

Video Prediction with Appearance and Motion Conditions

The paper "Video Prediction with Appearance and Motion Conditions" proposes a novel approach to predict future frames of a video sequence using appearance and motion conditions. The authors leverage the intuition that future frames are highly dependent on the past frames, and propose to use both appearance and motion conditions as guidance to generate more accurate predictions.

The proposed method consists of two stages: a motion encoder and a decoder. The motion encoder is responsible for encoding motion information from previous frames, while the decoder takes the encoded information and generates future frames. The authors also introduce an appearance encoder that encodes the appearance information of the previous frames, which is used as a conditional input to the decoder.

To evaluate their approach, the authors conduct experiments on two benchmark datasets and show that their method outperforms several state-of-the-art methods in terms of prediction accuracy. They also conduct ablation studies to demonstrate the effectiveness of each component of their proposed approach.

In summary, the key contributions of this paper are:

1. A novel approach to predict future frames of a video sequence using both appearance and motion conditions.
2. The introduction of an appearance encoder that encodes appearance information to be used as a conditional input to the decoder.
3. Experimental results that demonstrate the effectiveness of the proposed approach and its outperformance of several state-of-the-art methods.

To summarize, our contributions include:

- We propose AMC-GAN that can generate multiple different videos from a single image by manipulating input conditions. The code is available at <http://vision.snu.ac.kr/projects/amc-gan>.
- We develop a novel conditioning scheme that helps the training by varying appearance and motion conditions.
- We use perceptual triplet ranking to encourage videos of similar conditions to look similar. To our best knowledge, this has not been explored in video prediction.

The model proposed in "Video Prediction with Appearance and Motion Conditions" consists of two main components: a motion encoder and a decoder.

The motion encoder is responsible for encoding the motion information of previous frames. Specifically, it takes in a sequence of motion vectors, which represent the displacement of each pixel from the previous frame to the current frame. The motion encoder then applies convolutional layers to this sequence of motion vectors to generate a hidden representation, which is used as input to the decoder.

The decoder is responsible for generating future frames based on the encoded motion information and appearance information. The appearance information is encoded using an appearance encoder, which takes in a sequence of previous frames and generates a hidden representation that is also used as input to the decoder. The decoder then takes the concatenated encoded motion and appearance information and generates future frames using convolutional and deconvolutional layers.

During training, the model is trained to minimize the mean squared error between the predicted future frames and the ground truth future frames. The authors also introduce an adversarial loss term to encourage the generated frames to look realistic and coherent.

Overall, this model is designed to leverage both motion and appearance information to generate more accurate predictions of future video frames.

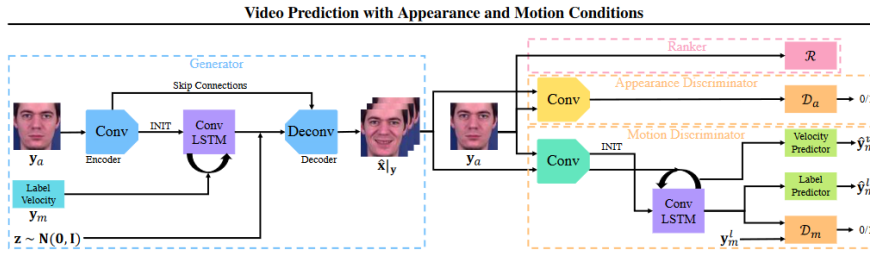


Figure 2. An overview of our AMC-GAN; we provide architecture and implementation details in the supplementary material.

The generator, also called the decoder, is responsible for generating future video frames based on the encoded appearance and motion information. It consists of several convolutional and deconvolutional layers, and takes as input the concatenated encoded appearance and motion information.

Specifically, the generator takes the encoded appearance information from the appearance encoder, and the encoded motion information from the motion encoder, and concatenates them along the channel dimension. This concatenated encoding is then passed through several convolutional layers, which are designed to extract features from the encoded input. The output of these convolutional layers is then passed through several deconvolutional layers, which are designed to upsample the features to generate a full-resolution video frame.

During training, the generator is trained to minimize the mean squared error between the generated frames and the ground truth frames, as well as the adversarial loss, which encourages the generated frames to look more realistic and coherent. Additionally, the generator is trained to produce frames that have high visual quality, as measured by the ranker component of the model.

Overall, the generator component of the model plays a crucial role in generating accurate and visually coherent predictions of future video frames, by leveraging both appearance and motion information and incorporating techniques such as adversarial training and quality-based loss functions.

Ranker: The ranker is a component introduced in the paper to further improve the visual quality of the generated frames. It is a convolutional neural network that is trained to rank the generated frames based on their visual quality. During training, the ranker is used to compute a quality score for each generated frame, which is then used to compute a quality-weighted loss. This loss encourages the model to generate frames that not only match the ground truth frames but also have high visual quality.

Appearance discriminator: The appearance discriminator is another component introduced in the paper to encourage the appearance encoder to produce more meaningful encodings. It is a convolutional neural network that is trained to distinguish between the appearance encodings produced by the appearance encoder and random noise. During training, the appearance encoder is trained to produce encodings that fool the appearance discriminator into thinking they are real appearance encodings, while the appearance discriminator is trained to correctly distinguish between real appearance encodings and fake ones.

Adversarial loss: The adversarial loss is a loss term introduced in the paper to encourage the generated frames to look more realistic and coherent. It is computed based on the output of the appearance discriminator, which is used to distinguish between real and generated frames. The adversarial loss encourages the generator (i.e., the decoder) to produce frames that fool the discriminator into thinking they are real frames, while the discriminator is trained to correctly distinguish between real and generated frames.

Overall, these components (the ranker, appearance discriminator, and adversarial loss) help to improve the visual quality and coherence of the generated frames, leading to even more accurate and realistic predictions.

The motion discriminator is a component of the model that is designed to encourage the motion encoder to produce more meaningful encodings of motion information. It is a convolutional neural network that is trained to distinguish between the motion encodings produced by the motion encoder and random noise.

During training, the motion encoder is trained to produce encodings that fool the motion discriminator into thinking they are real motion encodings, while the motion discriminator is trained to correctly distinguish between real motion encodings and fake ones. This encourages the motion encoder to produce more accurate and meaningful encodings of the motion information in the video sequence, which can lead to more accurate predictions of future frames.

While the "Video Prediction with Appearance and Motion Conditions" paper primarily focuses on predicting future video frames, the model proposed in the paper can also be used for semantic video generation to some extent.

By incorporating both appearance and motion information, the model can generate more realistic and coherent predictions of future video frames, which can be useful for generating videos with specific semantic content. For example, if the model is trained on a dataset of videos of people playing basketball, it can generate new videos of people playing basketball with realistic movements and appearances.

However, the model does not explicitly incorporate semantic information in the form of object or scene labels. Thus, the ability of the model to generate videos with specific semantic content is limited by the

ability of the model to learn and capture this semantic content from the training data. In other words, the model can generate videos that are visually consistent with the input data but may not be able to generate videos that conform to specific semantic constraints or requirements.

Overall, while the model proposed in the paper is not specifically designed for semantic video generation, it can be useful for generating videos with specific semantic content to some extent, especially when trained on datasets that contain such semantic content.