
Temporal Coherence based Criteria for Predicting the Future using Deep Multi-stage Generative Adversarial Networks

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Predicting the future from a sequence of video frames has been recently a sought
2 after yet challenging task in the field of computer vision and machine learning.
3 Although there have been efforts for tracking using motion trajectories and flow
4 features, the problem of generating unseen frames has not been studied extensively.
5 In this paper, we deal with this problem using convolutional models inside a multi-
6 stage Generative Adversarial Networks (GAN) framework. The proposed method
7 uses two stages of GANs to generate crisp and clear set of future frames. Although
8 GANs have been used in the past for predicting the future, none of the works
9 consider relation between subsequent frames in the temporal dimension. Our main
10 contribution lies in formulating two objective functions based on the Normalized
11 Cross Correlation (NCC) and the Pairwise Contrastive Divergence (PCD). This
12 method, coupled with the traditional L1 loss, is tested on three real world video
13 datasets *viz.* Sports-1M, UCF-101 and the KITTI. Performance analysis reveals
14 superior results over the recent state-of-the-art methods.

15

1 Introduction

16 Video frame prediction has always been one of the fundamental problems in computer vision as it
17 caters to a wide range of applications including self-driving cars, surveillance, robotics and inpainting.
18 However, the challenge lies in the fact that, real-world scenes tend to be complex, and predicting
19 the future events requires modelling of complicated internal representations of the ongoing events.
20 Recently, the work of [12] modeled this problem in the framework of Generative Adversarial Networks
21 (GAN). Generative models, as introduced by Goodfellow *et. al.*, [5] try to generate images from
22 random noise by simultaneously training a generator (G) and a discriminator network (D) in a process
23 similar to a zero-sum game. Mathieu *et. al.* [12] shows the effectiveness of this adversarial training
24 in the domain of frame prediction using a combination of two objective functions (along with the
25 basic adversarial loss) employed on a multi-scale generator network. This idea stems from the fact
26 that the original L_2 -loss tends to produce blurry frames. This was overcome by the use of Gradient
27 Difference Loss (GDL) [12], which showed significant improvement over the past approaches when
28 compared using similarity and sharpness measures. However, this approach, although producing
29 satisfying results for the first few predicted frames, tends to generate blurry results for predictions far
30 away (~ 6) in the future.

31 In this paper, we aim to get over this tendency of producing blurry predictions by taking into account
32 the relation between consecutive frames in the temporal dimension also. We propose two objec-
33 tive functions: (a) **Normalized Cross-Correlation Loss (NCCL)** and (b) **Pairwise Contrastive**
34 **Divergence Loss (PCDL)** for effectively capturing inter-frame relationships in the GAN framework.
35 NCCL maximizes the cross-correlation between neighbourhood patches from consecutive frames

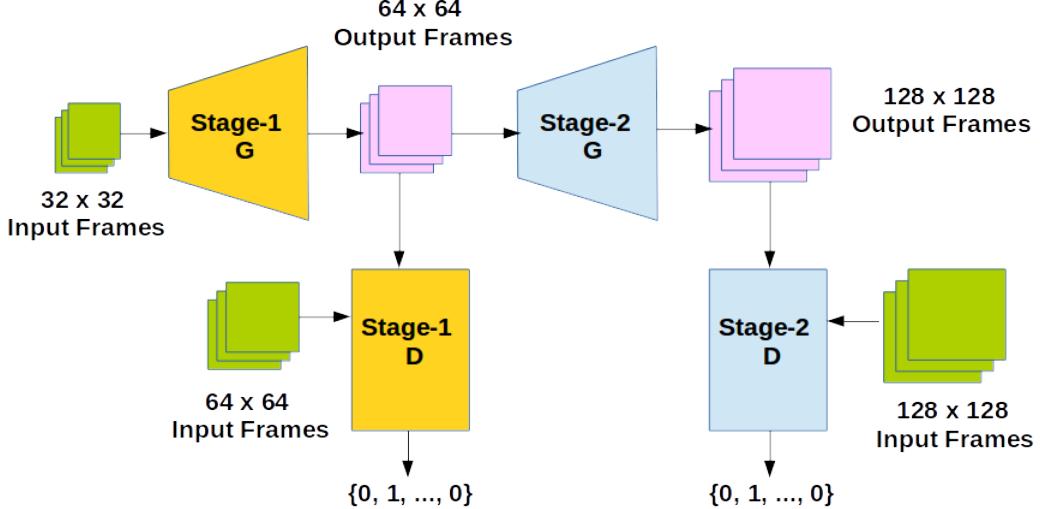


Figure 1: The proposed multi-stage GAN framework. The stage-1 generator network outputs a low-resolution version of predicted frames which are then fed through the stage-2 generator. Discriminators at both the stages predict 0 or 1 per predicted frame to denote its origin: synthetic or original.

