# LEARNING GRAMMAR OF COMPLEX ACTIVITIES VIA DEEP NEURAL NETWORKS

#### Mashaido Becky

Khoury College of Computer Sciences Northeastern University Boston, MA 02115 mashaido.r@northeastern.edu

January 11, 2021

#### **ABSTRACT**

Motivated by the growing amount of publicly available video data on online streaming services and an increased interest in applications that analyze continuous video streams such as autonomous driving [1], this technical report provides a theoretical insight into deep neural networks for video learning, under label constraints. I build upon previous work in video learning for computer vision, make observations on model performance and propose further mechanisms to help improve our observations.

#### 1 Introduction and Motivation

#### 1.1 What is the problem?

Video learning is increasingly becoming a huge part of automating computer vision tasks however with the growing amount of publicly available video data, comes the problem of unstructured video data. I address this problem by building upon a deep neural network architecture (proposed by Zijia Lu¹) that predicts the task of a video and actions happening in the video. Similar to multi-instance learning, video-level labels here are accessible, yet accurate action predictions require finding the frame-level labels. It is therefore relied upon this model to find the correspondence between the video-level and frame-level labels.

#### 1.2 Why is it interesting and important? Why is it hard? (E.g., why do naive approaches fail?)

This task is applicable to many other topics such as weakly supervised action segmentation. Weakly supervised action segmentation in test time, must analyse the relation between a test video with all possible action sequences from all tasks. This is computationally complex and prone to error. The model in this project is therefore proposed to reduce the complexity by predicting the most relevant tasks and narrowing down the range of possible action sequences. To this end, we experiment with different network architecture, motion features and visual features, and different training strategies.

# 1.3 Why hasn't it been solved before? (Or, what's wrong with previous proposed solutions? How does mine differ?)

In addition to points laid out in *section 1.2*, past research has explored the use of a NeuralNetwork-Viterbi [1] algorithm as shown in *Figure 1*, to generate frame labels. In this project however, our baseline model (Zijia Lu) depends on the attention mechanism, as shown in *Figure 2*.

<sup>&</sup>lt;sup>1</sup>Zijia Lu: https://www.khoury.northeastern.edu/people/zijia-lu/

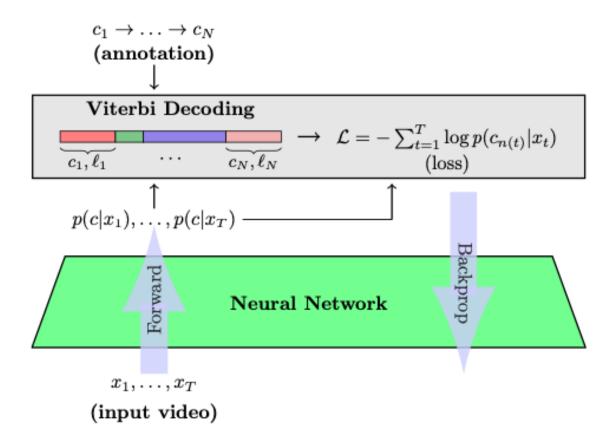


Figure 1: Input video is forwarded through the network and the Viterbi decoding is run on the output probabilities. The frame labels generated by the Viterbi algorithm are then used to compute a framewise cross-entropy loss based on which the network gradient is computed.

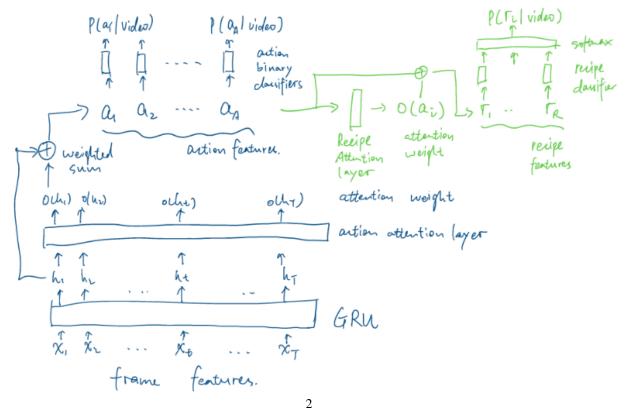


Figure 2: Our architecture for the recipe action classifier, which uses an attention mechanism.

