
Learning Differentiable Grammars for Continuous Data

AJ Piergiovanni¹ Anelia Angelova¹ Michael S. Ryoo¹

Abstract

This paper proposes a novel algorithm which learns a formal regular grammar from real-world continuous data, such as videos or other streaming data. Learning latent terminals, non-terminals, and productions rules directly from streaming data allows the construction of a generative model capturing sequential structures with multiple possibilities. Our model is fully differentiable, and provides easily interpretable results which are important in order to understand the learned structures. It outperforms the state-of-the-art on several challenging datasets and is more accurate for forecasting future activities in videos. We plan to open-source the code.

1. Introduction

Learning a formal grammar from continuous, unstructured data is a challenging problem. This is especially challenging when the elements (i.e., terminals) of the grammar to be learned are not symbolic or discrete (Chomsky, 1956; 1959), but are higher dimensional vectors, such as representations from real world data streams (e.g., videos).

Simultaneously, addressing such challenges is necessary for better automated understanding of streaming data. In video understanding, such as activity detection, a space-time convolutional neural network (CNN) (e.g., (Carreira & Zisserman, 2017)) generates a representation abstracting local spatio-temporal information at every time step, forming a temporal sequence of representations. Learning a grammar reflecting sequential changes in video representations will enable explicit and high-level modeling of temporal structure and relationships between multiple occurring events in videos. This not only allows for better recognition of human activities from videos by enforcing the learned grammar to local-level detections, but also enables forecasting of future representations based on the learned production rules. It also provides semantic interpretability of the video

recognition and prediction process.

In this paper, we propose a new approach of modeling a formal grammar in terms of learnable and differentiable neural network functions. The objective is to formulate not only the terminals and non-terminals of our grammar as learnable representations but also the production rules generating them as differentiable functions. We provide the loss function to train our differentiable grammar directly from data, and present methodologies to take advantage of it for recognizing and forecasting sequences. Rather than focusing on non-terminals and production rules to generate or parse symbolic data (e.g., text strings), our approach allows learning of grammar representations directly on top of higher-dimensional data stream (e.g., representation vector sequences). We confirm such capability experimentally by focusing on learning a differentiable regular grammar from continuous representations, which can be applied to any sequential data including outputs of 3-D space-time CNNs.

The primary contributions of our work are:

1. Design of a **fully differentiable** neural network that is able to learn the structure (terminals, non-terminals, and production rules) of a regular grammar.
2. The grammar model is easily **interpretable**, enabling understanding of the structures learned from data.
3. We confirm that the approach works on sequential real-world datasets, and outperforms the state-of-the-art on **challenging benchmarks**.
4. We show that the model is able to achieve better results on **future forecasting** of human activities which are to occur subsequently in videos.

The goal of this work is to provide to the research community a neural differentiable grammar-based matching and prediction for video analysis, which is also applicable to other domains. The results are interpretable which is very important for real-life decision making scenarios. Furthermore, it can predict with higher accuracy future events, which is crucial for anticipation and reaction to future actions, for example for an autonomous robot which interacts with humans in dynamic environments.

¹Google Brain. Correspondence to: AJ Piergiovanni <ajpiergi@google.com>, Anelia Angelova <anelia@google.com>, Michael Ryoo <mryoo@google.com>.

2. Background

A formal grammar G is defined with four elements: $G = (V, \Sigma, P, S)$ where V is a finite set of non-terminals, Σ is a finite set of terminals, P is a finite set of production rules, and S is the starting non-terminal.

In a regular grammar, the production rules P are in the following forms:

$$\begin{aligned} A &\rightarrow aB \\ A &\rightarrow a \\ A &\rightarrow \epsilon \end{aligned} \quad (1)$$

where A and B are non-terminals in V , a is any terminal in Σ , and ϵ denotes an empty string. A regular grammar is a type 3 formal grammar in the Chomsky hierarchy.

In this paper, we follow this traditional regular grammar definition, while extending it by making its terminals, non-terminals, and production rules represented in terms of differentiable neural network functions. Our differentiable grammar could be interpreted as a particular form of recurrent neural network (RNN). The main difference to the standard RNNs such as LSTMs and GRUs (Hochreiter & Schmidhuber, 1997; Cho et al., 2014) is that our grammar explicitly maintains a set of non-terminal representations (in contrast to having a single hidden representation in standard RNNs) and learns multiple distinct production rules per non-terminal. This not only makes the learned model more semantically interpretable, but also allows learning of temporal structures with multiple sequence possibilities. Our grammar, learned with a randomized production rule selection function, considers multiple transitions between abstract non-terminals when matching it with the input sequences as well as when generating multiple possible future sequences.

We also experimentally compare our grammar with previous models including LSTMs (Hochreiter & Schmidhuber, 1997) and Neural Turing Machines (NTMs) (Graves et al., 2014) in the experiments section.

3. Approach

3.1. Formulation

We model our formal grammar in terms of latent representations and differentiable functions mapping to representations. The parameters of our functions define production rules, which are learned together with the terminal and non-terminal representations.

