# Looking Ahead: Anticipating Pedestrians Crossing with Future Frames Prediction

**Mohamed Chaabane**[*]
Department of Computer Science
Colorado State University
chaabane@colostate.edu

**Ameni Trabelsi** [*]
Department of Computer Science
Colorado State University
ameni.trabelsi@colostate.edu

**Nathaniel Blanchard**
Department of Computer Science
Colorado State University
nathaniel.blanchard@colostate.edu

**Ross Beveridge**
Department of Computer Science
Colorado State University
ross.beveridge@colostate.edu

## Abstract

In this paper, we present an end-to-end future-prediction model that focuses on pedestrian safety. Specifically, our model uses previous video frames, recorded from the perspective of the vehicle, to predict if a pedestrian will cross in front of the vehicle. The long term goal of this work is to design a fully autonomous system that acts and reacts as a defensive human driver would — predicting future events and reacting to mitigate risk. We focus on pedestrian-vehicle interactions because of the high risk of harm to the pedestrian if their actions are miss-predicted. Our end-to-end model consists of two stages: the first stage is an encoder/decoder network that learns to predict future video frames. The second stage is a deep spatio-temporal network that utilizes the predicted frames of the first stage to predict the pedestrian's future action. Our system achieves state-of-the-art accuracy on pedestrian behavior prediction and future frames prediction on the Joint Attention for Autonomous Driving (JAAD) dataset. [2]

## 1 Introduction

For the last decade, researchers and companies alike have been striving to achieve Level 5 autonomy for self-driving vehicles (i.e., autonomous operation with no human intervention) [1]. One benefit of widespread Level 5 autonomy would be a reduction in vehicular accidents caused by human error, which kills around $1.35$ million people a year [2]. The feasibility of Level 5 has been buoyed by research overcoming several milestones, such as autonomous driving on highways [3], rough terrains [4], and urban environments [5]. These breakthroughs have led several companies to invest in consumer-grade autonomous vehicles. Nonetheless, there are still many hurdles left to overcome before Level 5 autonomy is reached. One issue is that these vehicles are not replicating the behavior of human drivers — specifically, the ways human drivers communicate with each other, and surrounding agents [6, 7]. For example, pedestrians use nonverbal cues such as eye gaze or hand gestures to communicate their crossing intent to human drivers. Conversely, in autonomous vehicles, any misunderstanding of pedestrians gestures or misprediction of their intents will more likely cause traffic accidents. In Figure 1, we present an uncertain situation where a pedestrian may or may not cross in front of the vehicle. Incorrectly predicting that the pedestrian will not cross would likely lead to pedestrian injury.

This work focuses on the understanding and anticipation of pedestrian action as a step toward an autonomous vehicle capable of understanding such nonverbal communication. To this end, we employ the Joint Attention for Autonomous Driving (JAAD) dataset [8], which is specifically designed to capture the visual and behavioral complexity between pedestrians and drivers from the perspective of the driver. JAAD consists of a multitude of complexities and conditions

---

[*]These two authors contributed equally

[2]Accepted as an oral presentation in WACV-20

Figure 1: Will the pedestrian step in the vehicle's path? Human drivers and pedestrians interact through nonverbal, physical communication. For example, a pedestrian might make eye contact with a driver. From this communication, the driver and pedestrians predict each other's actions. The driver may anticipate the pedestrian will attempt to cross the road, and slow down to facilitate. Fully autonomous vehicles are expected to identify the signals pedestrians exude, predict potential actions, and react appropriately. This work presents a pedestrian prediction model that embodies these principles.

(e.g., weather, location, etc.,) that allow us to thoroughly test the robustness of our model. The model itself operates in two stages: first, the network takes $N$ video frames (the past) and predicts the next $N$ video frames (the future). Specifically, the input frames are encoded into features using spatio-temporal 3D-Convolution layers. The features are then decoded into predictions of the next future frames using depth-wise separable convolutional LSTM layers. The second stage of the model utilizes the predicted $N$ future frames to classify if the pedestrian will step in front of the vehicle. This stage processes the predicted frames with a supervised action network to predict pedestrians crossing actions. The full model is trained in an end-to-end fashion to minimize loss on both future frames prediction and pedestrian crossing prediction. Our model is able to achieve state-of-the-art accuracy on both future frames prediction and pedestrian future action prediction. An ablation study of the model's components shows our model is capable of capturing the key aspects of both egocentric and external movement required for accurate pedestrian prediction. The robustness of the future frames prediction component allows the action classification component to achieve state-of-the-art accurately, further indicating the predicted frames have the level of detail required to make such predictions.

In the future, these models may prove useful to the understanding and prediction of other environmental variables, such as predicting the behavior of other vehicles. The model may also be useful for predicting pedestrian intent, rather than action, something that can be inferred from additional labels provided by the JAAD dataset. For now, we focus on anticipating one of the most dangerous circumstances: a pedestrian stepping in front of a vehicle.

In summary, this paper makes the following contributions:

- We present a future video frames prediction encoder/decoder network that operates in unsupervised manner to predict $N$ future frames of a video using $N$ initial frames.
- We propose an end-to-end model that predicts the future video frames and uses the predicted frames as input for supervised action recognition network to predict when pedestrians will step in front of the vehicle.
- We achieve state-of-the-art performance on both future frames prediction and predicting pedestrian future crossing action on JAAD dataset.
- We conduct a thorough ablation study that shows the model components are robust, efficient, and effective across a multitude of weather conditions, locations, and other variables.

## 2    Related Work

Our work is related to two avenues of previous work: future video frames prediction and pedestrian future action prediction.