36 whereas, CDL applies a penalty when subsequent generated frames are predicted wrongly by the
 37 discriminator network (D), thereby separating them far apart in the feature space.

38 The rest of the paper is organized as follows: section 2 describes the multi-stage generative adversarial
 39 architecture used, the sections 3 - 5 introduce the different loss funtions employed: the adversarial loss
 40 (AL), L_2 -loss (L2) and most importantly NCCL and CDL. We show the results of our experiments
 41 on Sports-1M [9], UCF-101 [16] and KITTI [4] and compare them with state-of-the-art techniques in
 42 section 6. Finally, we conclude our paper highlighting the key points and future direction of research
 43 in section 7.

44 2 Multi-stage Generative Adversarial Model

45 Generative Adversarial Networks (GAN) [5] are composed of two networks: (a) the Generator (G)
 46 and (b) the Discriminator (D). The generator G tries to generate realistic images by learning to
 47 model the true data distribution p_{data} and thereby trying to make the task of differentiating between
 48 original and generated images by the discriminator difficult. The discriminator D, in the other hand,
 49 is optimized to distinguish between the synthetic and the real images. In essence, this procedure of
 50 alternate learning is similar to the process of two player min-max games [5]. Overall, the GANs try
 51 to minimize the following objective function

$$\min_G \max_D v(D, G) = \mathbb{E}_{x \sim p_{data}} [\log(D(x))] + \mathbb{E}_{z \sim p_z} [\log(1 - D(G(z)))] \quad (1)$$

52 where, x is a real image from the true distribution p_{data} and z is a vector sampled from the distibution
 53 p_z , usually uniform or Gaussian. The adversarial loss employed in this paper is slightly different
 54 from equation 1, as, the input to our network is a sequence of frames of a video, instead of a noise
 55 vector z .

56 As convolutions account only for short-range relationships, pooling layers are used to garner in-
 57 formation from wider range. But, this process generates low resolution images. To overcome this,
 58 Mathieu *et. al.* [12] uses a multi-scale generator network, equivalent to the reconstruction process
 59 of a Laplacian pyramid [15], coupled with discriminator networks to produce high-quality output
 60 frames of size 32×32 . There are two shortcomings of this approach:

- 61 a. Generating image output at higher dimensions *viz.* (128×128) or (256×256) , requires
 62 multiple use of some static upsampling operator applied on the output of the generators. In
 63 our proposed model, this upsampling is handled by the generator models implicitly through

64 the use of consecutive unpooling operations, thereby generating predicted frames at much
 65 higher resolution in lesser number of scales.

66 b. As the generator network parameters are not learned with respect to any objective function
 67 which captures the temporal relationship effectively, the output becomes blurry after ~ 4
 68 frames.

69 To overcome the first issue, we propose a multi-stage (2-stage) generative adversarial network.

70 2.1 Stage-1

71 Generating the output frame(s) directly often produces blurry outcomes. Instead, we simplify the
 72 process by first generating crude, low-resolution version of the frame(s) to be predicted. The stage-1
 73 generator (G_1) consists of a series of convolutional layers coupled with unpooling layers [18] which
 74 upsample the frames. We used ReLU non-linearity in all but the last layer, in which case, tanh
 75 was used following the scheme of [15]. The inputs to G_1 are m number of consecutive frames
 76 of dimension $W_0 \times H_0$, whereas the outputs are n predicted frames of size $W_1 \times H_1$, where,
 77 $W_1 = W_0 \times 2$ and $H_1 = H_0 \times 2$. These outputs, stacked with the upsampled version of the original
 78 input frames, produce the input of dimension $(m + n) \times W_1 \times H_1$ to the stage-1 discriminator (D_1).
 79 D_1 applies a chain of convolutional layers followed by multiple fully-connected layers to finally
 80 produce an output vector of dimension $(m + n)$, consisting of 0's and 1's.