Each non-terminal in V is a latent representation with fixed dimensionality, whose actual values are learned based on the training data. Each terminal in Σ corresponds to a video representation that could be obtained at every time step,

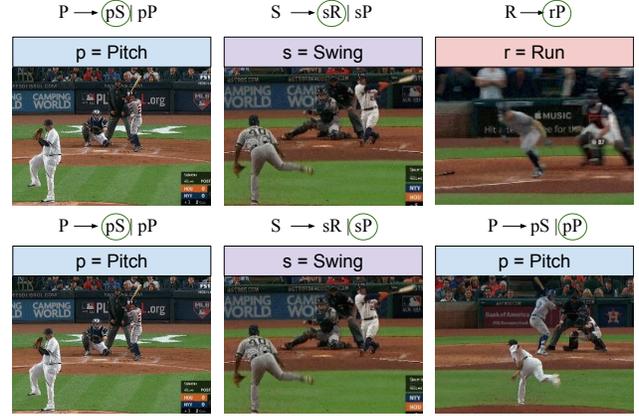


Figure 1. Example regular grammar giving the sequence of possible activities in a baseball video. For example, a swing only occurs after a pitch.

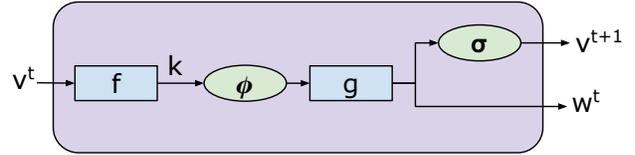


Figure 2. Illustration of the connection between functions in the grammar model. ϕ is the gumbel-softmax function and σ is the softmax function.

such as a vector with activity class predictions. This has to be learned as well. Our production rules are represented as a pair of two functions:

- f : a function that maps each non-terminal in V (e.g., A) to a subset of production rules (i.e., the rules for the current non-terminal) $\{p_i\} \subset P$.
- g : a function that maps each rule p_i to a terminal (e.g., a) and the next non-terminal (e.g., B).

$$\begin{aligned} f &: V \rightarrow \{P\} \\ g &: P \rightarrow (V, \Sigma). \end{aligned} \quad (2)$$

The combination of the two functions effectively captures multiple production rules per non-terminal, such as “ $A \rightarrow aB$ ” and “ $A \rightarrow aA$ ”. The starting non-terminal S is learned to be one of the latent representations in V . The functions are learned from data.

These form a straight forward (recursive) generative model, which starts from the starting non-terminal $S = v^0$ and iteratively generates a terminal at every time step. Representing our production rules as functions allows us to model the generation of a sequence (i.e., a string) of terminals as the

repeated application of such functions. At every time step t , let us denote the first function mapping each non-terminal to a set of production rules as $k = f(v^t; \theta_1)$, and the second function mapping each rule to a non-terminal/terminal pair as $(v^{t+1}, w^t) = g(p_i; \theta_2)$ where $v \in V$, and $w \in \Sigma$. k is a latent vector describing the production rule activations corresponding to v^t .

In its simplest form, we can make our grammar rely only on one production rule by applying the softmax function (σ) to the activation vector k : $p_i = \sigma(k)$. This formulation makes p_i a (soft) one-hot indicator vector selecting the i -th production rule. Our sequence generation then becomes:

$$(v^{t+1}, w^t) = g(\sigma(f(v^t; \theta_1)); \theta_2). \quad (3)$$

We represent each $v \in V$ as a N -dimensional soft one-hot vector where N is the number of non-terminals. In the actual implementation, this is constrained by having a softmax function as a part of g_{θ_2} to produce v^{t+1} . Each $w \in \Sigma$ is a T -dimensional representation we learn to generate, where T is the dimensionality of the sequential representation at every time step. This process is shown in Fig. 2.

We further extend Eq. 3 to make the grammar consider multiple production rules in a randomized fashion during its learning and generation. More specifically, we use the Gumbel-Softmax trick (Jang et al., 2017; Maddison et al., 2017) to replace the softmax in Eq. 3. Treating the activation vector k as a distribution over production rules, the Gumbel-Softmax (ϕ) allows sampling of different production rules:

$$(v^{t+1}, w^t) = g(\phi(f(v^t; \theta_1)); \theta_2). \quad (4)$$

In our case, this means that we are learning the grammar production rules which could be selected/sampled differently even for the same non-terminal (i.e., v^t) while still maintaining a differentiable process.

The idea behind our grammar formulation is to allow direct training of the parameters governing generation of the terminals (e.g., video representations in our case), while representing the process in terms of explicit (differentiable) production rules. This is in contrast to traditional work that attempted to extract grammar from already-trained standard RNNs (Gers & Schmidhuber, 2001) or more recent neural parsing works using discrete operators (Dyer et al., 2016) and memory-based RNNs (Graves et al., 2014). Our formulation also adds interpretability/explainability to our temporal models learned from data streams, as we confirm more in the following subsections.

