



Caption Generation

Trent Kindvall
DTSA 5511 Final Presentation



Caption Generation

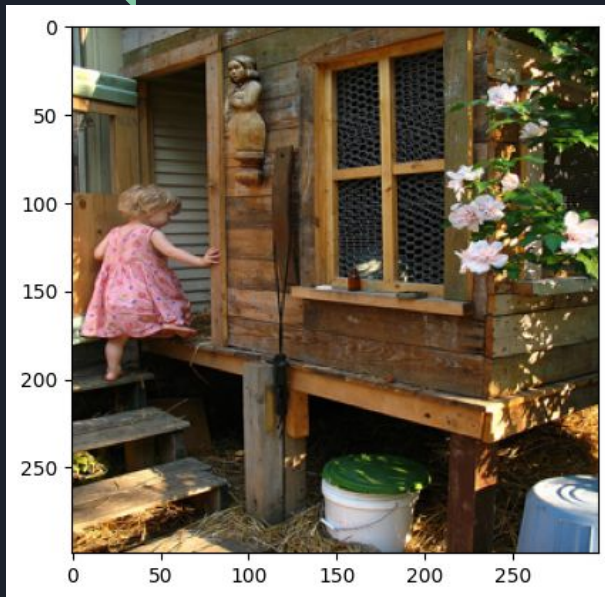
- Computer vision
- Natural Language Processing (NLP)
- Can assist people with vision impairment



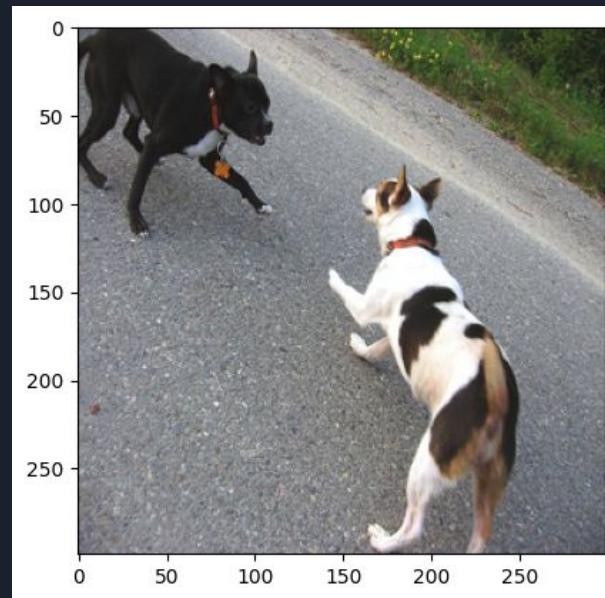
Flickr 8k Dataset

- ~8000 images with 5 captions describing the image

Images



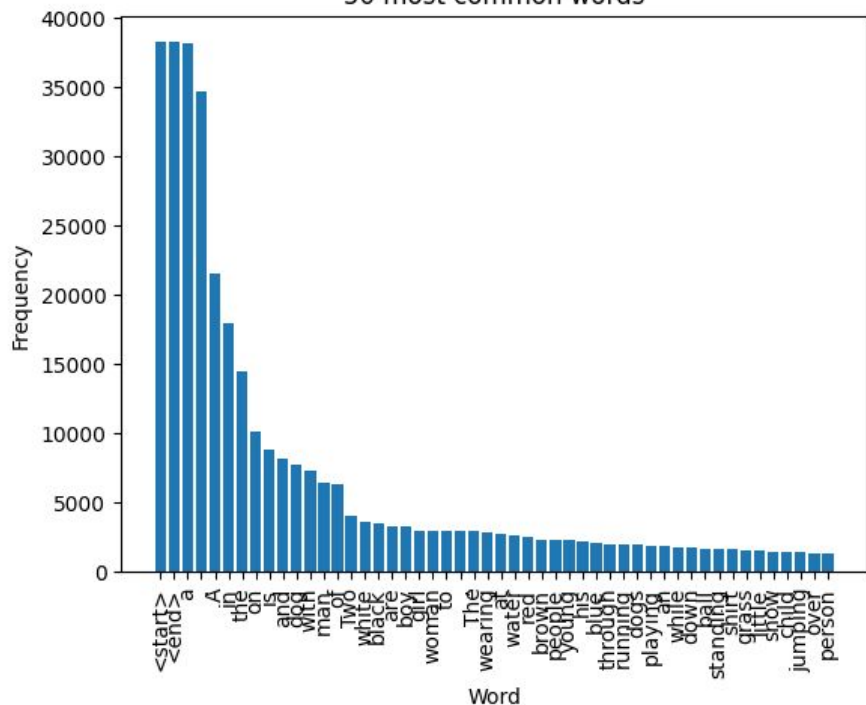
Caption:
A girl going into a wooden building .



