO: common objects in context. *CoRR*, abs/1405.0312, 2014.
- [6] K. Papineni, S. Roukos, T. Ward, and W.-J. Zhu. Bleu: a method for automatic evaluation of machine translation. In P. Isabelle, E. Charniak, and D. Lin, editors, *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA, July 2002. Association for Computational Linguistics.
- [7] O. Vinyals, A. Toshev, S. Bengio, and D. Erhan. Show and tell: A neural image caption generator. *CoRR*, abs/1411.4555, 2014.
- [8] K. Xu, J. Ba, R. Kiros, K. Cho, A. C. Courville, R. Salakhutdinov, R. S. Zemel, and Y. Bengio. Show, attend and tell: Neural image caption generation with visual attention. *CoRR*, abs/1502.03044, 2015.

7 Appendix

7.1 Sample Outputs

Examples of generated captions from each model, along with corresponding images and reference captions.

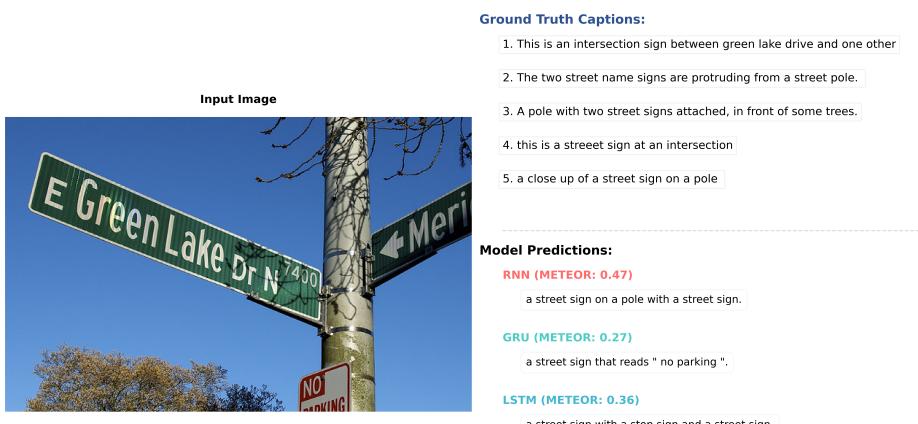
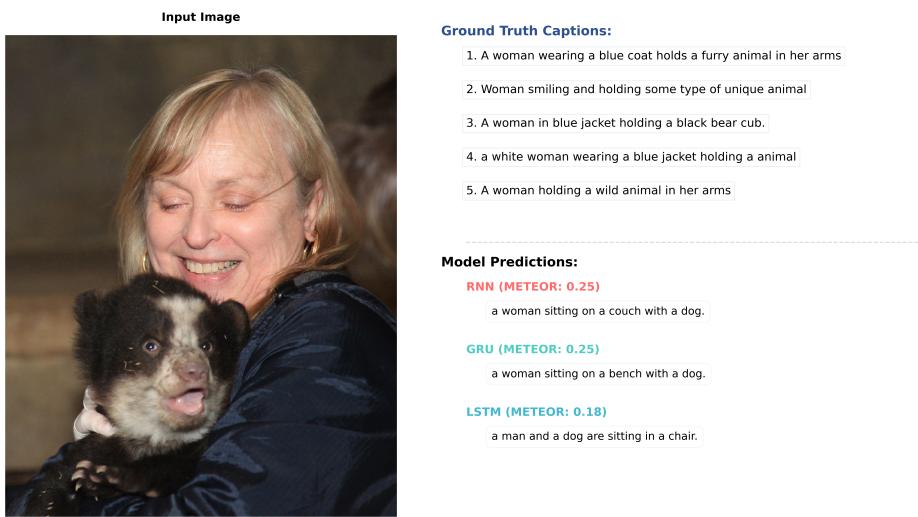


Figure 6: Sample outputs (1/4)

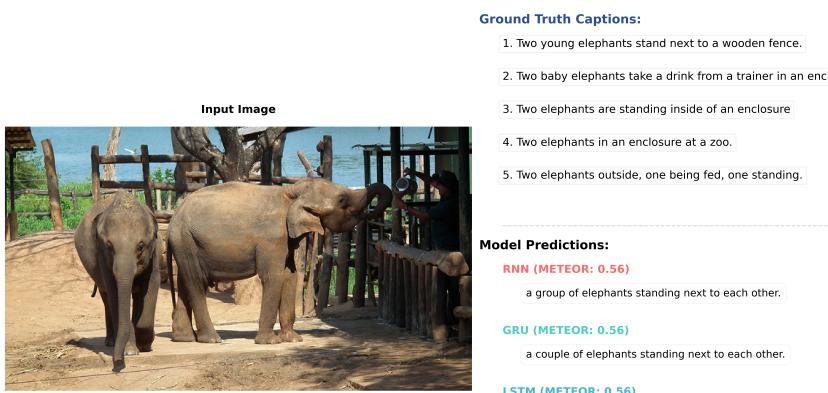
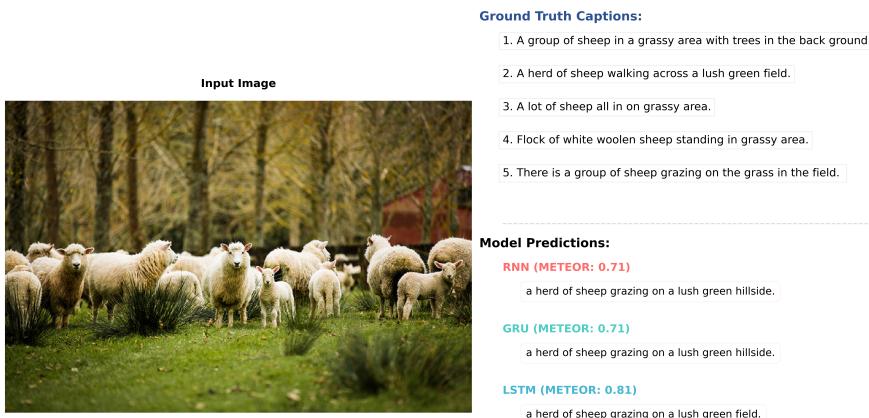
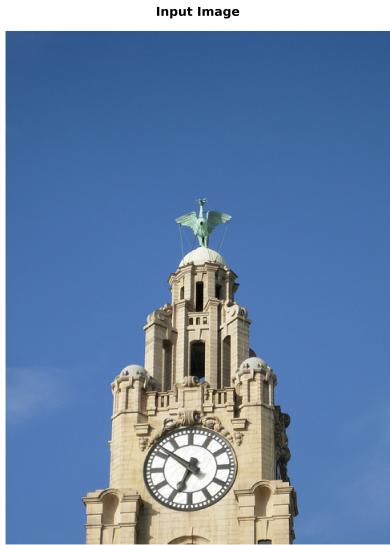