**Future Video Frames Prediction:** In recent years, Recurrent Neural Networks (RNN) have been widely used in future video frames prediction [9, 10]. Much of this work focuses on predicting one frame into the future, with an option to modify the network for long range prediction by taking in the predicted frame as input. Some examples of this include [11, 12, 9, 13, 10]. Our work differentiates itself by focusing on an architecture explicitly designed to predict many frames into the future.

Other studies have applied stochastic variational methods to future frames prediction, such as [14] and [15]. Several others have worked with generative adversarial networks [16, 17, 18]. These studies focus on the sharpness of the generated video frames, treating it as a major characteristic to distinguish between real and fake video frames. However, models from these techniques are often unstable and difficult to properly evaluate. We focus the training and evaluation of our model on traditional dataset-driven performance, which allows us to accurately understand the applications and limitations of our results.

Recently, Gujjar et al. [19] studied the prediction of urban pedestrian actions by predicting future frames of traffic scene videos. Of these related works, this study's approach is the most similar to ours. The key difference is that our network provides better spatio-temporal representations for next frames prediction. This distinction is due to our network's deeper nature and the use of residual connections, which enable it to extract more complex features without falling into vanishing gradient problem. Our network also has a reduced running time and improved performance, thanks to the use of depth-wise separable convLSTMs rather than standard convLSTMs (see Section 4.1.3). Furthermore, our network uses lateral connections to reduce the blur in future frames, something especially important for long term prediction.

**Pedestrian Future Action Prediction:** In this work, we are mainly concerned with modeling pedestrian future action in the context of autonomous driving cars. Many existing approaches are based on Hidden Markov Model (HMM) where the pedestrian's intent is represented in the hidden state [20, 21]. These approaches have been extended to combining the motion models of all the intentions together into a single Mixed Observability Markov Decision Process (MOMDP), which is a structured variant of the more common Partially Observable Markov Decision Process (POMDP) [22]. Although models utilizing the Markovian process are known for their fast adaptability, their assumption can be restrictive due to insufficient prior conditioning. Thus, the main limitation of the presented methods is their lack of memory. Our approach overcomes this limitation by using RNN, known for its good memory retention qualities, in our model, thus extending its long-term memory.

Other approaches for prediction of time series assume they are samples from a process generated by a linear process driven by a white, zero-mean, Gaussian input [23, 24]. Although they can be more accurate, Gaussian processes have shown to be slower than Markov models since they use the entire observed trajectory to predict the future state [25]. Switching linear dynamical models, applied in constrained environments, were introduced as extensions to these models [23, 26]. Kooij et al. proposes a Dynamic Bayesian Network for pedestrian path prediction which incorporates environment parameters such as pedestrian situational awareness and head orientation with a Switching Linear Dynamical System to predict changes in the dynamics of pedestrians [26]. These motion models require accurate and precise segmentation and tracking of pedestrians to be efficient. Such assumption can be challenging due to the difficulty of extracting reliable image features for segmentation and tracking [27].

Consequently, many approaches, including ours, study pedestrian activity models that are extracted directly from the image space of the captured scenes. Hasan et al. [28] treat the prediction of adverse pedestrian actions as an anomaly detection problem. They built a fully convolutional autoencoder to learn the local features followed by a classifier to capture the regularities. Rasouli et al. [8] extract context features extracted from input frames using AlexNet [29] and train a linear SVM model to predict future crossing action of pedestrians on JAAD dataset. These approaches are limited because they focus only on spatial appearances, ignoring the temporal coherence in long-term motions. To solve this issue, Gujjar et al. [19] processes the crossing actions classification by feeding the predicted frames of their future frame prediction network to a C3D based network [30] which takes into account the temporal dynamics in addition to the spatial appearances. As mentioned before, this work is similar to ours, but varies in the training strategy employed, the experimental study of network components, and our network's higher performance with shorter running time (see Section 4.1).

# 3   Methods

Here, we detail the end-to-end model (Section 3.1) and enumerate experimental details (Section 3.2), including an overview of the dataset (Section 3.2.1) and model search procedure (Section 3.2.2).