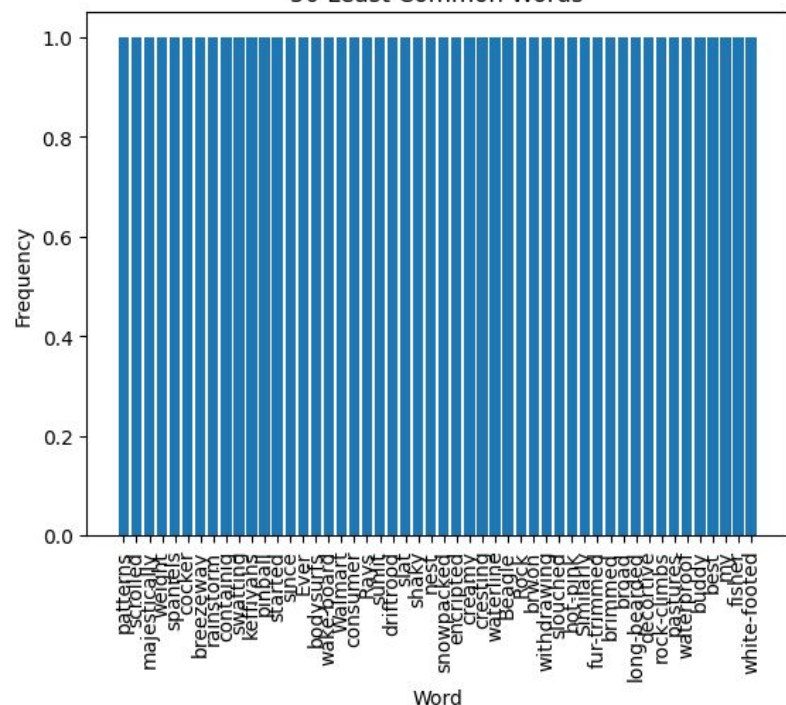
Caption:
A black dog and a tri-colored dog
playing with each other on the road .

