
Deep Predictive Coding for Multimodal Spatiotemporal Representation Learning

Anonymous Author(s)

Abstract

1 Common sense reasoning relates to the capacity of *learning representations* that
2 disentangle hidden factors behind spatiotemporal sensory data. In this work, we
3 hypothesize that the predictive coding theory of perception and learning from
4 neuroscience literature may be a promising candidate for implementing such
5 common-sense inductive biases. We build upon the *PredNet* implementation by
6 Lotter, Kreiman, and Cox (2016) and extend its application to the challenging
7 task of inferring abstract, everyday human actions such as *cooking* and *diving*.
8 Our transfer learning experiments demonstrate good generalization of learned
9 representations on the UCF-101 action classification dataset for both visual and
10 auditory modalities.

11 1 Motivation and Methods

12 The *PredNet* model by Lotter et al. (2016) was shown to learn representations that disentangle latent
13 variables correlated to the movement of objects in synthetic and natural images. We extend their study
14 to address the following questions: (1) Can unsupervised predictive coding models learn higher-level
15 spatiotemporal concepts, namely quotidian activities such as *driving* or *exercising*? (2) Are predictive
16 coding inductive biases general enough so that these models can also learn from auditory information?

17 **Predictive coding networks** Inspired by the predictive coding theory (Friston & Kiebel, 2009),
18 the *PredNet* model relies on the idea that to predict the next video frame, a model needs to capture
19 latent structure that explains the image sequences. The model architecture consists of recurrent
20 convolutional layers (Xingjian et al., 2015) that propagate bottom-up prediction errors, which are
21 used by the upper-level layers to generate new predictions. For implementation details, please refer
22 to the *PredNet* architecture description by Lotter et al. (2016).

23 **Unsupervised training** We evaluate predictive coding models trained on different quantities of
24 unlabeled videos. The main idea is that the more data we use to train the model, the more "common
25 sense" it should get about how events unfold in the world and, as a consequence, it should be better at
26 disentangling latent explanatory factors. Using as starting point a *PredNet* pre-trained on the KITTI
27 dataset (Geiger, Lenz, Stiller, & Urtasun, 2013), we further train the model with unlabeled videos (67
28 hours of visual data and 37 hours of auditory data) from the Moments in Time dataset (Monfort et al.,
29 2018), a large-scale activity recognition dataset.

30 **Supervised action recognition** For each sequence of ten frames in the input (video frame or audio
31 spectrogram), the *PredNet* activations for each layer are spatially pooled to match the higher-level
32 layer dimensions and concatenated to form one tensor representation with dimensions (16, 20, 339)
33 corresponding to a one-second spatiotemporal pattern. Those representations are then flattened and
34 used as input to an action classifier consisting of a Long Short-Term Memory (LSTM) (Hochreiter
35 & Schmidhuber, 1997) layer (64 hidden units) and a fully connected layer followed by a softmax
36 activation that outputs a probability distribution over the UCF-101 action classes.

37 2 Results and Discussion

38 **Action recognition using visual data** The predictive coding model with random weights gives a
39 poor top-1 accuracy of 1.64%, which is slightly above the random baseline (Table 1). However, when
40 we train the classifier with features generated by the 67-hour predictive coding model, the accuracy
41 increases to 51.9%, which is competitive with results from the unsupervised "tuple verification" by
42 Misra, Zitnick, and Hebert (2016) and an LSTM classifier using the Inception convolutional network
43 (Carreira & Zisserman, 2017). It is worth to note, however, that in both of those approaches, the
44 convolutional models are fine-tuned end-to-end using the UCF-101 labels. In our case, the predictive
45 coding weights were kept fixed, and only the weights from the LSTM classifier were optimized for
46 the specific task.