Detailed implementation of production rule functions: Although any other differentiable functions could be used for modeling our functions f and g , we use matrix operations to implement them. Given a matrix of production

rules, W , a $N \times (R \cdot N)$ matrix, where R is the maximum number of production rules per non-terminal, we obtain the activation vector k with size $R \cdot N$ as:

$$k = f(v) = vW \quad (5)$$

We constrain W so that its each column is a vector with only one non-zero element (i.e., each production rule may originate from only one non-terminal). In the actual implementation, W is obtained by modeling it as a $N \times R$ matrix and then inflating it with zeros to have the form of a block diagonal matrix of size $N \times (R \cdot N)$ with the block size $1 \times R$.

Similarly, the function g mapping each production rule to the next non-terminal and corresponding terminal is implemented using a $(R \cdot N) \times N$ matrix H_1 , and a $(R \cdot N) \times T$ matrix H_2 :

$$(v^{t+1}, w^t) = g(v^t) = (\sigma(FH_1), FH_2) \quad (6)$$

where $F = \phi(f(v^t))$. With this implementation, learning the grammar production rules is done by learning the matrices W , H_1 , and H_2 directly. Figure 4 describes an example.

3.2. Learning

We train our grammar model to minimize the following binary cross entropy loss:

$$\mathcal{L} = \sum_{t,c} z_c^t \log(w_c^t) + (1 - z_c^t) \log(1 - w_c^t) \quad (7)$$

where z^t is the ground truth label vector at time t with dimensionality $|c|$ and w^t is the output of the grammar model (terminal). In the case where the grammar is used to predict discrete class labels, z^t becomes a one-hot vector. Training of our functions f and g (or matrices W , H_1 , and H_2) can be done with a straight forward backpropagation for the simple production rule case of Eq. 3, as it becomes a deterministic function per non-terminal at each time step. Backpropagating through the entire sequential application of our functions also allow learning of the starting non-terminal representation $S = v^0$.

Learning multiple production rules: In general, our function f maps a non-terminal to a ‘set’ of production rules where different rules could be equally valid. This means that we are required to train the model by generating many sequences, by taking b rules at each step (b is the branching factor).

We enumerate through multiple productions rules by randomizing the production rule selection by using the Gumbel-Softmax trick (Jang et al., 2017; Maddison et al., 2017) as suggested in the above subsection. This allows for weighted

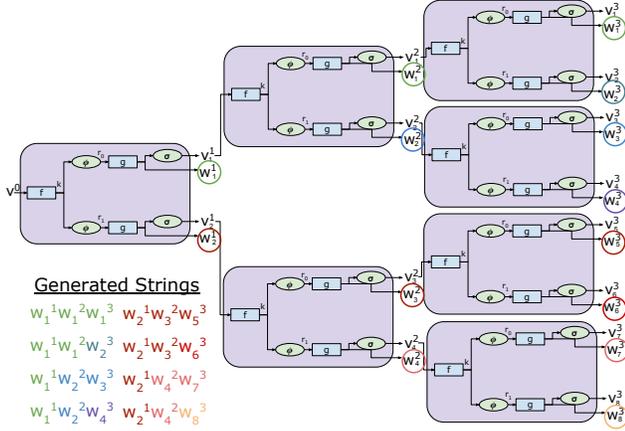


Figure 3. Visualization of training the grammar with branching. The output (r_i) of the Gumbel-Softmax (ϕ) is different for each branch, producing different strings.

Algorithm 1 The training of the grammar, with multiple branches

Input: sequence s
 Set initial nonterminal v^0
for $t = 0$ **to** T **do**
 for $c = 0$ **to** current total branches **do**
 Get rules for current nonterminal: $k = f(v^t)$
 for $b = 0$ **to** Number of branches **do**
 Randomly select a rule: $p = \phi(k)$
 Get next non-terminal and terminal
 $(v_b^{t+1}, w_b^t) = g_{\theta_2}(p)$
 end for
 end for
end for
 $loss = \min_b \mathcal{L}(s, w_b)$, min over all branches

random selection of the rules based on the learned rule probabilities. In order to train our grammar model with the Gumbel-Softmax, we maintain multiple different ‘branches’ of non-terminal selections and terminal generations, and measure the loss by considering all of them. Algo. 1 and Figure 3 illustrate the training and branching process. When generating many branches, we compute the loss for each generated sequence, then take the minimum loss over the b branches, effectively choosing the branch that generated the most similar string:

$$\mathcal{L} = \min_b \sum_{t,c} z_{t,c} \log w_{b,c}^t + (1 - z_{t,c}) \log(1 - w_{b,c}^t) \quad (8)$$

where $w_{b,c}^t$ is the output of the grammar model (terminals) at time t for class c and branch b . Branches are pruned to make the process computationally tractable, limiting the total number of branches we maintain.