Figure 7: Sample outputs (2/4)



Ground Truth Captions:

1. A tall clock tower with a statue on top.
2. There is a clock in the top of a tall tower
3. A large clock tower with a gargoyle atop sits in front of a clear blue sky.
4. A large clock tower with a statue on the top.
5. Clock tower with a bronze statue on top on a sunny day.

Model Predictions:

RNN (METEOR: 0.63)

a large clock tower with a clock on its side.

GRU (METEOR: 0.69)

a tall clock tower with a clock on it's side.

LSTM (METEOR: 0.70)

a tall clock tower with a clock on it.



Ground Truth Captions:

1. A great shot of a very nice and large city somewhere.
2. A river is seen from afar and the river is next to a city that features old buildings.
3. a tall tower with a clock on top with buildings near by
4. Big Ben towering over the city of London England
5. An aerial view of the city of London on a cloudy day.

Model Predictions:

RNN (METEOR: 0.17)

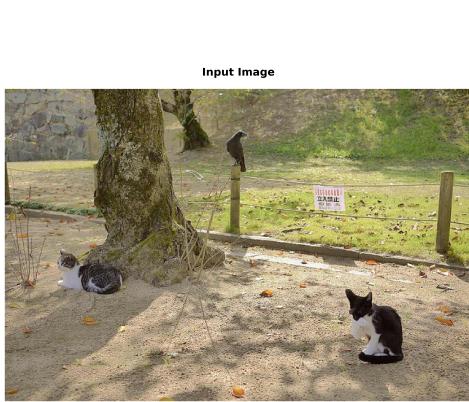
a group of people standing on a dock near a river.

GRU (METEOR: 0.31)

a large passenger jet sitting on top of a runway.

LSTM (METEOR: 0.13)

a train traveling through a rural country side.



Ground Truth Captions:

1. Two cats sitting in the dirt under a tree.
2. A couple of cats are by a tree
3. an image of two cats in a sand lot
4. Two cats by a tree in the dirt
5. Several cats sitting near a tree and a bird on a fence post.

Model Predictions:

RNN (METEOR: 0.27)

a brown and white dog standing next to a fence.

GRU (METEOR: 0.31)

a dog is standing on a bench near a tree.

LSTM (METEOR: 0.21)

a dog standing on a leash in a park.

Figure 8: Sample outputs (3/4)

Input Image



Ground Truth Captions:

1. A cross country skier on the slopes on a sunny day.
2. A man riding skis across a snow covered slope.
3. A skier about to head down the slope.
4. A skier posing on a snowy ski slope.
5. a man on skis standing on a snowy terrain

Model Predictions:

RNN (METEOR: 0.78)
a person riding skis down a snow covered slope.

GRU (METEOR: 0.90)
a man riding skis down a snow covered slope.

LSTM (METEOR: 0.66)
a person riding a snow board on a snow covered slope.

Input Image



Ground Truth Captions:

1. The woman sits on the curb writing on a notepad near a fire hydrant.
2. There is a woman sitting down with a notepad in her hand.
3. a woman sitting next to a fire hydrant with a notebook in hand
4. A young woman sitting on a curb next to a fire hydrant writing on a notepad.
5. there is a woman sitting outside writing and drinking coffee

Model Predictions:

RNN (METEOR: 0.29)
a woman holding a cell phone while standing in a room.

GRU (METEOR: 0.60)
a woman sitting on a bench next to a man with a hat on top.

LSTM (METEOR: 0.29)
a woman holding a pink umbrella over her head.

Input Image



Ground Truth Captions:

1. A woman and a man sitting on a couch next to each other.
2. ai man and a woman playing with a wii
3. A man and woman sitting on a couch playing a game.
4. A couple sitting on a couch holding Wii remotes.
5. Two people sitting together on a couch playing wii

Model Predictions:

RNN (METEOR: 0.64)
a man sitting on a couch with a laptop.

GRU (METEOR: 0.27)
a group of people sitting around a table with a laptop.

LSTM (METEOR: 0.27)
a group of people sitting around a table with a laptop.

Figure 9: Sample outputs (4/4)