

Attention based image to caption using a Transformer based network

Willem Diepeveen, Niek van der Laan, Daphne van Tetering, Paul Verkooijen, Kasper Wendel

TU Delft

Introduction

The canonical method for sequential modelling is the Recurrent Neural Network (RNN). In recent years, other architectures have been proposed and were shown to perform better in several cases. Especially attention to model recurrence got more popular after the paper by Vaswani et al. [1] that proposes the Transformer network. The Transformer network is something special, as it does not use any recurrences nor convolutions making it very attractive due to its high degree of parallelism. The network is purely built up from attention mechanisms and fully connected layers. Although these networks have been tested against the accepted RNN variations, these transformer networks are not the standard yet.

Research objective

Xu et al. [2] used a VGGNet and a LSTM to make an algorithm that turns an image into a caption in their paper *Show, Attend and Tell*. Inspired by this image-to-caption task, our main objective is:

- Implement a trainable VGGNet-Transformer combination for the image-to-caption task;
- Compare the performance with the VGGNet-LSTM combination.

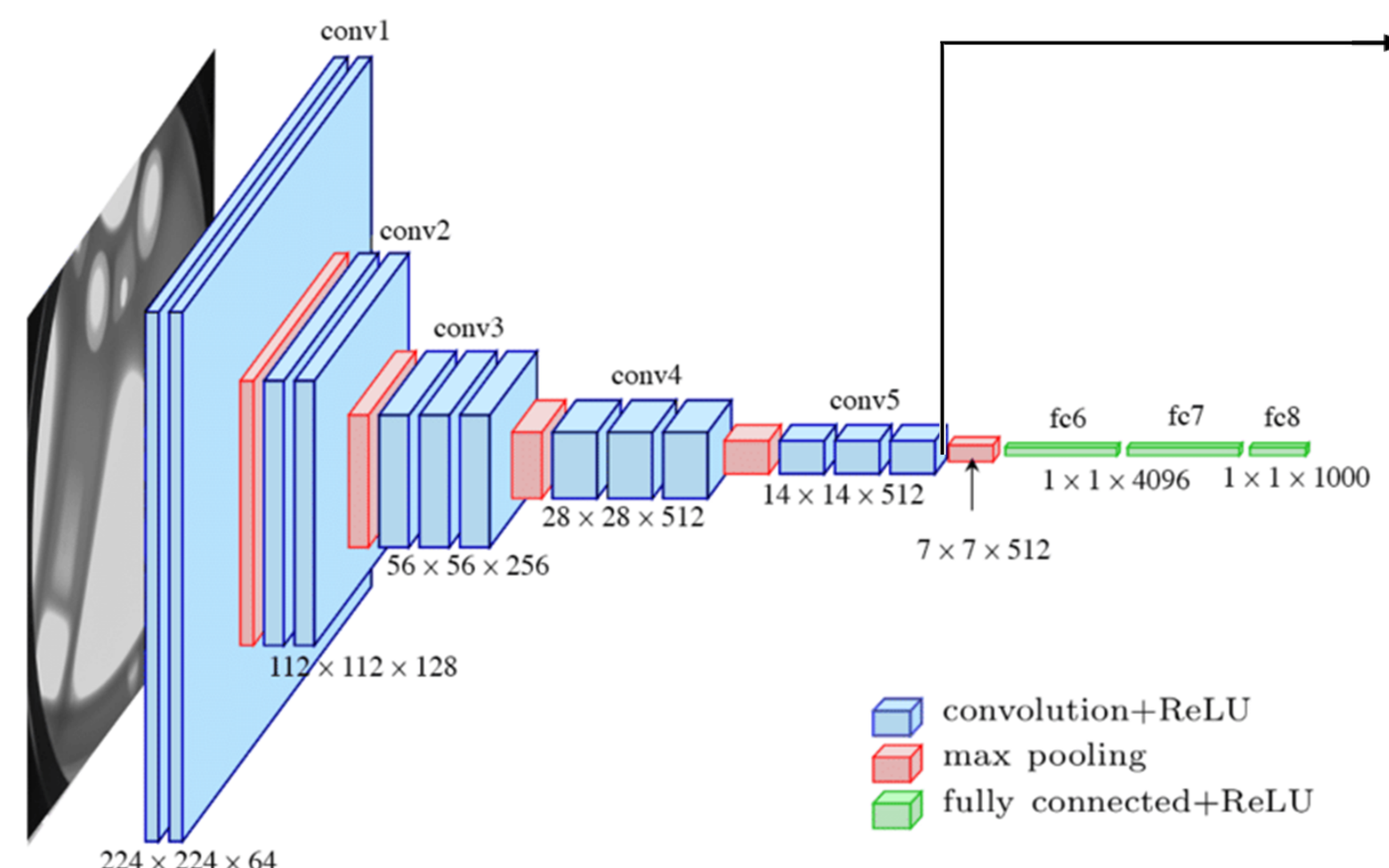


Figure 1: The feature maps obtained after the 5th convolutional layer will be flattened to get the shape 512×196 used as the VGG output of the proposed network.