3.3. Interpretability

As our model is constrained to use a finite set of non-terminals, terminals and production rules, it allows for easy interpretability of the learned grammar structure. We can conceptually convert the learned production rule matrices W , H_1 , and H_2 into a discrete set of symbolic production rules by associating symbols with the learned terminal (and non-terminal) representations. The matrix W describes the left-hand side non-terminal of the production rule following the regular grammar (e.g., $\mathbf{A} \rightarrow aB$), the matrix H_2 describes the terminal of the production rule (e.g., $A \rightarrow \mathbf{a}B$), and the matrix H_1 corresponds to the right-hand side non-terminal of the rule (e.g., $A \rightarrow a\mathbf{B}$). Element values of the matrix W in particular suggests the probability associated with the production rule (i.e., it governs the probability of the corresponding production rule being randomly selected with Gumbel-Softmax). Fig. 4 shows how we can construct a grammar from the learned matrices.

Fig. 8 illustrates examples of such interpreted grammar, learned from a raw baseball video dataset. This was done by associating symbols with w^t and v^t .

3.4. Application to video datasets

While application of our differentiable grammar learning to 1-D data is rather straightforward, when applying this to more complex continuous data with various contents such as videos, certain extensions are needed.

To apply the grammar model to videos, we make a few key changes. The initial non-terminal is learned based on the video representation. We learn a function ψ that maps from the video representation to the initial non-terminal: $S = v_0 = \psi(q)$, where q is the output of a video CNN (e.g., I3D (Carreira & Zisserman, 2017)). We then train the grammar model as above, where the ground truth is the sequence of one-hot vector based on the activity labels in the video.

During inference (which is about predicting frame-level activity labels), we generate a sequence by selecting the rule that best matches the CNN predicted classes. We then multiply the predictions from the grammar with the predictions from the CNN. To predict future, yet-unseen actions, we generate a sequence following the most likely production rules.

4. Experiments

4.1. Toy Examples

We first confirm that our model is able to learn the rules of simple, hand-crafted grammars and show how we can easily

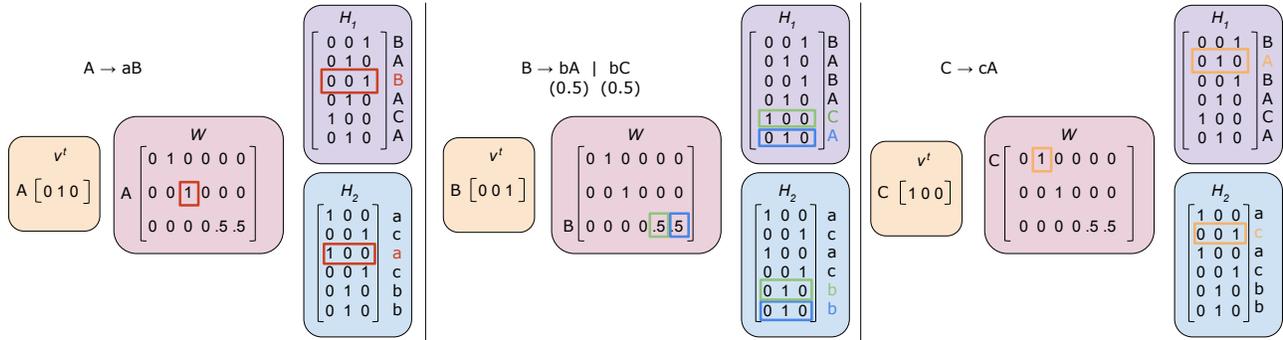


Figure 4. Visualization of the learned toy grammar and how we can construct the grammar from the learned matrices. The non-terminal v^t gives an soft-index into the rule matrix W , which gives probabilities over the rules. The rules give a soft-index into the non-terminal matrix (H_1) and terminal matrix (H_2).

interpret the learned model. Given the simple grammar:

$$\begin{aligned} A &\rightarrow aB \\ B &\rightarrow bC \mid bA \\ C &\rightarrow cA \end{aligned}$$

We train a model with 3 terminal symbols (a , b , and c), 3 non-terminal symbols (A , B and C), and 2 production rules per non-terminal. We can then examine the learned grammar structure, shown in Fig. 4. We observe that the learned starting non-terminal corresponds to ‘A’, and by following the learned rules, we end up with ‘aB’. From non-terminal ‘B’, the learned rules go to ‘bA’ or ‘bC’ with 50% probability. From non-terminal ‘C’, the learned rules go to ‘cA’. This confirms that the model is able to learn grammar rules and can easily be interpreted.

4.2. Air Pollution Timeseries Dataset

We further test the algorithm on a timeseries dataset in order to demonstrate its use to timeseries domain.

The Air Pollution prediction dataset (Liang et al., 2016) is intended to predict urban pollution levels. The data provides measurements hourly, for 24 hours a day and spans several years of sensing. This dataset contains several environmental factors as input features and measures the overall air pollution (‘PM2.5 concentration’). It contains about 43,000 examples which are split consecutively into train and test with a 50:50 ratio.