## 3.1 Architecture

Our end-to-end model consisted of two stages: the first stage was an unsupervised encoder/decoder network that generated predicted future video frames. The second stage was a deep spatio-temporal action recognition network that utilized the generated video frames to predict pedestrian action — specifically, if the pedestrian would cross in front of the vehicle.
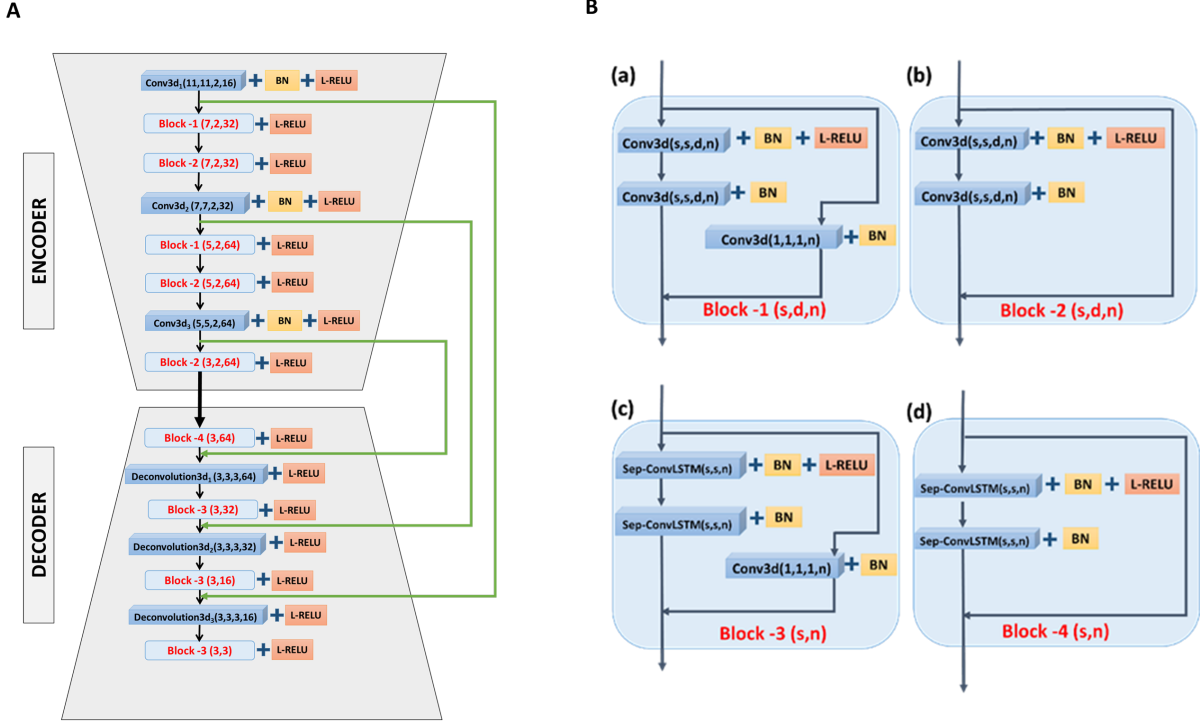
### 3.1.1 Future Frames Prediction Component



Figure 2: A- Overview of the proposed encoder/decoder network used in our approach to predict the next $N$ video frames (the future) using the first $N$ video frames (the past) as input. B- The 4 different residual blocks used in the architecture in A. (a) and (b) are the residual blocks used in the encoder. (c) and (d) are used in the decoder.

The future frames prediction component of the model was an unsupervised encoder/decoder that predicted future frames in a video sequence. $N$ consecutive video frames were input into the model, and the model predicted $N$ frames into the future. Figure 2A is a visual representation of the encoder/decoder architecture. The encoder mapped the input sequence of frames into a low-dimensional feature space with different resolutions. The decoder mapped the low-dimensional representation space of the input frames to an output image space. We define the input as a sequence $x = \{x_1, \ldots, x_N\}$ where $x_i$ refers to the frame of index $i$ in the input. The encoder yielded dense representations $z = \{z_1, ..., z_N\}$ from the input. The decoder then output the following frames, denoted as $y' = \{y'_{N+1}, \ldots, y'_{2N}\}$, as a prediction of the ground-truth frames $y = \{y_{N+1}, \ldots, y_{2N}\}$.

Next, we will detail the precise structure of the encoder and decoder.

**Encoder:** The encoder was a spatio-temporal neural network composed of three-dimensional convolutional layers. 3D convolutions modeled both the spatial and sequential relationships of temporal connections across frames. $N$ RGB frames were the encoder input. The size of input is $3 \times N \times H \times W$. The temporal duration of the outputted feature maps matched the input images.

The main components of the encoder were Block-1, Block-2 and 3D convolutions (conv3d) layers (as shown in Figure 2). Each convolutional operation was followed by batch normalization [31] and a Leaky-ReLU activation function. The residual blocks (Block-1 and Block-2 in Figure 2B) consisted of two 3D convolutions with no stride. Residual connections addition operations were identity shortcuts for Block-2 and $1 \times 1 \times 1$ 3D convolution operations matched

the input and output dimensions in Block-1. Inputs were downsampled in $conv3d_1$, $conv3d_2$, and $conv3d_3$ with a stride of 2. The filters of the last Block-2 of our encoder were time-dilated filters which captured the temporal dependency with various resolutions for the input sequences.

**Decoder:** The decoder was composed of convLSTMs layers interspersed with up-sampling layers. The main components of the decoder architecture are Block-3 and Block-4 shown in Figure 2B. The blocks were composed of two consecutive depth-wise separable convLSTMs with a residual connection connecting the input to the output. The residual connection in Block-4 was a simple identity shortcut, while Block-3 was $1 \times 1 \times 1$ 3D convolution operations matching the input and output dimensions. For the up-sampling layers , we used deconvolution [32] as it uses learnable parameters rather than predefined up-sampling.

**Encoder/Decoder Connections:** Lateral skip connections crossed from same-sized parts in the encoder to the decoder (the green lines shown in Figure 2). The lateral connections increased the amount of frames details, resulting in increased detail in the predicted frames.

### 3.1.2 Pedestrian Action Prediction Component

The second stage of the model consisted of a fine-tuned early action recognition network, the 'Temporal 3D ConvNets' (T3D) [33]. This stage predicted if the pedestrians would cross the street in the scene. $N$ predicted frames, generated from the encoder/decoder, were input into the network. The last classification layer of the T3D network was replaced with a fully connected layer, which produced a single output followed by sigmoidal activation. The component was trained with binary-cross entropy loss.

### 3.1.3 End-to-end Loss

The model had two main tasks that could be used for training: the future frames prediction loss and the pedestrian crossing loss. The full model was trained end-to-end with a multi-task learning objective:

$$L_{recog} = \lambda L_{pred} + L_{ce}(Y, \hat{Y}), \tag{1}$$

where $L_{ce}$ was the cross-entropy loss for crossing action classification, $\hat{Y}$ and $Y$ were high-level predictions and corresponding ground truth classes. The weight factor of the loss was $\lambda$. $L_{pred}$ was the future frame prediction loss, the pixel-wise loss between the pixels of the $N$ predicted frames and the $N$ ground-truth frames, defined as:

$$L_{pred} = \frac{1}{P} \left( \sum_{t=N+1}^{2N} \sum_{i=1}^{P} (y_{t,i} - y'_{t,i})^2 + \sum_{t=N+1}^{2N} \sum_{i=1}^{P} |y_{t,i} - y'_{t,i}| \right) \tag{2}$$

where $P = H \times W$, the number of pixels per frame. We use a combination of $l_1$ and $l_2$ norm losses for regularization purposes.