81 One of the key differences of our proposed GAN framework is that, the discriminator network
 82 produces decision output for multiple frames, instead of a single 0/1 outcome. This is exploited by
 83 one of the proposed objective functions, the CDL, which is described in later sections.

84 2.2 Stage-2

85 The second stage network closely resembles the stage-1 architecture, with difference only in the input
 86 and output dimensions. The input to the stage-2 generator (G_2) is formed by stacking the predicted
 87 frames and the upsampled inputs of G_1 , thereby having dimension of $(m + n) \times W_1 \times H_1$. The
 88 output of G_2 are n predicted high-resolution frames of size $W_2 \times H_2$, where, $W_2 = W_0 \times 4$ and
 89 $H_2 = H_0 \times 4$. The stage-2 discriminator (D_2), works in similar fashion as D_1 , producing output
 90 vector of length $(m + n)$.

91 Effectively, the multi-stage model can be represented by the following recursive equations:

$$Y_k = \begin{cases} G_k(Y_{k-1}, X_{k-1}), & \text{for } k \geq 2 \\ G_k(X_{k-1}) & \text{for } k = 1 \end{cases} \quad (2)$$

92 where, Y_k is the set of predicted frames, X_k are the input frames at the k th stage of the generator
 93 network G_k .

94 2.3 Training the multi-stage GAN

95 The training procedure of the multi-stage GAN model follows that of the original generative adversarial
 96 networks with minor variations. The training of the discriminator and the generator are described
 97 as follows:

98 **Training of the discriminator** Considering the input to the discriminator (D) as X (series of
 99 m frames) and the target output to be Y (series of n frames), D is trained to distinguish between
 100 synthetic and original inputs by classifying (X, Y) into class 1 and $(X, G(X))$ into class 0. Hence,
 101 for each of the stages k , we train D with target $\vec{1}$ (Vector of 1's with dimension m) for (X, Y) and
 102 target $\vec{0}$ (Vector of 0's with dimension n) for $(X, G(X))$. The loss function for training D can be
 103 described as follows:

$$\mathcal{L}_{adv}^D = \sum_{k=1}^{N_{stages}} L_{bce}(D_k(X_k, Y_k), \vec{1}) + L_{bce}(D_k(X_k, G_k(X_k)), \vec{0}) \quad (3)$$

104 where, L_{bce} , the binary cross-entropy loss can be defined as:

$$L_{bce}(Y, Y') = - \sum_{i=1}^{|Y|} Y'^i \log(Y^i) + (1 - Y'^i) \log(1 - Y^i), Y^i \in \{0, 1\}, Y'^i \in [0, 1] \quad (4)$$

105 **Training of the generator** We perform an optimization step on the generator network (G), keeping
 106 the weights of D fixed, by feeding a set of consecutive frames X sampled from the training data with
 107 target \vec{Y} (set of ground-truth output frames) and minimize the following adversarial loss:

$$\mathcal{L}_{adv}^G = \sum_{k=1}^{N_{stages}} L_{bce}(D_k(X_k, G_k(X_k)), \vec{1}) \quad (5)$$

108 By minimizing the above two losses (eqn. 3, 5), G tries to make the discriminator believe that, the
 109 source of the generated frames is the input data space itself. Although this alternate optimization of
 110 D and G is theoretically correct, in practical purposes, this produces an unstable system where G can
 111 produce samples that consecutively move far away from the original input space and in consequence
 112 D distinguishes them easily. To overcome this instability inherent in the GAN principle, and to make
 113 much clear and crisp predicted frames at high resolution, we add two more objective functions: (a)
 114 Normalized Cross Correlation Loss (NCCL) and (b) Pairwise Contrastive Divergence Loss (PCDL) to
 115 the original adversarial loss (refer to eqn. 3 and 5).

116 3 Normalized Cross-Correlation Loss (NCCL)

117 The main advantage of video over image data is the fact that, it offers a far richer space of data
 118 distribution by adding the temporal dimension along with the spatial one. Convolutional Neural
 119 Networks (CNN) can only capture short-range relationships, a small part of the vast available
 120 information, from the input video data, that too in the spatial domain. Although this can be somewhat
 121 alleviated by the use of 3D convolutions [8], this increases the number of learnable parameters by a
 122 large scale. Normalized cross-correlation has been used since long time in the field of video analytics
 123 [1, 2, 14, 11] to model the space-time relationships present in the data.