#### 1.4 What are the key components of my approach and results? Also include any specific limitations.

- 1. Baseline model consists of a gated recurrent unit (GRU) and an attention (fully connected) layer as shown in *Figure 2*
- 2. Video frames are fed into the GRU layer as features, and consequently through the action attention layer
- 3. A weighted sum is generated and resulting action features later passed through the recipe attention layer
- 4. We also generate a temporal attention vector as an output of the action attention layer
- 5. We compute a softmax to turn the recipe logits (from the recipe attention layer) into probabilities
- 6. We use the temporal attention vetor from *part 4* to analyze performance of our attention mechanism, as detailed in the *results section 2.3*
- 7. We train and test our model on two dataset-features outlined in the *dataset section 2.1* and with each experiment, make observations on our accuracy
- 8. We try different methods to improve the accuracy of our model and note down our observations
- 9. Lastly, we make recommendations for next steps, based on our observations

Our model performance did not seem to improve. We will get into this, by the end of this paper.

#### 2 Experiments and Results

#### 2.1 Dataset

This project involves working on the breakfast dataset which consist of diverse tasks:

 48 actions and 10 recipes related to breakfast preparation, performed by 52 different individuals in 18 different kitchens

We work with two sets of features:

- Motion: Improved Dense Trajectory + PCA; I3D Motion Feature
- Visual: I3D Visual Feature; pretrained ResNet Feature

#### 2.2 Experimental Set-Up

For all experiments, code was written in pytorch and numpy. A Neural network consisting of 1 hidden layer with 64 hidden neurons was used. Refer to *Figure 3* for our network. We train using 10,000 epochs and a batch size of 16, to help with CUDA memory issues, though one could change these commands as seen in *Figure 3*. Our network.py is where we set up our model and train.py is where we train our model. We then save the best checkpoint to test our model later.

#### 2.3 Results and Analysis

During training we observed that our model overfits therefore we implemented various regularization techniques including dropout and weight decay, with results shown in *Figures 4* and 5. We added these to the GRU layer and made sure to turn off dropout during testing. We also introduced our second set of features outlined in the *dataset section 2.1* and performed experiments on both old (motion) and new (visual) set of features, as shown in *Figures 4* and 5. We generated accuracy and loss plots to help visualize the performance of our model as shown in the sample *Figures 7* and 8.

```
batch size: 16
data: breakfast
epoch: 10,000
exp: split1/recipeWeight1.000_firstrun
gpu: 0
hidden_size: 64
lr: 0.001
lr_decay_iter: 2500
new_feature = False
no_gru: False
recipe_weight: 1.0
resume: None
split: 1
Number of training data 1460
Net ( (primary_gru): GRU(64, 32, batch_first=True, bidirectional=True)
     (action_attention_layer): Linear(in_features=64, out_features=48, bias=True))
```

Figure 3: Training commands and model set-up. Network consists of a gru and attention layer, with 64 hidden neurons.

Results on Motion Features						
Experiment Run	Experiment Setting	Test Recipe Accuracy	Test Action Accuracy	Test Action F1		
with dropout	rate $p = 50\%$	0.466	0.909	0.27		
with dropout	rate $p = 40\%$	0.471	0.906	0.269		
with dropout	rate $p = 60\%$	0.487	0.914	0.295		
with weight decay i.e L2 regularization	value = 0.01	0.413	0.905	0.132		

Figure 4: Results: baseline and regularized model on motion (old set of) features.

Results on Visual Features						
Experiment Run	Experiment Setting	Test Recipe Accuracy	Test Action Accuracy	Test Action F1		
with dropout	rate p = 50%	0.551	0.918	0.299		
with dropout	rate $p = 40\%$	0.521	0.916	0.29		
with dropout	rate $p = 60\%$	0.497	0.917	0.28		
with weight decay	value = $0.01$	0.469	0.922	0.202		
baseline without any- thing	N/A	0.508	0.917	0.316		
with weight decay	value = $0.1$	0.19	0.907	0		
with weight decay	value = 0.001	0.512	0.919	0.311		
with weight decay	value = $0.0001$	0.504	0.919	0.314		
combining best dropout and best weight decay	P = 50% and value = 0.0001	0.52	0.916	0.314		

Figure 5: Results: baseline and regularized model on visual (new set of) features.