## 3.2 Experiments

In this subsection, we describe the dataset that was used to evaluate our model (3.2.1), the model selection search details (3.2.2), and the experimental details of our experiments (3.2.3).

### 3.2.1 Dataset

The Joint Attention for Autonomous Driving (JAAD) dataset was used for training and evaluation across all experiments [8]. JAAD was designed to study the behavior of traffic participants, with a focus on pedestrian and the driver behavior. Approximately 240 hours of driving videos were collected in several locations using two vehicles equipped with wide-angle video cameras. Cameras were mounted in the center of the cars' windshields, below the rear view mirror. The dataset consists of 346 high-resolution video clips that focus on pedestrian and driver behaviors at the point of crossing. Most of the data was collected in urban areas, but some clips were filmed in rural areas.

We picked JAAD because it contained a variety of interactions and complex situations that may impact the behavior of the traffic participants. Complex Interactions included pedestrians crossing in groups or individually, partially occluded pedestrians, and pedestrians walking parallel to the street. Complex situations included interactions with other drivers in parking lots or uncontrolled intersections. A variety of difficult conditions are represented in the dataset, such as weather, light conditions, and day/night conditions.

| Calibration parameters | search space |
|---|---|
| Spatial filter size of 3D Convs | [3,5,7,11] |
| Temporal dilation rate | [1,2,3,4] |
| Spatial filter size of sep-ConvLSTMs | [3,5,7] |
| Temporal filter size of 3D Convs | [2,3,4] |
| Temporal filter size of sep-ConvLSTMs | [2,3,4] |

Table 1: encoder/decoder network hyperparameters and search space. Note: temporal dilation rate is implemented only in the last Block of the encoder.

### 3.2.2 Model Search

We conducted a large-scale model search to identify the best performing encoder/decoder model for our future video frames prediction component. The search involved three steps, which are further detailed below. First, we selected a potential architecture design. Second, we trained 38 variations of the architecture with random sets of hyperparameters. Finally, the best architecture/hyperparameters combination was identified by lowest error on average future frames prediction, across $N$ frames.

The architecture and hyperparameters presented in Figure 2 represent the combination with the highest performance.

**Architecture Design:** The main encoder/decoder components were experimentally manipulated in order to test multiple architectural designs. The number of layers, the order of the layers, and the number of channels within layers were all varied. Across all variations, the encoder output remained unchanged because the spatial dimension of the input was consistently downsampled by 8. In the decoder, the convLSTM block(s)-deconvolution pattern was consistently used.

**Hyperparameters Selection:** For each selected architecture, 38 hyperparameter settings were randomly sampled [34]. Each parameter setting was evaluated using its average pixel-wise prediction $l_1$ error on the validation set. Details of the hyperparameter search spaces are summarized in Table 1.

### 3.2.3 Experimental Setup

We used the same training, validation, and test clips presented in [19], which allowed us to directly compare our performance. $60\%$ of the data was set aside for training, $10\%$ for validation and $30\%$ for testing. Clips were divided into $2N$-frame videos with a temporal stride of 1. The frames were resized to $128 \times 208$, with $N = 16$. Thus, the model input was $3 \times 16 \times 128 \times 208$.

The future frames prediction component was pre-trained on JAAD dataset and the T3D network was pre-trained on UCF101 dataset [35]. The model was fine-tuned on the end-to-end loss, as detailed in Section 3.1.3, with the Adagrad [36] optimizer. The weight factor $\lambda$ was set to 0.5. The end-to-end model was fine-tuned for 30 epochs with a learning rate of $10^{-4}$. Our end-to-end model was fine-tuned with the same JAAD splits described in [19, 8] for comparison. The test set is composed of 1257 16-frames videos of which 474 are labelled as crossing and 783 as not crossing.

All experiments were run on an Ubuntu server with a TITAN X GPU with 12 GB of memory.

## 4 Results

We evaluated our model on two predictive tasks: future frames and future pedestrian crossing. In this section, we present a complete description of our experiments and subsequent results. In summary, our model outperformed state-of-the-art models in predicting future pedestrian crossing action on JAAD dataset, yielding 86.7 Average Precision (AP) as compared to 81.14 AP achieved by the best performing approach [19]. The source code and trained models will be made available to the public.

### 4.1 Future Frames Prediction

Future Frames Prediction is the problem of accurately generating future frames given a set of consecutive previous frames. We quantitatively compared our model's performance on future frames prediction against other state-of-the-art methods (Section 4.1.1), detailed consistent architectural trends that we identified from our model search (Section 4.1.2), conducted an Ablation Study (Section 4.1.3), and performed a qualitative analysis on our model (Section 4.1.4).

| Model | $l_1$ loss $(\times 10^{-1})$ |
|---|---|
| Res-EnDec [19] | $1.37 \pm 0.37$ |
| PredRNN++ [10] | $1.61 \pm 0.35$ |
| PredNet [12] | $1.30 \pm 0.41$ |
| Ours | $\mathbf{1.12 \pm 0.32}$ |

Table 2: Comparison of our model with state-of-the-art methods on JAAD dataset. We report pixel-wise prediction $l_1$ error averaged over the 16 predicted frames.