Model

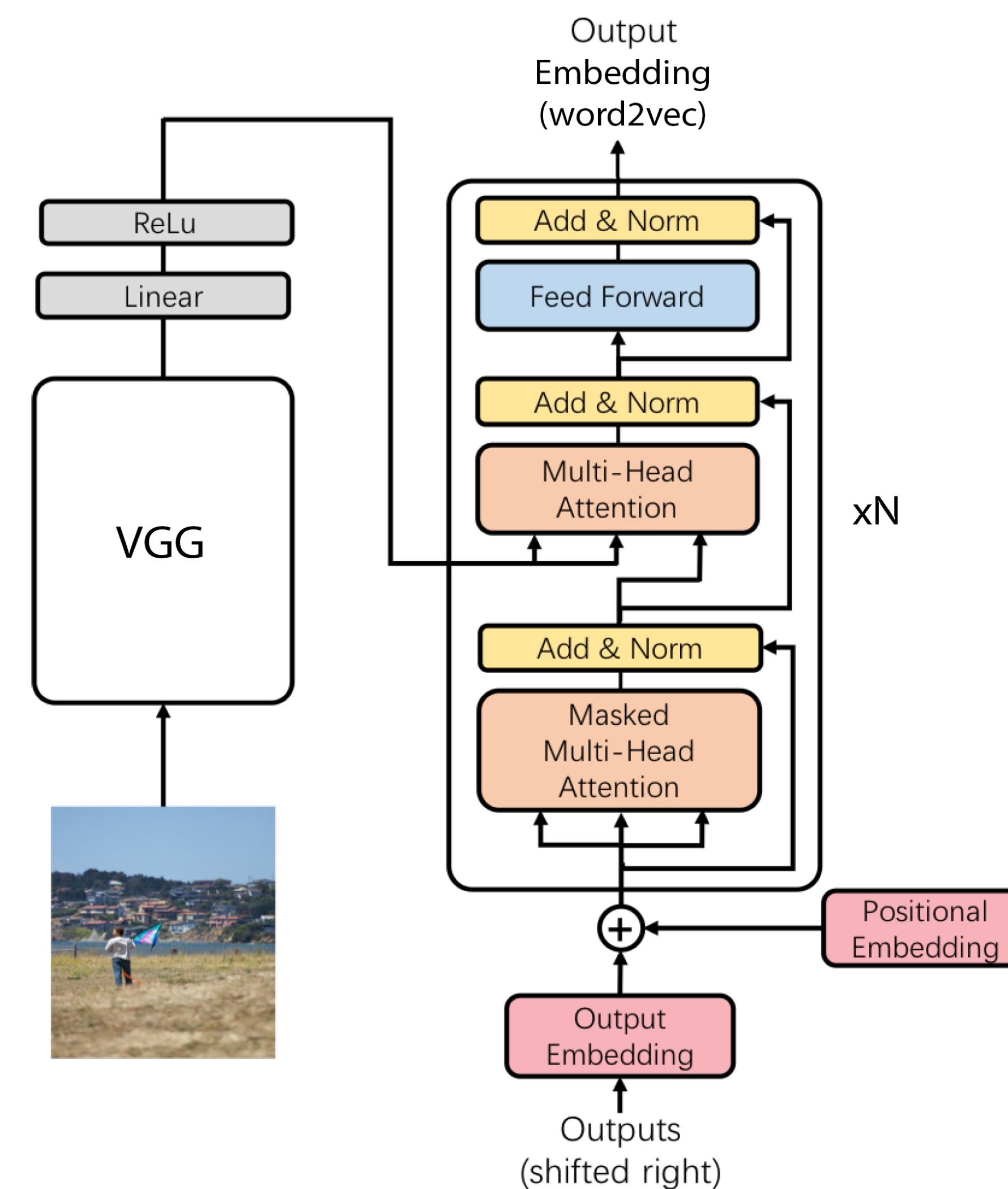


Figure 2: The proposed network inspired by the network of Zhu et al [3]. The VGG will output an array size 512×196 . The fully connected layer will reshape the image as input for the transformer. The Transformer initiates with the $\langle \text{start} \rangle$ vector and keeps predicting the next word until it guesses $\langle \text{stop} \rangle$.