47 **Predictive coding can also model auditory data** We trained the LSTM classification model using
48 predictive coding representations extracted using audio spectrograms from the 51 action classes of
49 the UCF-101 dataset that contain auditory information. The top-1 accuracy results are reported in
50 Table 2. As expected, the audio information is much less useful to distinguish action classes, as many
51 videos have soundtracks and other kinds of audio data that are completely unrelated to the activity.
52 Still, there was a significant improvement from the classifier trained on the features generated by the
53 random-weights model to the classifier based on the 37-hour pre-trained model. For comparison,
54 we also report the results of the Caffenet version by Wang, Yang, and Meinel (2016), which is a
55 convolutional network trained on audio spectrograms. Remarkably, our simple one-layer LSTM
56 classifier is competitive with their complex convolutional model trained end-to-end using action class
57 labels, which demonstrates the generality of the predictive coding inductive bias.

Table 1: Visual action recognition. Accuracies (top-1 percentage) for different pre-trained models on the test set of UCF-101 split 1. We also include results for the CNN tuple verification (Misra et al., 2016) and an LSTM classifier trained on top an Inception convolutional network trained from scratch (Carreira & Zisserman, 2017).

Features + Classifier	Accuracy (%)	Pre-training dataset
PredNet Video random + LSTM	1.64	-
PredNet Video 67h + LSTM	51.9	Moments in Time
CNN tuple verification	50.2	UCF-101
Inception + LSTM	54.2	-

Table 2: Auditory action recognition. Accuracies (top-1 percentage) for different models on the test set of UCF-101 split 1 (only videos from the 51 classes that contain audio).

Features + Classifier	Accuracy (%)	Pre-training dataset
PredNet Audio random + LSTM	22.7	-
PredNet Audio 37h + LSTM	24.8	Moments in Time
Caffenet (Wang et al., 2016)	25.2	-

58 3 Final Remarks

59 This work explores unsupervised learning from spatiotemporal data and uses video understanding
60 tasks as a proxy to evaluate the quality of learned representations. We focus on models that can
61 learn from large amounts of unlabeled videos and use this experience to solve downstream tasks
62 involving smaller labeled datasets. Therefore, we *do not* pursue the solution of the action recognition
63 problem itself, for which all the state-of-art approaches depend on a copious amount of labeled data
64 for pre-training (Carreira & Zisserman, 2017). Our results show that predictive coding representations
65 learned using the Moments in Time dataset (Monfort et al., 2018) (without in-domain fine-tuning) are
66 competitive with other unsupervised baselines when evaluated on UCF-101 action recognition task.

67 **References**

- 68 Carreira, J., & Zisserman, A. (2017). Quo vadis, action recognition? a new model and the kinetics
69 dataset. In *Computer vision and pattern recognition (cvpr), 2017 ieee conference on* (pp.
70 4724–4733).
- 71 Friston, K., & Kiebel, S. (2009). Predictive coding under the free-energy principle. *Philosophical
72 Transactions of the Royal Society B: Biological Sciences*, 364(1521), 1211–1221.
- 73 Geiger, A., Lenz, P., Stiller, C., & Urtasun, R. (2013). Vision meets robotics: The kitti dataset. *The
74 International Journal of Robotics Research*, 32(11), 1231–1237.
- 75 Hochreiter, S., & Schmidhuber, J. (1997). Long short-term memory. *Neural computation*, 9(8),
76 1735–1780.
- 77 Lotter, W., Kreiman, G., & Cox, D. (2016). Deep predictive coding networks for video prediction
78 and unsupervised learning. *arXiv preprint arXiv:1605.08104*.
- 79 Misra, I., Zitnick, C. L., & Hebert, M. (2016). Shuffle and learn: unsupervised learning using
80 temporal order verification. In *European conference on computer vision* (pp. 527–544).
- 81 Monfort, M., Zhou, B., Bargal, S. A., Andonian, A., Yan, T., Ramakrishnan, K., ... others
82 (2018). Moments in time dataset: one million videos for event understanding. *arXiv preprint
83 arXiv:1801.03150*.
- 84 Wang, C., Yang, H., & Meinel, C. (2016). Exploring multimodal video representation for action
85 recognition. In *Neural networks (ijcnn), 2016 international joint conference on* (pp. 1924–
86 1931).
- 87 Xingjian, S., Chen, Z., Wang, H., Yeung, D.-Y., Wong, W.-K., & Woo, W.-c. (2015). Convolutional
88 lstm network: A machine learning approach for precipitation nowcasting. In *Advances in
89 neural information processing systems* (pp. 802–810).