### 4.1.1 Quantitative Analysis

Most state-of-the-art methods for future frame prediction have published results on datasets that have static cameras and identical backgrounds, such as the Moving MNIST dataset [9] and the KTH action dataset [37]. In this work, we evaluated and compared our model to state-of-the-art models on the JAAD dataset, which consists of complex, real-world interactions and variability that we can expect autonomous vehicles to encounter regularly (we detail the complexities of the dataset in Section 3.2.1). We limited our quantitative comparison to state-of-the-art methods which have publicly available code, including PredRNN++ [10] and PredNet [12], and Res-EnDec [19], who trained and tested on the same data split we used (see Section 3.2.3).

PredRNN++ and PredNet were originally designed to predict one frame ahead, but they were modified to predict multiple future frames by treating their predicted frame as input and recursively iterating. We trained both models from scratch on JAAD dataset on the same train/validation/test described in our methods.

Res-EnDec [19] is, to the best of our knowledge, the best performing model with published results on JAAD dataset. The code for this model was not released, so we directly compared our results to the results reported in [19].

In Table2, we present the average performance of each model on the future frame prediction task across 16 time steps. Our model has the lowest error among the models we tested.

In Figure 3, we plotted the $l_1$ loss of each model between the predicted and the ground truth frames across multiple time steps, up to 16 frames into the future. This allowed us to evaluate each model's relative consistency across the predicted frames. The quality of the predicted frames degraded over time for all models. Res-EnDec model had a slight variation in this trend; the error was higher for time step 1, explained in [19] as a result caused by the reverse ordering of their inputs.

Our model outperformed the state-of-the-art methods for all time steps except the initial time step, where PredNet produced slightly better performance. PredRNN++ and PredNet produced reasonably accurate short term predictions, however, they broke down when extrapolating further into the future. Compared with our model and [19], their errors increased considerably over time. This is expected, since both PredRNN++ and PredNet are not explicitly designed and optimized to predict multiple future frames. The predicted frames unavoidably had different statistics than the natural images the models were optimized for [38]. Given that, it is unsurprising that our model and [19] have better performance for long term predictions.
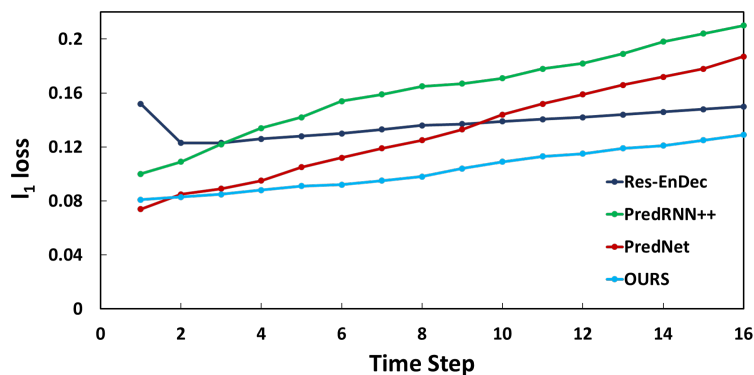


Figure 3: Comparison of our model with state-of-the-art models on JAAD dataset. We report frame-wise $l_1$ loss of generated sequences on the JAAD test set. All models are trained from scratch on JAAD training set.

| Model | $l_1$ loss ($\times 10^{-1}$) | SSIM | Time (ms) |
|---|---|---|---|
| Ours | $1.12 \pm 0.32$ | 0.924 | 88.74 |
| V-1 (Reg-convLSTM) | $1.13 \pm 0.39$ | 0.920 | 100.36 |
| V-2 (Spatial-convLSTM) | $1.11 \pm 0.31$ | 0.929 | 109.84 |
| V-3 (Depth-convLSTM) | $1.16 \pm 0.35$ | 0.911 | 80.69 |
| V-4 (w/o laterals) | $1.26 \pm 0.37$ | 0.883 | 86.69 |
| V-5 (w/o residuals) | $1.23 \pm 0.34$ | 0.894 | 85.31 |
| V-6 (Undilated) | $1.20 \pm 0.29$ | 0.905 | 87.43 |
| V-7 (w/o Deconv) | $1.18 \pm 0.43$ | 0.909 | 87.25 |

Table 3: Ablation study of different model variants on the JAAD dataset. We report the average pixel-wise prediction $l_1$ error, the per-frame structural similarity index measure (SSIM) [39] averaged over the 16 predicted frames and the running time.

### 4.1.2 Architectural Trends Identified from Model Search

The performance of our selected future frames prediction component, shown in Figure 3, resulted from the model selection criteria explained in Section 3.2.2. Through the model search process, we discovered several interesting and consistent architectural trends and their relationship to next frame prediction performance and speed. We report these trends here, so they can inform future model searches related to this topic.

We found that kernels with decreasing sizes in the spatial dimension ($11 \times 11 \to 7 \times 7 \to 5 \times 5 \to 3 \times 3$) and constant size in the time dimension exhibited higher performance across temporal variations, while still capturing the details from the scene. We also noted that ascendant dilation rates in the temporal dimensions at the end of the encoder network ($1 \to 2 \to 4$) performed best, and enabled our encoder network to have larger temporal receptive field without the need to add more layers.

### 4.1.3 Ablation Study

In Table 3 we present the results of an ablation study investigating how 7 architectural variants (V-1 - V-7) affected our model's performance and speed.