Results

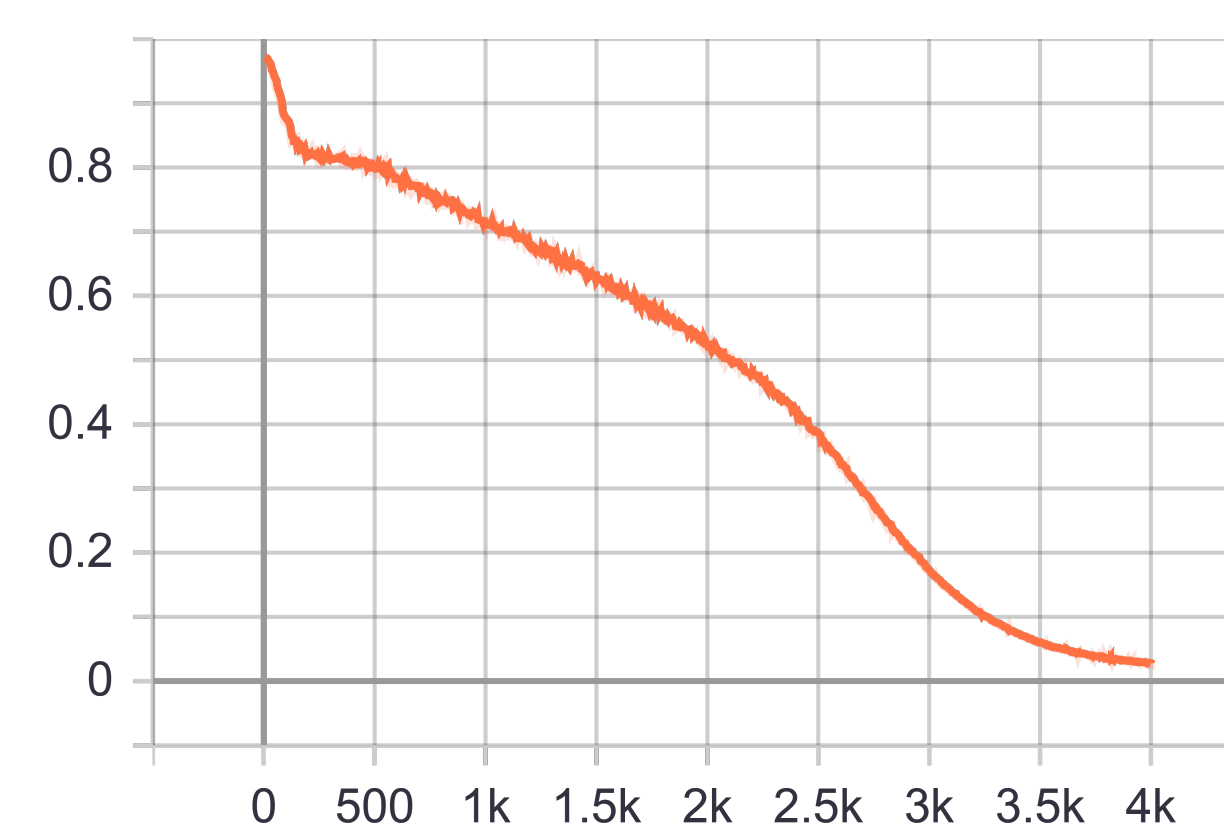


Figure 3: The cosine distance loss for the dev set.

$k = 2024$

Real: START man at use tire shop have white beard END
Pred: pail man boy boy girl girl girl girl boy boy boy boy boy
boy boy boy boy boy boy baby baby my you you you
you pack pack walk walk walk walk walk walk walk walk
sit sit sit sit sit sit sit sit sit sit

$k = 2676$

Real: START man in black leath jacke be sleep in subwa car
END
Pred: START (100x)

Method

We tried to **overfit** on a **dev** set for **trainability** and **compare** on a **full** set for **performance** testing.

1 Data

- Sampled Flickr8 dev (training) set of size 100
- Full Flickr8 dataset using a training, test and validation set with official split (6000/1000/1000) for performance testing

2 Training

- Hardware - Google Cloud CPU with capacity 7.5 GB & NVIDIA Tesla V100 GPU.
- Optimizer - Adam with Noam scheme learning rate schedule
- Regularization - Dropout
- Loss - Cosine distance loss

3 Validation

- BLEU scores
- SPICE scores

Conclusion

We implemented a trainable Transformer-VGGNet combination. The major drawback is probably the choice for embedding the special tokens $\langle \text{start} \rangle$, $\langle \text{stop} \rangle$ and $\langle \text{pad} \rangle$. We proposed a solution to this problem, which has not yet been tested.

Discussion