124 Normalized cross correlation (NCC) measures the similarity of two image patches as a function of
 125 the displacement of one relative to the other. This can be mathematically defined as

$$NCC(f, g) = \sum_{x,y} \frac{(f(x, y) - \mu_f)(g(x, y) - \mu_g)}{\sigma_f \sigma_g} \quad (6)$$

126 where, $f(x, y)$ is a subimage, $g(x, y)$ is the template to be matched, μ_f, μ_g denotes the mean of
 127 the subimage and the template respectively and σ_f, σ_g denotes the standard deviation of f and g
 128 respectively.

129 In the domain of video frame(s) prediction, we incorporate the NCC by first extracting small non-
 130 overlapping square patches of size $h \times h$ ($1 < h \leq 4$), denoted by a 3-tuple $P_t\{x, y, h\}$, where, x
 131 and y are the co-ordinates of the top-left pixel of a particular patch, from the predicted frame at time t
 132 and then calculating the cross-correlation score with the patch extracted from the ground truth frame
 133 at time $(t - 1)$, represented by $\hat{P}_{t-1}\{x - 2, y - 2, h + 4\}$.

134 In simpler terms, we aim to find the cross-correlation score between a small portion of the current
 135 predicted frame and the local neighborhood of that in the previous ground-truth frame. We assume
 136 that, the motion present in the full frame can be subdivided into smaller parts and can be effectively
 137 approximated by looking into small local neighborhoods in the temporal dimension. This stems
 138 from the fact that, unless the video contains significant jitter or unexpected random events like scene
 139 change, the motion features remain smooth over time. The step-by-step process for finding the
 140 cross-correlation score by matching local patches of predicted and ground truth frames is described
 141 in algorithm 1.

142 The idea of calculating the NCC score can be effectively modeled into an objective function for the
 143 generator network G , where, it tries to maximize the score over a batch of inputs. In essence, this
 144 objective function tries to model the temporal data distribution by smoothing the local motion features
 145 generated by the convolutional model. This loss function, \mathcal{L}_{NCC} can be defined as

$$\mathcal{L}_{NCC} = \sum_{batch=1}^N Score_{NCC} \quad (7)$$

146 where, $batch$ denotes a minibatch of input frames and $Score_{NCC}$, obtained using algorithm 1, is the
 147 average normalized cross-correlation score per batch. The generator tries to maximize \mathcal{L}_{NCC} along
 148 with the adversarial loss defined in section 2.

149 We also propose a variant of this objective function, termed as Smoothed Normalized Cross-Correlation
 150 Loss (SNCCL), where the patch similarity finding logic of NCCL is extended by convolving with
 151 Gaussian filters to suppress small sudden motion pattern changes. This algorithm is discussed in sec.
 152 A of the supplementary document for space restrictions.

Algorithm 1: Calculation of the normalized cross-correlation score for finding similarity between a set of predicted frame(s) and a set of ground-truth frame(s).

Input: Ground-truth frames (GT), Predicted frames ($PRED$)

Output: Cross-correlation score ($Score_{NCC}$)

```

1 Variables:
2  $h$  = height and width of an image patch
3  $t$  = current time;
4 Initialize:  $Score_{NCC} = 0$ ;
5 for  $t = 1$  upto  $T$  do
6   for  $i = 0$  upto  $H$ ,  $i \leftarrow i + h$  do
7     for  $j = 0$  upto  $H$ ,  $j \leftarrow j + h$  do
8        $P_t \leftarrow extract\_patch(PRED_t, i, j, h);$ 
9       \\ Extracts a patch from the predicted frame at time  $t$  of dimension  $h \times h$  starting from
10      the top-left pixel index  $(i, j)$ ;
11       $\hat{P}_{t-1} \leftarrow extract\_patch(GT_{t-1}, i - 2, j - 2, h + 4);$ 
12      \\ Extracts a patch from the ground-truth frame at time  $(t - 1)$  of dimension
13       $(h + 4) \times (h + 4)$  starting from the top-left pixel index  $(i - 2, j - 2)$ ;
14       $\mu_{P_t} \leftarrow avg(P_t);$ 
15       $\mu_{\hat{P}_{t-1}} \leftarrow avg(\hat{P}_{t-1});$ 
16       $\sigma_{P_t} \leftarrow standard\_deviation(P_t);$ 
17       $\sigma_{\hat{P}_{t-1}} \leftarrow standard\_deviation(\hat{P}_{t-1});$ 
18       $Score_{NCC} \leftarrow Score_{NCC} + Absolute\left(\sum_{x,y} \frac{(P_t(x,y) - \mu_{P_t})(\hat{P}_{t-1}(x,y) - \mu_{\hat{P}_{t-1}})}{\sigma_{P_t} \sigma_{\hat{P}_{t-1}}}\right);$ 
19    end
20  end
21 end
22  $Score_{NCC} \leftarrow avg(Score_{NCC});$ 
```