In the first three architecture variants, we experimented with convLSTM variants in the decoder. Our model used a depth-wise separable convLSTM. Variant one (V-1) consisted of standard convLSTMs, V-2 was spatially separable convLSTMs, and V-3 was depth-wise convLSTMs. SSIM performance across all three variants was similar, but V-2 (used spatial convLSTM layers) performed slightly better than the other variants. In the end, we opted for depth-wise separable convLSTMs because our network ran $21\,ms$ faster than V-2, with very similar performance.

The remaining variants removed or replaced network components, highlighting the importance of each component. In V-4, we removed the lateral connections, stymieing pixel information in the deconvolution. In V-5, residual connections were removed, underscoring the importance of this feature. In V-6, the temporal dilation was set to 1, limiting the model's temporal information. Finally, in V-7, we replaced deconvolution with interpolation. We suspect deconvolution can reconstruct the shape and boundaries more accurately than interpolation.
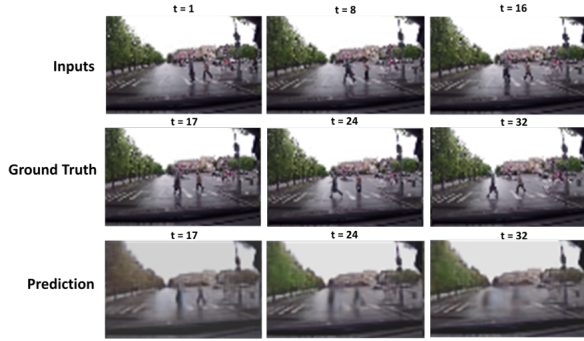
### 4.1.4 Qualitative Analysis

Sample predictions for our model on the JAAD dataset are shown in Figure 4. The model is able to make accurate predictions in a wide range of scenarios. In sequence 1 and sequence 3 of Figure 4, pedestrians are crossing the street in two different weather conditions (rain/snow), and the model successfully predicts their positions moving forward from frame to frame. Similarly in Sequence 2, a car is passing in the opposite direction, and the model, while not perfect, is able to predict its trajectory, as well as the position of the zebra stripes of the crossing area and other stationary objects as the driver moves forward. Sequence 4 illustrates that our model can distinguish between moving pedestrians and standing pedestrians as it accurately predicts the movement of the pedestrian who is crossing the street while the second pedestrian in the left-hand side of the scene was still standing in their same initial position.
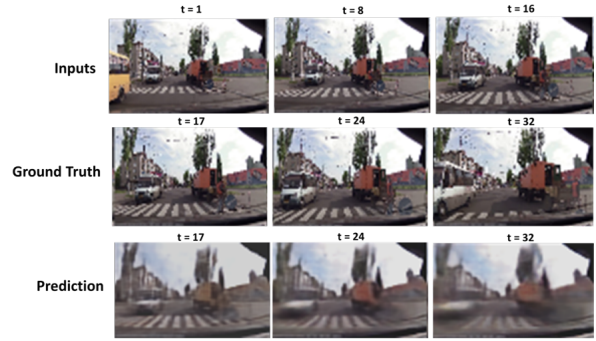
## 4.2 Pedestrian Future Crossing Action Prediction

As described in Section 3.1.2, we extended our future frames prediction component with the T3D network to classify a predicted future scene as pedestrian crossing/not crossing. We compared T3D with some variants where we replace T3D network with other 3D ConvNets networks for action recognition including C3D [30] and 3D ConvNets based
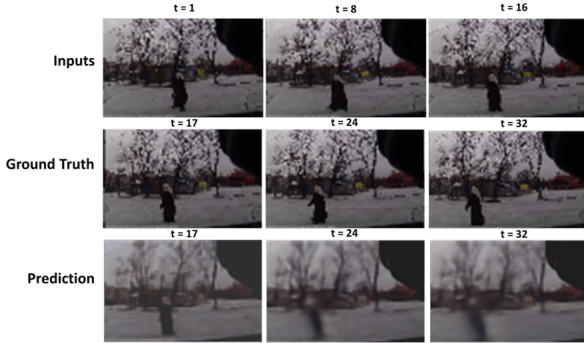
Figure 4: Future frames prediction examples on JAAD dataset. We predict 16 frames into the future by observing 16 frames.

| Model | AP | Time (ms) |
|---|---|---|
| C3D | 84.9 | **93.16** |
| Inception-3D | 82.2 | 104.68 |
| ResNet3D-50 | 79.6 | 99.45 |
| T3D | **86.7** | 102.48 |

Table 4: Evaluation results of future crossing action prediction with different action recognition networks as extension to our future frames prediction component. We report Average Precision (AP) and the running time. All action recognition networks are pre-trained on UCF101 dataset and trained using same multi-task training strategy.

on Inception [40] and ResNet-34 [41]. To make the comparison fair, all action recognition networks were pre-trained on the UCF101 dataset and then end-to-end trained to minimize the multi-task learning objective function shown in equation 1. Table 4 compares average precision (AP) and running time for predicting future crossing action using only the generated 16 frames. Running time in Table 4 corresponds to our model's 16-frame input and 16-frame prediction, as well as the prediction of future crossing action. T3D network outperformed other 3D convNets, likely because it modeled variable temporal 3D convolution kernel depths over shorter and longer time ranges.

| Model | AP |
|---|---|
| Action [8] | $39.24 \pm 16.23$ |
| Action + Context [8] | $62.73 \pm 13.16$ |
| Res-EnDec [19] | 81.14 |
| Ours ( Separate training ) | 85.8 |
| Ours ( Joint training ) | **86.7** |

Table 5: Comparison of average precision (AP) of our model (with two different training strategies) with state-of-the-art methods for predicting future crossing action on JAAD dataset

Using JAAD dataset, predicting 16 future frames at a frame rate of $30 fps$ corresponds to looking ahead $533ms$ in time. Any advantage gained from this is reduced by the running time of the model. The running time for our model is $102.48ms$ which provides a reliable maximum look-ahead time of $430ms$. This allows $81\%$ of the time before the occurrence of the final predicted frame to be utilized for defensive reaction.