We overfit on the test set for $k > 2500$, but the result before that is not very good. We observe that $\langle \text{stop} \rangle$ is never guessed and the algorithm overfits by guessing $\langle \text{start} \rangle$. Although this is not desired, we do know what happens:

- Most of the sentences are short and most of the training input is padded
- For Word2Vec $\langle \text{pad} \rangle$ is non-zero and has similarity
- Guessing $\langle \text{start} \rangle$ has cosine distance $1 - 51/52 \approx 0.02$ with $\langle \text{pad} \rangle$, which is the loss we observe

Proposed solution: different embeddings for $\langle \text{start} \rangle$, $\langle \text{stop} \rangle$ and $\langle \text{pad} \rangle$ promoting dissimilarity

- $\langle \text{start} \rangle$ as $[\text{rand}([-1, 1])_{50}, 1, 0]$
- $\langle \text{stop} \rangle$ as $[\text{rand}([-1, 1])_{50}, 0, 1]$
- $\langle \text{pad} \rangle$ as $[\text{rand}([-1, 1])_{50}, -1, -1]$

References

- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in Neural Information Processing Systems*, pages 5998–6008, 2017.
- Kelvin Xu, Jimmy Ba, Ryan Kiros, Kyunghyun Cho, Aaron Courville, Ruslan Salakhudinov, Rich Zemel, and Yoshua Bengio. Show, attend and tell: Neural image caption generation with visual attention. In *International conference on machine learning*, pages 2048–2057, 2015.
- Xinxin Zhu, Lixiang Li, Jing Liu, Haipeng Peng, and Xinxin Niu. Captioning transformer with stacked attention modules. *Applied Sciences*, 8(5):739, 2018.