153 4 Pairwise Contrastive Divergence Loss (PCDL)

154 As discussed in sec. 3, the proposed method tries to capture motion features that vary slowly over
 155 time. The NCCL aims to achieve this using local similarity measures. To complement this in a global
 156 scale, we use the idea of pairwise contrastive divergence over the input frames. The idea of exploiting
 157 this temporal coherence for learning motion features has been studied in the recent past [6, 7, 13].

158 By assuming that, motion features vary slowly over time, we describe \hat{Y}_t and \hat{Y}_{t-1} as a temporal
 159 pair, where, \hat{Y}_t and \hat{Y}_{t-1} are the predicted frames at time t and $(t - 1)$ respectively, if the outputs of
 160 the discriminator network D for both these frames are 1. With this notation, we model the slowness
 161 principle of the motion features using an objective function as

$$\begin{aligned} \mathcal{L}_{PCDL} &= \sum_{i=0}^T D_\delta(\hat{Y}_i, \hat{Y}_{i+1}, p_i) \\ &= \sum_{i=0}^T p_i d(\hat{Y}_i, \hat{Y}_{i+1}) + (1 - p_i) \max(0, \delta - d(\hat{Y}_i, \hat{Y}_{i+1})) \end{aligned} \tag{8}$$

162 where, T is the time-duration of the frames predicted, p_i is the output decision ($p_i \in \{0, 1\}$) of the
 163 discriminator, $d(x, y)$ is a distance measure (L_2 in this paper) and δ is the margin. Equation 8, in
 164 simpler terms, tries to minimize the distance between frames that have been predicted correctly and
 165 encourages the distance in the negative case, upto a margin δ .

166 PCDL can be extended upto higher order versions, taking into account triplets or n number of
 167 predicted frames instead of the general pairwise case. We discuss about higher order versions
 168 (especially when $n = 3$) in sec. B of the supplementary document.

169 5 Combined Loss

170 Finally, we combine the objective functions described in eqn. 5 - 8 along with the general $L1$ -loss
 171 with different weights as

$$\begin{aligned} \mathcal{L}_{Combined} = & \lambda_{adv} \mathcal{L}_{adv}^G(X, Y) + \lambda_{L1} \mathcal{L}_{L1}(X, Y) \\ & + \lambda_{NCCL} \mathcal{L}_{NCCL}(X, Y) + \lambda_{PCDL} \mathcal{L}_{PCDL}(X, Y) \end{aligned} \quad (9)$$

172 For the sake of simplicity, all the weights *viz.* λ_{L1} , λ_{NCCL} and λ_{PCDL} have been set as 1, while
 173 λ_{adv} equals 0.05 for all the experimental studies. This overall loss is minimized during the training
 174 stage of the multi-stage GAN using Adam optimizer [10].

175 6 Experiments

176 Experimental studies of our video frame(s) prediction model have been carried out on video clips
 177 from Sports-1m, UCF-101 [16] and KITTI [4]. The input-output configuration used for training
 178 the system is as follows: **input:** 4 frames and **output:** 6 frames. We compare our results to recent
 179 state-of-the-art methods by computing two popular metrics: (a) **Peak Signal to Noise Ratio (PSNR)**
 180 and (b) **Structural Similarity Index Measure (SSIM)** [17].

181 6.1 Datasets

182 **Sports-1M** A large collection of sports videos collected from YouTube spread over 487 classes.
 183 The main reason for chosing this dataset is the amount of movement in the frames. Being a collection
 184 of sports videos, this has sufficient amount of motion present in most of the frames, making it an
 185 efficient dataset for training the prediction model.

186 **UCF-101** This dataset contains 13320 annotated videos belonging to 101 classes having 180
 187 frames/video on average. The frames in this video do not contain as much movement as the Sports-
 188 1m and hence this is used only for testing purpose.