Captions

50 most common words



50 Least Common Words

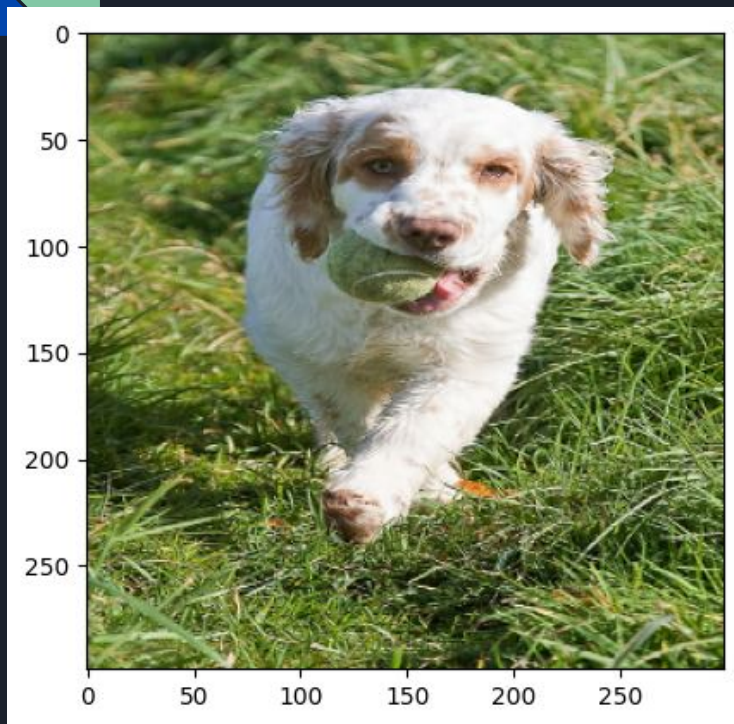




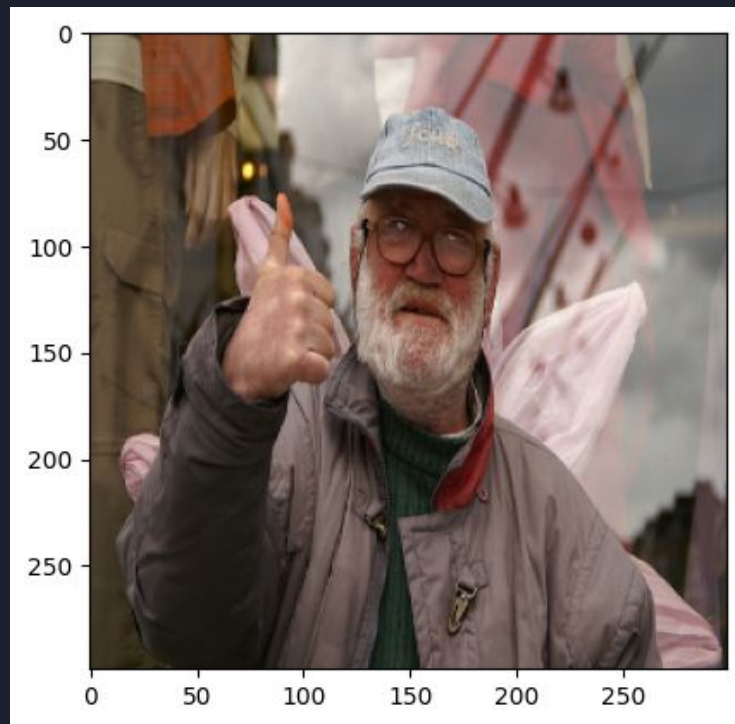
Model Architecture

- CNN
 - Image processing
 - Transfer learning
 - EfficientNet-b0
 - Can categorize images into a 1000 different categories
- Transformer Encoder
 - Translator between the cnn and the decoder
- Transformer Decoder
 - Generates the caption

Results

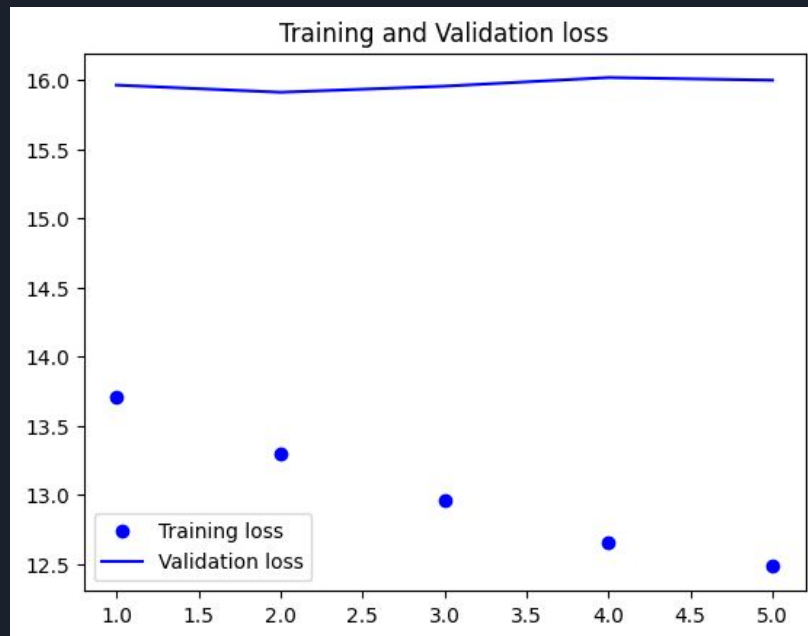
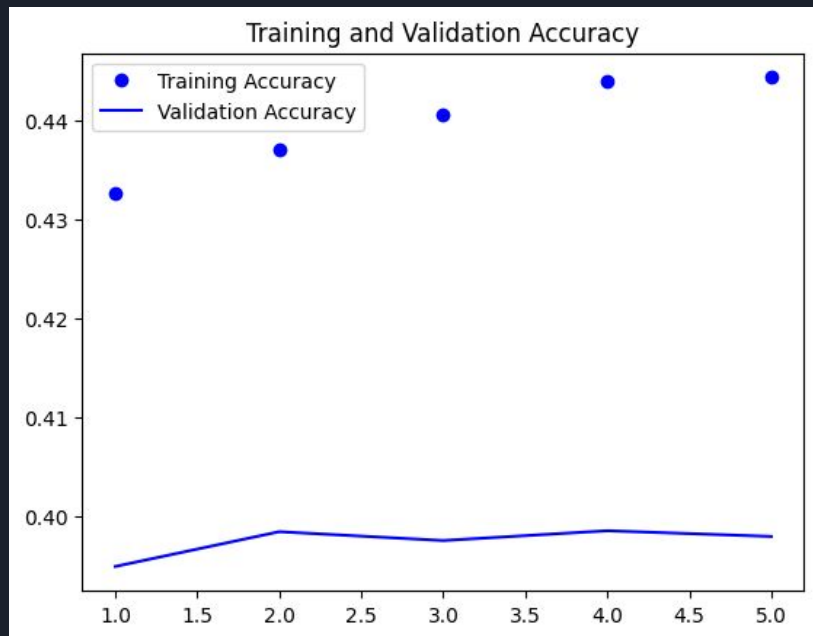


Predicted Caption: a white dog is running through the grass



Predicted Caption: a man wearing a hat and sunglasses is holding a sign

Results





BLEU Score

$$\text{Bleu}(N) = \text{Brevity Penalty} \cdot \text{Geometric Average Precision Scores}(N)$$

I achieved a score of .59 which was higher than a lot of the other attempts I saw on Kaggle



Conclusion

- Good results
- Fun use of deep neural networks and generative applications