

# Sync-DRAW: Automatic Video Generation using Deep Recurrent Attentive Architectures

Gaurav Mittal\*  
gaurav.mittal.191013@gmail.com

Tanya Marwah\*  
IIT Hyderabad  
ee13b1044@iith.ac.in

Vineeth N Balasubramanian  
IIT Hyderabad  
vineethnb@iith.ac.in

## ABSTRACT

This paper introduces a novel approach for generating videos called Synchronized Deep Recurrent Attentive Writer (Sync-DRAW). Sync-DRAW can also perform text-to-video generation which, to the best of our knowledge, makes it the first approach of its kind. It combines a Variational Autoencoder (VAE) with a Recurrent Attention Mechanism in a novel manner to create a temporally dependent sequence of frames that are gradually formed over time. The recurrent attention mechanism in Sync-DRAW attends to each individual frame of the video in synchronization, while the VAE learns a latent distribution for the entire video at the global level. Our experiments with Bouncing MNIST, KTH and UCF-101 suggest that Sync-DRAW is efficient in learning the spatial and temporal information of the videos and generates frames with high structural integrity, and can generate videos from simple captions on these datasets.<sup>1</sup>

## CCS CONCEPTS

•Computing methodologies → Computer vision; Neural networks;

## KEYWORDS

Deep Learning, Video Generation, Text-to-Video Generation

## 1 INTRODUCTION

Over the years, several generative methods have been proposed to capture and model the latent representations of high-dimensional data such as images and documents. In recent years, with the success of deep learning methods, Variational Auto-Encoders (VAEs) [11] and Generative Adversarial Networks (GANs)[4] have emerged to be the most promising deep learning-based generative methods for performing tasks such as image generation [6, 16].

Existing efforts (discussed in Section 2) in the last couple of years have primarily focused on generation of images using the aforementioned methods, by modeling the spatial characteristics of a given training dataset [6, 15, 24], as well as extending these methods to more novel applications such as adding texture/style to images [3]. However, there has been very little effort in the

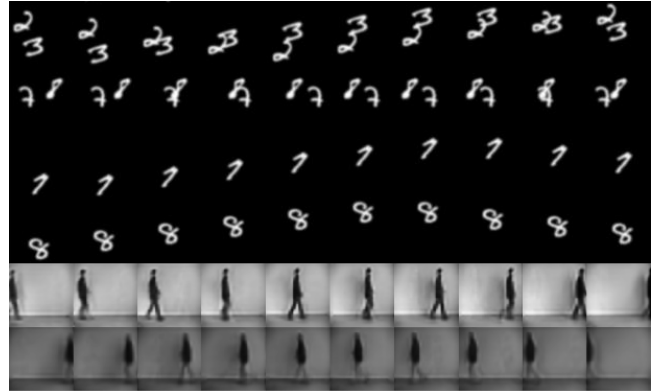


Figure 1: Videos generated automatically by proposed Sync-DRAW for Single-Digit Bouncing MNIST, Two-Digit Bouncing MNIST, and KTH Dataset.

automated generation of videos. While there have been a few efforts in video prediction [22], the recent work by Vondrick et al. [25] was the first to generate videos in an unsupervised manner using GANs. Despite acceptable performance, the generated videos had several drawbacks, as highlighted later in this work. Our work is an attempt to provide a different perspective to video generation using VAEs, which overcomes a few drawbacks of [25]; and more importantly, provides a methodology to generate videos from captions for the first time.

We propose a novel network called Sync-DRAW, which uses a recurrent VAE to automatically generate videos. It utilizes a local attention mechanism which attends to each frame individually in ‘synchronization’ and hence, the name *Sync-DRAW* (*Synchronized Deep Recurrent Attentive Writer*). Our experiments show that the proposed attention mechanism in Sync-DRAW plays a significant role in generating videos that maintain the structure of objects to a large extent. Sync-DRAW also takes a step further by learning to associate videos with captions, and thus learning to generate videos from just text captions. To the best of our knowledge, this work is the first effort to generate video sequences from captions, and is also the first to generate videos using VAEs. Figure 1 shows examples of video frames generated by trained Sync-DRAW models on different datasets.

The key contributions of this paper lie in: (i) the design of an efficient and effective deep recurrent attentive architecture to generate videos that preserves the structural integrity of objects in each frame; (ii) the adaptation of the proposed architecture to generate videos from just text captions; and (iii) a comprehensive review and experimental study of video generations on datasets of varying complexity. We now discuss earlier efforts related to this work.

\*Equal Contribution

<sup>1</sup>Accepted in ACM Multimedia 2017 to be held at Mountain View, CA, USA.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

MM’17, October 23–27, 2017, Mountain View, CA, USA.

© 2017 ACM. ISBN 978-1-4503-4906-2/17/10...\$15.00

DOI: <http://dx.doi.org/10.1145/3123266.3123309>

## 2 BACKGROUND AND RELATED WORK

In order to effectively generate meaningful structured data such as an image, a learning model should be capable of learning the latent representations for the given data, thus explaining the success of deep architectures [1], which are today synonymous with representation learning, as generative models. Initial efforts in generating such data using deep architectures were based on Deep Belief Networks [7] and Deep Boltzmann Machines [18], which contained many layers of latent variables and millions of parameters. However, their overwhelming dependence on Markov Chain Monte Carlo-based sampling methods made them difficult to scale. Over the last few years, newer deep learning approaches such as Generative Adversarial Networks (GANs) [4] and Variational Auto-Encoders (VAEs) [11] have been proven to perform quite effectively on generative tasks, and have been shown to scale to large datasets.