189 **KITTI** This consists of high-resolution video data from different road conditions. We have taken
 190 raw data from two categories: (a) city and (b) road.

191 6.2 Architecture of the network

Network	Stage-1 (G)	Stage-2 (G)	Stage-1 (D)	Stage-2 (D)
Number of feature maps	64, 128, 256U, 128, 64	64, 128, 256, 512U, 256, 128, 64	64, 128, 256	128, 256, 512, 256, 128
Kernel sizes	5, 3, 3, 3, 5	5, 5, 5, 5, 5	3, 5, 5	7, 5, 5, 5, 5
Fully connected	N/A	N/A	1024, 512	1024, 512

Table 1: Network architecture details. G and D represents the generator and discriminator networks respectively. U denotes an unpooling operation which upsamples an input of dimension $h \times h$ into $(h \times 2) \times (h \times 2)$.

192 The architecture details for the generator (G) and discriminator (D) networks for experimental studies
 193 is shown in table 1. All the convolutional layers except the terminal one in both stages of G are
 194 followed by ReLU non-linearity. The last layer is tied with tanh activation function. In both the
 195 stages of G , we use unpooling layers to upsample the image into higher resolution in magnitude of 2
 196 in both dimensions (height and width). The learning rate is set to 0.003 for G , which is gradually
 197 decreased to 0.0004 over time. The discriminator (D) uses ReLU non-linearities and is trained with a
 198 learning rate of 0.03. We use minibatches of 8 clips for training the overall network.

199 **6.3 Evaluation metric**

200 Assement of the quality of the predicted frames is done by two methods: (a) Peak Signal to Noise
 201 Ratio (PSNR) and (b) Structural Similarity Index Measure (SSIM). **PSNR** measures the quality of
 202 the reconstruction process through the calculation of Mean-squared error beteen the original and
 203 the reconstructed signal in logarithmic decibel scale [1]. **SSIM** is also an image similarity measure
 204 where, one of the images being compared is assumed to be of perfect quality [17].
 205 As the frames in videos are composed of foreground and background, and in most cases the back-
 206 ground is static (not the case in the KITTI dataset, as it has videos taken from camera mounted on
 207 a moving car), we extract random sequences of 32×32 patches from the frames with significant
 208 motion. Calculation of motion presence is done by the use of optical flow method of *Brox et. al.* [3].

209 **6.4 Comparison**

210 We compare the results by testing on videos from UCF-101 using model trained on the Sports-
 211 1M dataset in table 2. Superiority of our method over the most recent work [12] can be clearly
 212 detrmined from the comparison. We followed similar choice of test set videos as in [12] to make a
 213 fair comparison. One of the impressive facts in our model is that, it can produce acceptably good
 214 predictions even in the 4th frame, which is a significant result, considering the compared approach
 215 uses separate smaller models for achieving this feat. It is to be noted that, even though the metrics
 216 for the first predicted frame do not differ by a large margin compared to the results from [12], the
 217 values decrease much slowly for the models trained with the proposed objective functions (rows 6-8
 218 of table 2). Our assumption for this phenomenon is the incorporation of the objective functions based
 219 on temporal relations also rather than only learning in the spatial domain.

Methods	1st frame prediction scores		2nd frame prediction scores		4th frame prediction scores	
	PSNR	SSIM	PSNR	SSIM	PSNR	SSIM
L1	28.7	0.88	23.8	0.83	-	-
GDL L1	29.4	0.90	24.9	0.84	-	-
GDL L1*	29.9	0.90	26.4	0.87	-	-
Adv + GDL fine-tuned*	32.0	0.92	28.9	0.89	-	-
Adv + NCCL + L1	35.4	0.94	32.9	0.92	22.7	0.64
Adv + NCCL + PCDL	37.1	0.95	34.5	0.92	23.4	0.69
Adv + NCCL + PCDL + L1	37.3	0.95	34.7	0.92	23.6	0.69
Adv + SNCCL + PCDL + 3- PCDL	38.2	0.95	35.8	0.92	24.2	0.69

Table 2: Comparison of different methods for the UCF-101 dataset. The first four rows report the results from [12]. (*) indicates models fine tuned on patches of size 64×64 [12]. (-) denotes unavailability of data. GDL stands for Gradient Difference Loss [12].

