

Joshua Lumpkin

Training setup:

- Epochs = 50
- Learning rate = .001
- Batch size = 10

```
S2VTModel(
    (encoder_lstm): LSTM(4096, 500, num_layers=2, batch_first=True, dropout=0.5)
    (decoder_lstm): LSTM(500, 500, num_layers=2, batch_first=True, dropout=0.5)
    (fc): Linear(in_features=500, out_features=3529, bias=True)
)
```

Results:

I am very displeased with the results. I could not get it to not be fixated on singular terms. It seemed to always converge on a similar sentence structure, like what is shown in the picture. I attempted 5-6 model variations but could not get past this. Some variations I tried: model with attention, singular feature to caption, extensive training with all captions per video, simple models, stepLR. I did not have time to implement beam search.

[illegible]

Example video features shape: (80, 4096)

Epoch: 1

[illegible][illegible]

Epoch [1/30], Loss: 7.303551539130833

Example video features shape: (80, 4096)

Epoch [2/30], Loss: 5.13406401095183

Example video features shape: (80, 4096)

Epoch [3/30], Loss: 4.891730940860251

Example video features shape: (80, 4096)

Epoch [4/30], Loss: 4.754633126051529

Example video features shape: (80, 4096)

Epoch [5/30], Loss: 4.623692668002585

Epoch: 6

Predicted

[illegible][illegible]

Example video features shape: (80, 4096)

Epoch [6/30], Loss: 4.595051039820132