After implementing regularization techniques we still encountered overfitting issues and therefore decided to take a look under the hood i.e fully understand our attention mechanism. Here, we analyzed our attention mechanism as documented on attention.py and visualized our attention distribution for each test video as shown in the sample *Figure 9*. We also compute a quantitative score to help measure the quality of our attention as outlined in attention.py, using *Figure 6* formula below.

$$Score_v = \frac{1}{A} \sum_{i=1}^{A} \frac{1}{T} \cdot l_i \cdot P(a_i)$$
 (1)

Figure 6: Computing a quantitative score to measure the quality of attention for each test video, using the one-hot form for the ground truth label in each frame and the corresponding action attention vector as mentioned in section 1.5.  $P(a_i)$  is the attention vector for action  $a_i$  and  $l_i$  is the one-hot label for action  $a_i$ .

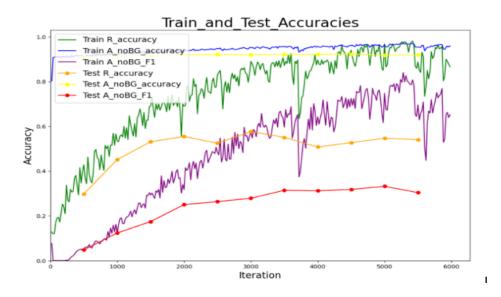
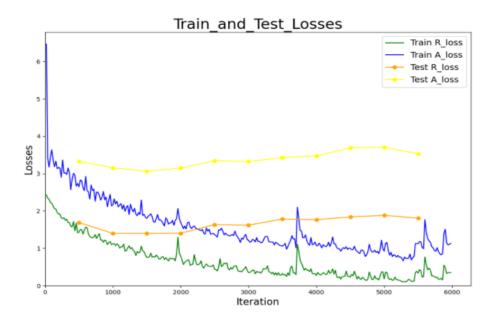


Figure 7: Visualizations: accuracy on regularized model (best weight decay and best dropout values) on new set of features.



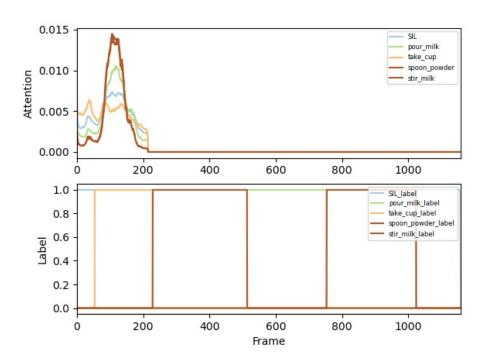


Figure 8: Visualizations: loss on regularized model (best weight decay and best dropout values) on new set of features.

Figure 9: Visualizations: attention mechanism on regularized model (best weight decay and best dropout values) on new set of features.

Gathering the above experimental results, we observe the following:

- 1. Initially when we had the subplots from *Figure 9* on one plot, the magnitude of our attention was surprisingly small and illegible
- 2. After separating into two subplots (as shown in *Figure 9*) our attention (subplot) mechanism seemed to concentrate only on the first part of a given test video, thus missing a lot of information on the rest of the video
- 3. Furthermore, different actions in a given test video seemed to correlate with each other as shown in the subplot attention *Figure 9*, where the attention is both high and low within the same frame. This means that our attention focuses on the same location in a test video

### 3 Conclusion and Future Directions

Based on our results and analysis, we propose the following next steps to help improve our model:

- Change our network structure and/or add a penalization loss term
- Penalize the attention such that it does not focus on the same video frame

Video learning under label constraints is a challenging task but with the huge traction it has gained in computer vision tasks, very promising indeed.

## References

[1] Alexander Richard, Hilde Kuehne, Ahsan Iqbal and Jurgen Gal. "NeuralNetwork-Virtebi: A Framework for Weakly Supervised Video Learning" *arXiv:1805.06875v1[cs.CV] 17 May 2018*