Figure 5 visualizes the prediction results for our model. As seen it is correctly approximating the actual values (a portion of the training data is also shown). We evaluate the model by measuring root mean squared error. By simply predicting the last seen value for the remaining data, we get RMSE of 36.45. Using our grammar model, we get RMSE

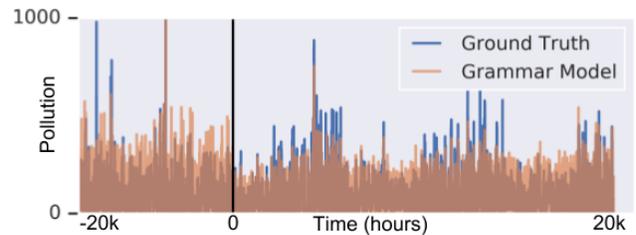


Figure 5. Results on the AirPollution timeseries data. The left of the black line is training data, right of it is unseen prediction.



Figure 6. Examples of videos in Charades dataset.

of 22.14. An LSTM-based model performs at 27.17.

4.3. Activity Detection Experiments

We further confirm that our method works on 3 real-world, challenging activity detection datasets: MLB-YouTube (Piergiovanni & Ryoo, 2018a), Charades (Sigurdsson et al., 2016b), and MultiTHUMOS (Yeung et al., 2015). All datasets are evaluated by per-frame mAP. The datasets are described as follows:

MultiTHUMOS: The MultiTHUMOS dataset (Yeung et al., 2015) is a large scale video analysis dataset which has frame-level annotations for activity recognition. It is a chal-



Figure 7. Examples of videos in MultiTHUMOS dataset.

lenging dataset and supports dense multi-class annotations (i.e. per frame), which are also used here for both prediction and ground truth. It contains 400 videos or about 30 hours of video and 65 action classes. Examples are shown in Fig. 7.

Charades: The Charades dataset (Sigurdsson et al., 2016b) is a challenging dataset with unstructured activities in videos. The videos are everyday activities in a home environment. It contains 9858 videos and spans 157 classes. Examples are shown in Fig. 6.

MLB-YouTube: The MLB-YouTube dataset (Piergiiovanni & Ryoo, 2018a) is a challenging video activity recognition dataset collected from live TV broadcast baseball games (with many challenges, such as the small resolution of activities in question). It further offers the challenge of fine-grained activity recognition as all potential activities are encountered in the same context and environment, unlike many other datasets which feature more diverse activities which may also use context for recognition. It has 4290 videos in 42 hours of video. Additionally, baseball games follow a rigid structure, making it ideal to evaluate the learned grammar. Some example frames are shown in Fig. 1.

Implementation Details We implemented our models in PyTorch. The learning rate was set to 0.1, decayed every 50 epochs by 10, and the models were trained for 400 epochs. We pruned the number of branches to 2048 by random selection.

4.4. Results on MLB-Youtube

Table 1 shows the results of the proposed algorithm on the MLB-YouTube dataset, compared to all state-of-the-art algorithms including RNNs such as LSTMs and Neural Turing Machines (NTM). We evaluated the methods in two different settings: 1) learning grammar on top of features learned from I3D and 2) on top of a recently proposed super-events method. The result clearly shows that our differentiable grammar learning is able to better capture temporal/sequential information in videos. We also compare to LSTMs and NTMs using both CNN features (e.g., I3D) as

Table 1. Results on the MLB-YouTube dataset (mAP).

| Model | mAP |
|--|-------------|
| Random | 13.4 |
| I3D | 34.2 |
| I3D + LSTM | 39.4 |
| I3D + NTM (Graves et al., 2014) | 36.8 |
| I3D class prob + LSTM | 37.4 |
| I3D class prob + NTM (Graves et al., 2014) | 36.8 |
| I3D with Grammar (ours) | 43.4 |
| I3D + super-events (Piergiiovanni & Ryoo, 2018b) | 39.1 |
| I3D + super-events with Grammar (ours) | 44.2 |

Table 2. Results on the MultiTHUMOS dataset (mAP).

| Method | mAP |
|--|-------------|
| Two-stream (Yeung et al., 2015) | 27.6 |
| Two-stream + LSTM (Yeung et al., 2015) | 28.1 |
| Multi-LSTM (Yeung et al., 2015) | 29.6 |
| Predictive-corrective (Dave et al., 2017) | 29.7 |
| I3D baseline | 29.7 |
| I3D + LSTM | 29.9 |
| I3D + NTM (Graves et al., 2014) | 29.8 |
| I3D class prob + LSTM | 29.8 |
| I3D class prob + NTM (Graves et al., 2014) | 29.7 |
| I3D with Grammar (ours) | 32.3 |
| I3D + super-events (Piergiiovanni & Ryoo, 2018b) | 36.4 |
| I3D + super-events with Grammar (ours) | 37.7 |

input and using the predicted class probabilities as input, as that is more comparable to our grammar model. We find that the use of class probabilities slightly degrades performance for LSTMs and NTMs.