We also compared AP scores of our model with the results presented by Rasouli et al. in [8] and Gujjar et al. in [19] in Table 5. Our model outperformed Action [8], Action + Context [8] and Res-EnDec [19] by about $47.5\%$, $24\%$ and $5.6\%$, respectively. Our model and Res-EnDec model outperformed models in [8] with more than $18\%$, which reflects the effectiveness of predicting future video frames and using those frames to predict future crossing action. Using same action recognition network (C3D) as in [19], our model achieves $84.9$ AP (Table 4) which is $3.8\%$ higher than Res-EnDec model. Likely, some of our improvements can be attributed to the increased quality of our generated future frames. The remaining improvements stem from the multi-task training strategy we employ — our model's AP rose $0.9\%$ when we implemented this approach, in agreement with previous findings [42].

## 5  Conclusion

In this paper, we introduced an end-to-end future-prediction model that uses previous video frames to predict the pedestrian's future crossing action. Our multi-task model is composed of two stages. The first stage consists in processing $N$ input video frames using an encoder/decoder network to predict $N$ future frames. The second stage utilizes the predicted frames to classify the pedestrian's future actions. Our end-to-end model predicts a future pedestrian crossing action with an effective look-ahead of $430 \, ms$ on JAAD dataset.

Our end-to-end model achieves state-of-the-art performance in model speed ($102.48 \, ms$), predicting future frames ($0.924 \, SSIM$), and predicting pedestrians crossing ($86.7 \, AP$) on JAAD dataset. Our ablation study demonstrates the effectiveness of the different model's components on the model performance. Further quantitative and qualitative experiments have shown the ability of our proposed model across the various weather conditions, locations and other variables included in the JAAD dataset.

While this is just one step in replicating human-like behavior for self-driving cars, the safety of pedestrians is an essential first step to tackle, and our work has promising implications for the future of replicating human-like behavior in the context of fully autonomous driving.

## References

[1] Claudine Badue, Rânik Guidolini, Raphael Vivacqua Carneiro, Pedro Azevedo, Vinicius Brito Cardoso, Avelino Forechi, Luan Ferreira Reis Jesus, Rodrigo Ferreira Berriel, Thiago Meireles Paixão, Filipe Mutz, et al. Self-driving cars: A survey. *arXiv preprint arXiv:1901.04407*, 2019.

[2] *Global status report on road safety 2018.* Geneva: World Health Organization, 2018.

[3] Ernst Dieter Dickmanns and Alfred Zapp. A curvature-based scheme for improving road vehicle guidance by computer vision. In *Mobile Robots I*, volume 727, pages 161–168. International Society for Optics and Photonics, 1987.

[4] Sebastian Thrun, Mike Montemerlo, Hendrik Dahlkamp, David Stavens, Andrei Aron, James Diebel, Philip Fong, John Gale, Morgan Halpenny, Gabriel Hoffmann, et al. Stanley: The robot that won the darpa grand challenge. *Journal of field Robotics*, 23(9):661–692, 2006.

[5] Chris Urmson, Joshua Anhalt, Drew Bagnell, Christopher Baker, Robert Bittner, MN Clark, John Dolan, Dave Duggins, Tugrul Galatali, Chris Geyer, et al. Autonomous driving in urban environments: Boss and the urban challenge. *Journal of Field Robotics*, 25(8):425–466, 2008.

[6] Karthik Mahadevan, Sowmya Somanath, and Ehud Sharlin. Communicating awareness and intent in autonomous vehicle-pedestrian interaction. In *Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems*, page 429. ACM, 2018.

[7] Matus Sucha, Daniel Dostal, and Ralf Risser. Pedestrian-driver communication and decision strategies at marked crossings. *Accident Analysis & Prevention*, 102:41–50, 2017.

[8] Amir Rasouli, Iuliia Kotseruba, and John K Tsotsos. Are they going to cross? a benchmark dataset and baseline for pedestrian crosswalk behavior. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 206–213, 2017.

[9] Nitish Srivastava, Elman Mansimov, and Ruslan Salakhudinov. Unsupervised learning of video representations using lstms. In *International conference on machine learning*, pages 843–852, 2015.

[10] Yunbo Wang, Zhifeng Gao, Mingsheng Long, Jianmin Wang, and S Yu Philip. Predrnn++: Towards a resolution of the deep-in-time dilemma in spatiotemporal predictive learning. In *International Conference on Machine Learning*, pages 5110–5119, 2018.

[11] Chelsea Finn, Ian Goodfellow, and Sergey Levine. Unsupervised learning for physical interaction through video prediction. In *Advances in neural information processing systems*, pages 64–72, 2016.

[12] William Lotter, Gabriel Kreiman, and David Cox. Deep predictive coding networks for video prediction and unsupervised learning. In *Proceedings of the International Conference on Learning Representations*, 2017.

[13] Ilya Sutskever, Oriol Vinyals, and Quoc V Le. Sequence to sequence learning with neural networks. In *Advances in neural information processing systems*, pages 3104–3112, 2014.

[14] Mohammad Babaeizadeh, Chelsea Finn, Dumitru Erhan, Roy H Campbell, and Sergey Levine. Stochastic variational video prediction. *arXiv preprint arXiv:1710.11252*, 2017.