220 We also trained our model on the KITTI dataset and report the findings in table I of the supplementary
 221 document due to space restrictions.
 222 Finally, we show the prediction results obtained on both the UCF-101 and KITTI in fig. 2. The
 223 results show only the first predicted frame for clear visibility. It is evident from the subfigures that,
 224 our proposed objective functions produce impressive quality frames while the models trained with L1
 225 loss tends to output blurry reconstruction.

226 **7 Conclusion**

227 In this paper, we modified the Generative Adversarial Networks (GAN) framework with the use
 228 of unpooling operations and introduced two objective functions based on the normalized cross-
 229 correlation and the contrastive divergence estimate, in the domain of video frame(s) prediction.
 230 Studies show clear improvement of the proposed methods over the recent best methods. These
 231 objective functions can be used with more complex networks involving 3D convolutions and recurrent
 232 neural networks. In the future, we aim to learn weights for the cross-correlation such that it focuses
 233 adaptively on areas involving varying amount of motion.

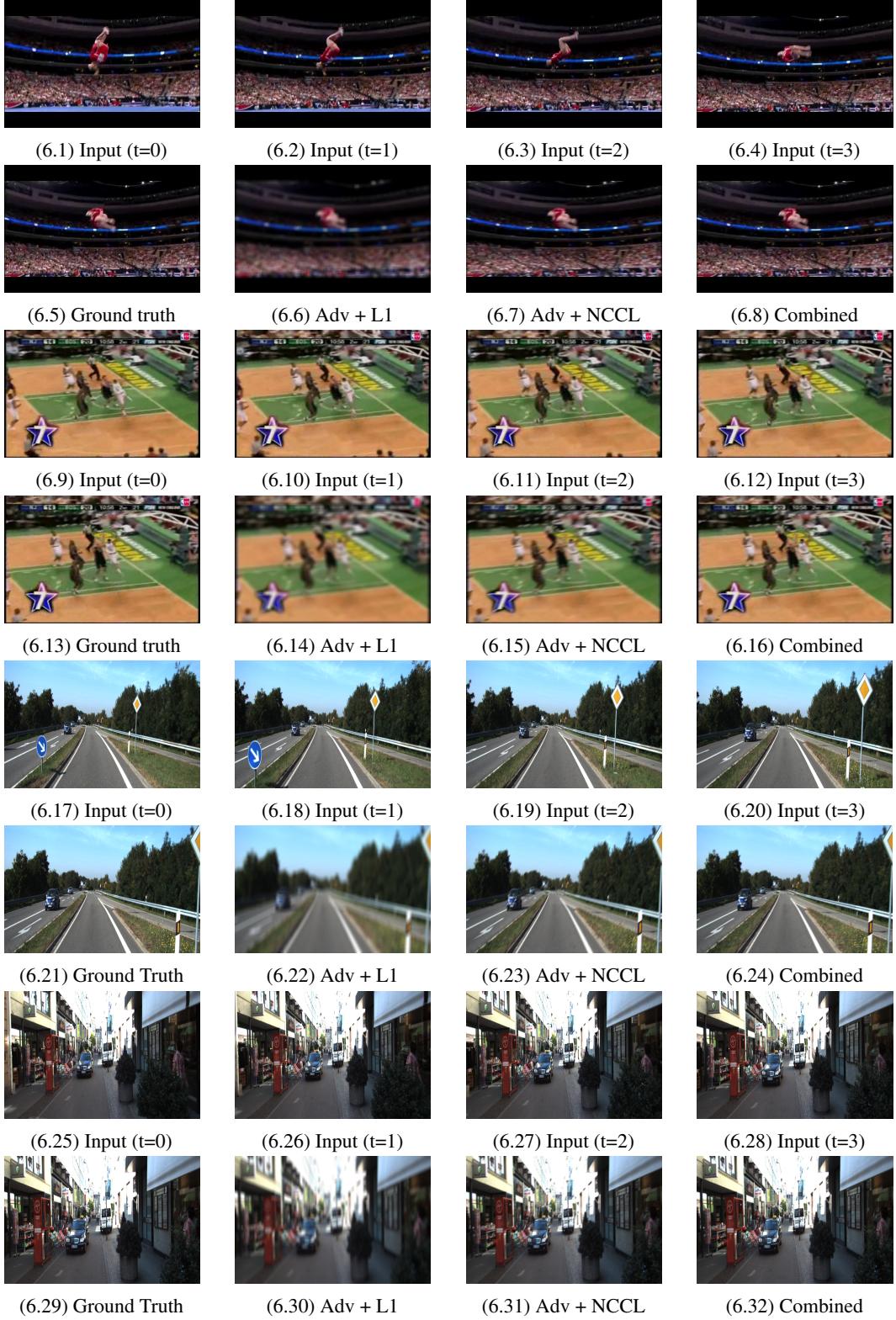


Figure 2: Comparison of applying different objective functions for training. Figures (i)-(xvi) are from UCF-101 [16], whereas, (xvii)-(xxxii) are from the KITTI [4] dataset. '*Combined*' stands for the combined loss described in section 5. Best viewed in color.

234 **References**

- 235 [1] A. C. Bovik. *The essential guide to video processing*. Academic Press, 2009.
236 [2] K. Briechle and U. D. Hanebeck. Template matching using fast normalized cross correlation. In
237 *Aerospace/Defense Sensing, Simulation, and Controls*, pages 95–102. International Society for Optics and
238 Photonics, 2001.
239 [3] T. Brox and J. Malik. Large displacement optical flow: descriptor matching in variational motion estimation.
240 *IEEE transactions on pattern analysis and machine intelligence*, 33(3):500–513, 2011.
241 [4] A. Geiger, P. Lenz, C. Stiller, and R. Urtasun. Vision meets robotics: The kitti dataset. *International*
242 *Journal of Robotics Research (IJRR)*, 2013.
243 [5] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio.
244 Generative adversarial nets. In *Advances in neural information processing systems*, pages 2672–2680,
245 2014.