4.5. Results on MultiTHUMOS

Table 2 shows results comparing two common methods with and without the proposed grammar. We also test both settings as above and compare to the state-of-the-art. In both settings we can see that use of the learned grammar outperforms previously known methods.

4.6. Results on Charades

Table 3 has results comparing the proposed grammar to other prior techniques on the Charades dataset (v1_localize setting). As seen, this dataset is quite challenging since recently its detection accuracy was below 10 percent mAP. Our results here too outperform the state-of-the-art, increasing the accuracy on this dataset to over 20 percent mAP. We note that there are consistent improvements in both

Table 3. Detection results on the Charades dataset (Charades_v1_localize setting).

| Method | mAP |
|---|-------------|
| Predictive-corrective (Dave et al., 2017) | 8.9 |
| Two-stream (Sigurdsson et al., 2016a) | 8.94 |
| Two-stream+LSTM (Sigurdsson et al., 2016a) | 9.6 |
| R-C3D (Xu et al., 2017) | 12.7 |
| Sigurdsson et al. (Sigurdsson et al., 2016a) | 12.8 |
| I3D baseline | 17.2 |
| I3D + LSTM | 18.1 |
| I3D + NTM (Graves et al., 2014) | 17.5 |
| I3D class prob + LSTM | 17.6 |
| I3D class prob + NTM (Graves et al., 2014) | 17.4 |
| I3D with Grammar (ours) | 18.5 |
| I3D + super-events (Piergiovanni & Ryoo, 2018b) | 19.4 |
| I3D + super-events with Grammar (ours) | 20.3 |

settings, similar to the results on MultiTHUMOS and MLB-YouTube. In particular, the differentiable grammar learning outperformed previous RNNs including LSTMs and NTMs.

4.7. Future Prediction/Forecasting

As our grammar model is generative, we can apply it to predict the future, unseen activities. Future prediction is important, especially for autonomous systems (e.g., robots) as they need to anticipate potential future activities to respond to. Once the grammar is learned, future sequences containing unseen activities can be generated by selecting the most probable production rule at every (future) time step.

For this experiment we consider predicting at short-term horizons (in the next 2 seconds), mid-term horizons (next 10 seconds), and more longer-term horizons (in the next 20 seconds). We compare to baselines such as random guessing, repeatedly predicting the last seen frame, and an LSTM approach (using I3D features) which has been commonly used for future frame forecasting. We evaluate these methods using per-frame mAP.

Table 4 shows the results for future prediction for the MultiTHUMOS dataset. We confirm the proposed method is more accurate at future prediction at all future horizons considered. We note that 10-20 seconds in the future is a very challenging setting to try to predict especially in the context of multi-label datasets.

Table 5 shows the results for future prediction for the Charades dataset. Here too, we can see the proposed grammar approach is more accurate at future frame prediction, with predictions at 10 seconds in the future outperforming state-of-the-art for 2 seconds only. This data is more challenging

Table 4. Future prediction on the MultiTHUMOS dataset for various time horizons.

| Method | 2 sec | 10 sec | 20 sec |
|----------------------|-------------|------------|------------|
| Random | 2.6 | 2.6 | 2.6 |
| Last frame | 6.2 | 5.8 | 2.8 |
| I3D + LSTM | 8.5 | 6.6 | 2.9 |
| I3D + Grammar (ours) | 10.4 | 8.3 | 3.5 |

Table 5. Future prediction on the Charades dataset for various time horizons.

| Method | 2 sec | 10 sec | 20 sec |
|----------------------|------------|------------|------------|
| Random | 2.4 | 2.4 | 2.4 |
| Last frame | 6.8 | 3.3 | 2.4 |
| I3D + LSTM | 6.5 | 4.6 | 2.5 |
| I3D + Grammar (ours) | 8.6 | 7.3 | 5.5 |

by itself, which makes the future prediction even harder.

4.8. Visualization of Learned Grammars

In Figure 4, we illustrate how we convert from the learned matrices to the grammar and production rules. From the training data, we know the mapping from terminal symbol to label. We can then examine the rule matrix, W and the non-terminals, H_1 to construct the rules.

We also visualize the learned grammar for the MLB-YouTube dataset, in which interestingly the typical baseball sequences are learned. Figure 8 is the conceptual visualization of the learned regular grammar. In Figure 9, we illustrate the actual learned matrices corresponding to one of the production rules. In Figure 10, we illustrate how all the learned rules are inferred from the learned matrices.

In Figure 8, we illustrate the learned grammar. We see that that the learned grammar matches the structure in a baseball game and the probabilities are similar to the observed data, confirming that our model is able to learn the correct rule structure. For example, an activity starts with a pitch which can be followed by a swing, bunt or a hit. After a hit, foul, or strike, another pitch follows. The learned grammar is illustrated with probabilities for each rule in parenthesis.