[15] Emily Denton and Rob Fergus. Stochastic video generation with a learned prior. In *International Conference on Machine Learning*, pages 1182–1191, 2018.

[16] Emily L Denton et al. Unsupervised learning of disentangled representations from video. In *Advances in neural information processing systems*, pages 4414–4423, 2017.

[17] Sergey Tulyakov, Ming-Yu Liu, Xiaodong Yang, and Jan Kautz. Mocogan: Decomposing motion and content for video generation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1526–1535, 2018.

[18] Carl Vondrick, Hamed Pirsiavash, and Antonio Torralba. Generating videos with scene dynamics. In *Advances In Neural Information Processing Systems*, pages 613–621, 2016.

[19] Pratik Gujjar and Richard Vaughan. Classifying pedestrian actions in advance using predicted video of urban driving scenes. In *2019 International Conference on Robotics and Automation (ICRA)*, pages 2097–2103. IEEE, 2019.

[20] Richard Kelley, Alireza Tavakkoli, Christopher King, Monica Nicolescu, Mircea Nicolescu, and George Bebis. Understanding human intentions via hidden markov models in autonomous mobile robots. In *Proceedings of the 3rd ACM/IEEE international conference on Human robot interaction*, pages 367–374. ACM, 2008.

[21] Qiuming Zhu. Hidden markov model for dynamic obstacle avoidance of mobile robot navigation. *IEEE Transactions on Robotics and Automation*, 7(3):390–397, 1991.

[22] Tirthankar Bandyopadhyay, Chong Zhuang Jie, David Hsu, Marcelo H Ang, Daniela Rus, and Emilio Frazzoli. Intention-aware pedestrian avoidance. In *Experimental Robotics*, pages 963–977. Springer, 2013.

[23] Vasiliy Karasev, Alper Ayvaci, Bernd Heisele, and Stefano Soatto. Intent-aware long-term prediction of pedestrian motion. In *2016 IEEE International Conference on Robotics and Automation (ICRA)*, pages 2543–2549. IEEE, 2016.

[24] Nicolas Schneider and Dariu M Gavrila. Pedestrian path prediction with recursive bayesian filters: A comparative study. In *German Conference on Pattern Recognition*, pages 174–183. Springer, 2013.

[25] David Ellis, Eric Sommerlade, and Ian Reid. Modelling pedestrian trajectory patterns with gaussian processes. In *2009 IEEE 12th International Conference on Computer Vision Workshops, ICCV Workshops*, pages 1229–1234. IEEE, 2009.

[26] Julian Francisco Pieter Kooij, Nicolas Schneider, Fabian Flohr, and Dariu M Gavrila. Context-based pedestrian path prediction. In *European Conference on Computer Vision*, pages 618–633. Springer, 2014.

[27] Benjamin Völz, Karsten Behrendt, Holger Mielenz, Igor Gilitschenski, Roland Siegwart, and Juan Nieto. A data-driven approach for pedestrian intention estimation. In *2016 IEEE 19th International Conference on Intelligent Transportation Systems (ITSC)*, pages 2607–2612. IEEE, 2016.

[28] Mahmudul Hasan, Jonghyun Choi, Jan Neumann, Amit K Roy-Chowdhury, and Larry S Davis. Learning temporal regularity in video sequences. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 733–742, 2016.

[29] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pages 1097–1105, 2012.

[30] Du Tran, Lubomir D Bourdev, Rob Fergus, Lorenzo Torresani, and Manohar Paluri. C3d: generic features for video analysis. *CoRR, abs/1412.0767*, 2(7):8, 2014.

[31] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International Conference on Machine Learning*, pages 448–456, 2015.

[32] Vincent Dumoulin and Francesco Visin. A guide to convolution arithmetic for deep learning. *arXiv preprint arXiv:1603.07285*, 2016.

[33] Ali Diba, Mohsen Fayyaz, Vivek Sharma, Amir Hossein Karami, Mohammad Mahdi Arzani, Rahman Yousefzadeh, and Luc Van Gool. Temporal 3d convnets: New architecture and transfer learning for video classification. *arXiv preprint arXiv:1711.08200*, 2017.

[34] Ameni Trabelsi, Mohamed Chaabane, and Asa Ben-Hur. Comprehensive evaluation of deep learning architectures for prediction of DNA/RNA sequence binding specificities. *Bioinformatics*, 35(14):i269–i277, 07 2019.

[35] Khurram Soomro, Amir Roshan Zamir, and M Shah. A dataset of 101 human action classes from videos in the wild. *Center for Research in Computer Vision*, 2012.

[36] John Duchi, Elad Hazan, and Yoram Singer. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of Machine Learning Research*, 12(Jul):2121–2159, 2011.

[37] Christian Schuldt, Ivan Laptev, and Barbara Caputo. Recognizing human actions: a local svm approach. In *Proceedings of the 17th International Conference on Pattern Recognition, 2004. ICPR 2004.*, volume 3, pages 32–36. IEEE, 2004.

[38] Samy Bengio, Oriol Vinyals, Navdeep Jaitly, and Noam Shazeer. Scheduled sampling for sequence prediction with recurrent neural networks. In *Advances in Neural Information Processing Systems*, pages 1171–1179, 2015.

[39] Zhou Wang, Alan C Bovik, Hamid R Sheikh, Eero P Simoncelli, et al. Image quality assessment: from error visibility to structural similarity. *IEEE transactions on image processing*, 13(4):600–612, 2004.

[40] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1–9, 2015.

[41] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.

[42] Sebastian Ruder. An overview of multi-task learning in deep neural networks. *arXiv preprint arXiv:1706.05098*, 2017.