246 [6] R. Goroshin, J. Bruna, J. Tompson, D. Eigen, and Y. LeCun. Unsupervised learning of spatiotemporally
247 coherent metrics. In *Proceedings of the IEEE International Conference on Computer Vision*, pages
248 4086–4093, 2015.
249 [7] R. Hadsell, S. Chopra, and Y. LeCun. Dimensionality reduction by learning an invariant mapping. In
250 *Computer vision and pattern recognition, 2006 IEEE computer society conference on*, volume 2, pages
251 1735–1742. IEEE, 2006.
252 [8] S. Ji, W. Xu, M. Yang, and K. Yu. 3d convolutional neural networks for human action recognition. *IEEE*
253 *transactions on pattern analysis and machine intelligence*, 35(1):221–231, 2013.
254 [9] A. Karpathy, G. Toderici, S. Shetty, T. Leung, R. Sukthankar, and L. Fei-Fei. Large-scale video classification
255 with convolutional neural networks. In *CVPR*, 2014.
256 [10] D. Kingma and J. Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
257 [11] J. Luo and E. E. Konofagou. A fast normalized cross-correlation calculation method for motion estimation.
258 *IEEE transactions on ultrasonics, ferroelectrics, and frequency control*, 57(6):1347–1357, 2010.
259 [12] M. Mathieu, C. Couprie, and Y. LeCun. Deep multi-scale video prediction beyond mean square error.
260 *arXiv preprint arXiv:1511.05440*, 2015.
261 [13] H. Mobahi, R. Collobert, and J. Weston. Deep learning from temporal coherence in video. In *Proceedings*
262 *of the 26th Annual International Conference on Machine Learning*, pages 737–744. ACM, 2009.
263 [14] A. Nakhmani and A. Tannenbaum. A new distance measure based on generalized image normalized cross-
264 correlation for robust video tracking and image recognition. *Pattern recognition letters*, 34(3):315–321,
265 2013.
266 [15] A. Radford, L. Metz, and S. Chintala. Unsupervised representation learning with deep convolutional
267 generative adversarial networks. *CoRR*, abs/1511.06434, 2015.
268 [16] K. Soomro, A. R. Zamir, and M. Shah. Ucf101: A dataset of 101 human actions classes from videos in the
269 wild. *arXiv preprint arXiv:1212.0402*, 2012.
270 [17] Z. Wang, A. C. Bovik, H. R. Sheikh, and E. P. Simoncelli. Image quality assessment: from error visibility
271 to structural similarity. *IEEE transactions on image processing*, 13(4):600–612, 2004.
272 [18] M. D. Zeiler and R. Fergus. Visualizing and understanding convolutional networks. In *European conference*
273 *on computer vision*, pages 818–833. Springer, 2014.