5. Related work

Chomsky grammars (Chomsky, 1956; 1959) are designed to represent functional linguistic relationships. They have found wide applications in defining programming languages, natural language understanding, and understanding of images and videos (Socher et al., 2011).

There are early works exploring extracting grammars/state machines from trained RNNs (Kolen, 1994; Bodén & Wiles, 2000; Tiño et al., 1998). Other works have attempted to

- Brendel, W., Fern, A., and Todorovic, S. Probabilistic event logic for interval-based event recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2011.
- Carreira, J. and Zisserman, A. Quo vadis, action recognition? a new model and the kinetics dataset. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.
- Chen, D. and Manning, C. A fast and accurate dependency parser using neural networks. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, pp. 740–750, 2014.
- Cho, K., van Merriënboer, B., Gulcehre, C., Bahdanau, D., Bougares, F., Schwenk, H., and Bengio, Y. Learning phrase representations using rnn encoder-decoder for statistical machine translation. In *EMNLP*, 2014.
- Chomsky, N. Three models for the description of language. In *IRE Transactions on Information Theory (2)*, pp. 113124, 1956.
- Chomsky, N. On certain formal properties of grammars. In *Information and Control. 2 (2)*, pp. 137167, 1959.
- Das, S., Giles, C. L., and Sun, G.-Z. Learning context-free grammars: Capabilities and limitations of a recurrent neural network with an external stack memory. In *Proceedings of The Fourteenth Annual Conference of Cognitive Science Society*, pp. 14, 1992.
- Dave, A., Russakovsky, O., and Ramanan, D. Predictive-corrective networks for action detection. *arXiv preprint arXiv:1704.03615*, 2017.
- Dyer, C., Kuncoro, A., Ballesteros, M., and Smith, N. A. Recurrent neural network grammars. In *NAACL-HLT*, 2016.
- Gers, F. A. and Schmidhuber, J. Lstm recurrent networks learn simple context-free and context-sensitive languages. *IEEE Transactions on Neural Networks*, 12(6):1333–1340, 2001.
- Giles, C. L., Horne, B. G., and Lin, T. Learning a class of large finite state machines with a recurrent neural network. *Neural Networks*, 8(9):1359–1365, 1995.
- Graves, A., Wayne, G., and Danihelka, I. Neural turing machines. *arXiv preprint arXiv:1410.5401*, 2014.
- Hochreiter, S. and Schmidhuber, J. Long short-term memory. *Neural Computation*, 9:1735–80, 12 1997. doi: 10.1162/neco.1997.9.8.1735.
- Ivanov, Y. A. and Bobick, A. F. Recognition of visual activities and interactions by stochastic parsing. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 22(8):852–872, 2000.
- Jang, E., Gu, S., and Poole, B. Categorical reparameterization with gumbel-softmax. In *International Conference on Learning Representations (ICLR)*, 2017.
- Joo, S.-W. and Chellappa, R. Attribute grammar-based event recognition and anomaly detection. In *Computer Vision and Pattern Recognition Workshop*. IEEE, 2006.
- Kolen, J. F. Fool’s gold: Extracting finite state machines from recurrent network dynamics. In *Advances in Neural Information Processing Systems (NIPS)*, 1994.
- Kwak, S., Han, B., and Han, J. H. On-line video event detection by constraint flow. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 36(6):1174–1186, 2014.
- Liang, X., Zou, T., Guo, B., Li, S., Zhang, H., Zhang, S., Huang, H., and Chen, S. X. Assessing beijing’s pm2.5 pollution: severity, weather impact, apec and winter heating. In *Proceedings of the Royal Society A*, 471, 20150257, 2016.
- Maddison, C. J., Mnih, A., and Teh, Y. W. The concrete distribution: A continuous relaxation of discrete random variables. In *International Conference on Learning Representations (ICLR)*, 2017.
- Mayberry, M. R. and Miikkulainen, R. Sardsrn: A neural network shift-reduce parser. In *Proceedings of the 16th Annual Joint Conference on Artificial Intelligence*, 1999.
- Moore, D. and Essa, I. Recognizing multitasked activities from video using stochastic context-free grammar. In *Proceedings of the American Association for Artificial Intelligence (AAAI)*, 2002.
- Piergiovanni, A. and Ryoo, M. S. Fine-grained activity recognition in baseball videos. In *CVPR Workshop on Computer Vision in Sports*, 2018a.
- Piergiovanni, A. and Ryoo, M. S. Learning latent super-events to detect multiple activities in videos. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018b.
- Pirsiavash, H. and Ramanan, D. Parsing videos of actions with segmental grammars. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2014.
- Ryoo, M. S. and Aggarwal, J. K. Semantic representation and recognition of continued and recursive human activities. *International Journal of Computer Vision (IJCV)*, 82(1):1–24, 2009.

- Si, Z., Pei, M., Yao, B., and Zhu, S.-C. Unsupervised learning of event and-or grammar and semantics from video. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2011.
- Sigurdsson, G. A., Divvala, S., Farhadi, A., and Gupta, A. Asynchronous temporal fields for action recognition. *arXiv preprint arXiv:1612.06371*, 2016a.
- Sigurdsson, G. A., Varol, G., Wang, X., Farhadi, A., Laptev, I., and Gupta, A. Hollywood in homes: Crowdsourcing data collection for activity understanding. In *Proceedings of European Conference on Computer Vision (ECCV)*, 2016b.
- Socher, R., Lin, C. C., Manning, C., and Ng, A. Y. Parsing natural scenes and natural language with recursive neural networks. In *Advances in Neural Information Processing Systems (NIPS)*, 2011.
- Sun, G.-Z., Giles, C. L., Chen, H.-H., and Lee, Y.-C. The neural network pushdown automaton: Model, stack and learning simulations. *arXiv preprint arXiv:1711.05738*, 2017.
- Tiño, P., Horne, B. G., Giles, C. L., and Collingwood, P. C. Finite state machines and recurrent neural network-automata and dynamical systems approaches. In *Neural networks and pattern recognition*, pp. 171–219. Elsevier, 1998.
- Xu, H., Das, A., and Saenko, K. R-c3d: Region convolutional 3d network for temporal activity detection. *arXiv preprint arXiv:1703.07814*, 2017.
- Yeung, S., Russakovsky, O., Jin, N., Andriluka, M., Mori, G., and Fei-Fei, L. Every moment counts: Dense detailed labeling of actions in complex videos. *International Journal of Computer Vision (IJCV)*, pp. 1–15, 2015.
- Zeng, Z., Goodman, R. M., and Smyth, P. Discrete recurrent neural networks for grammatical inference. *IEEE Transactions on Neural Networks*, 5(2):320–330, 1994.

Learning Differentiable Grammars

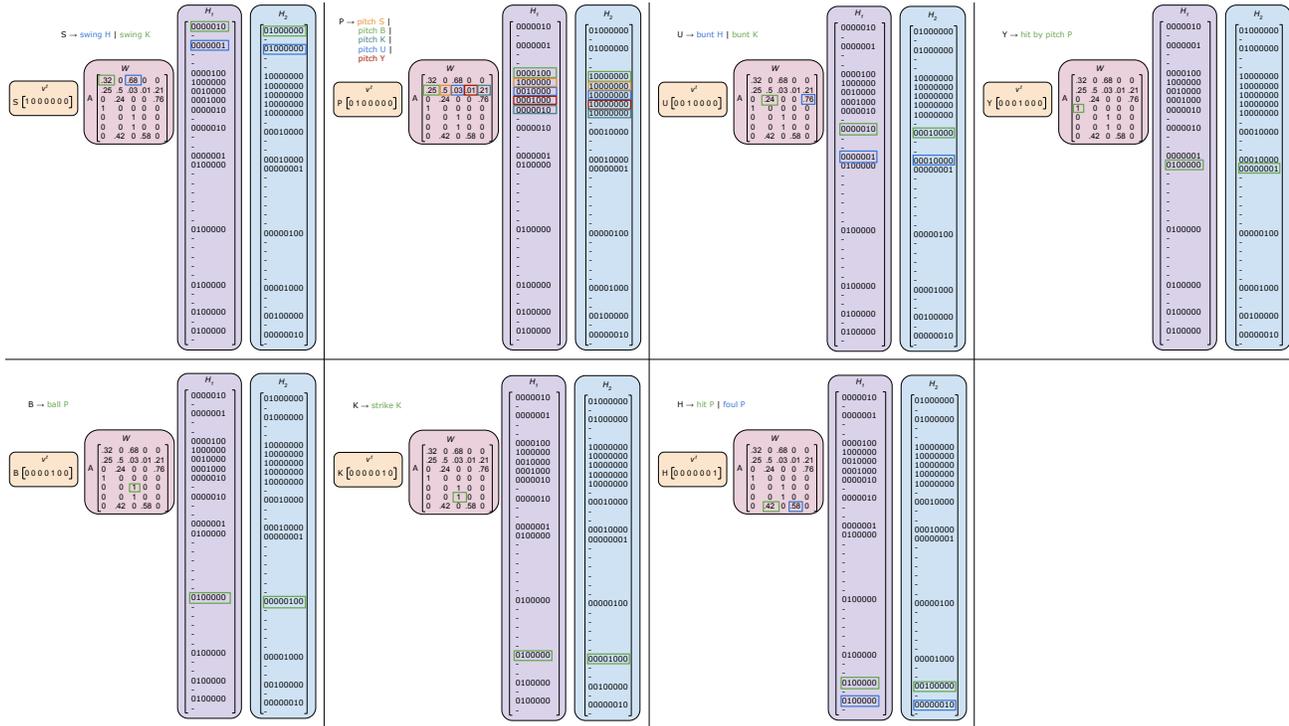


Figure 10. Visualization of the learned grammar for MLB-YouTube videos. Note, for simplicity, we only visualize the applicable rules in W . ‘-’ is used for terminals and non-terminals that are never used.