The Deep Recurrent Attentive Writer (DRAW) [6] was the earliest work to utilize a Recurrent-VAE (R-VAE) to learn to generate images by progressively refining an image canvas over time using a sequence of samples generated from a learnt distribution. DRAW further combines this R-VAE with an attention mechanism [14], which allows the generated image to focus on certain parts of the image at each timestep. This notion of attention resembles human visual perception, and has been used recently for related applications such as generating captions for images [27] or even generating images from captions [15, 17]. However, spatial attention differs from spatiotemporal attention, and there has been very little effort on using spatio-temporal latent representations to generate image sequences or videos, which we address in this work. Further, to the best of our knowledge, there has been no prior work in text-to-video generation, which is one of the key objectives of this work.

The state-of-the-art in the domain of video generation is a recently proposed method called VGAN by Vondrick et al. [25], which attempts to generate videos using a spatio-temporal GAN. This GAN employs a two-stream strategy with one stream focusing on the foreground, and another on the static background. While the method delivers promising results, closer observation shows that its generations are fairly noisy as well as lack ‘objectness’ (the structural integrity of the object), as shown in Section 4.7. In this work, we propose a new and different perspective to video generation based on Recurrent VAEs (instead of GANs), which uses a novel attention mechanism to generate videos. We show in this work that the videos generated by the proposed approach are perceived well by human users, as well as maintain the objects’ structural integrity to a significant extent on visual inspection. Moreover, we extend the R-VAE to perform text-to-video generation, which to the best of our knowledge, is the first effort to do so.

Other efforts that can be considered as related to the proposed work include efforts in the unsupervised domain that have attempted the task of video prediction (predict a future frame given a video sequence). One of the first efforts in this regard, [22], seeks to learn unsupervised video representations using Long Short-Term Memory units (LSTMs) [8] to predict future frames in videos. However, such efforts are fundamentally different from the objectives of this work, since they do not attempt to generate complete videos from scratch. We now present the proposed methodology.

## 3 SYNC-DRAW METHODOLOGY

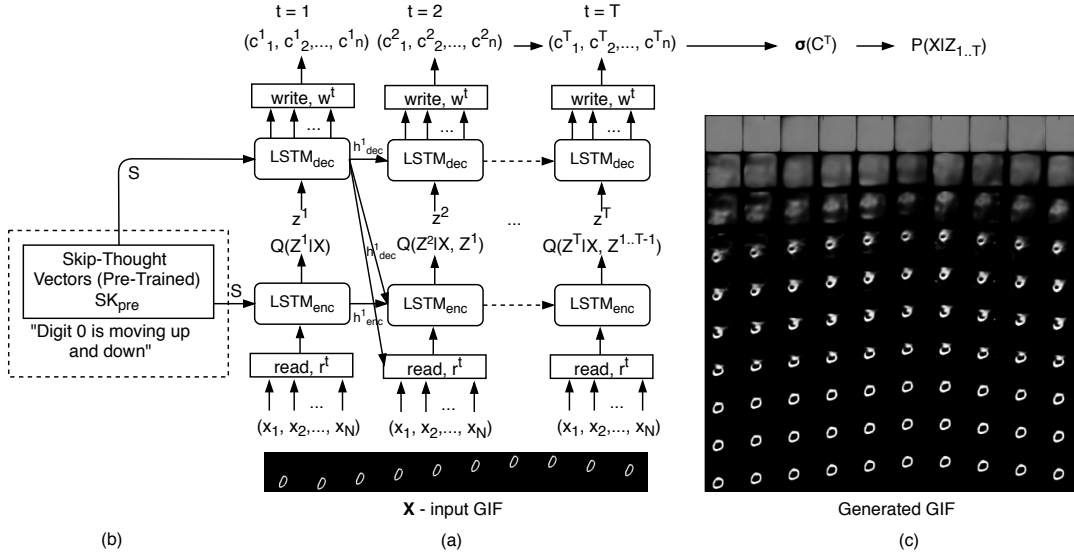
Simply speaking, a video sequence can be perceived as comprising of objects undergoing/performing action(s) across multiple frames over a background. Hence, we primarily require our model to: (i) generate a video sequence; (ii) model individual objects in the video sequence; and (iii) capture the motion of the objects in the sequence.

To generate a video as a combination of several objects and to model their motion separately, it seems intuitive to generate the video incrementally with the network focusing and forming objects one at a time. A recurrent visual attention mechanism can prove to be an effective tool to *attend* to relevant parts of a video sequence, thus allowing the network to learn latent representations of coherent parts of spatio-temporal data (objects). A video comprises of unique artifacts that are not associated with images (such as camera motion, e.g. zoom in/out), which have hitherto not been considered in image generation efforts. We hence propose a new read/write attention mechanism which operates at a frame level, and uses a grid of Gaussian filters with learnable variance and stride in each frame, which allows the network to automatically learn coherent regions of video sequences (inspired from [5, 6]). We also use a Recurrent VAE (R-VAE), which has been used for image generation in [6] (but not video generation), to embed this frame-wise attention mechanism and generate a video over a set of time steps.

A trivial approach to extend the way R-VAE is used for image generation in [6] to video generation would be to have a separate R-VAE for generating each frame. However, such an approach ignores the temporal relationship between frames, and the number of parameters to be learnt across the VAEs may explode exponentially with increasing number of frames. Therefore, to learn the video data distribution efficiently but effectively, the proposed Sync-DRAW uses a single R-VAE for the entire video (Figure 2). This global R-VAE, firstly, allows us to model the temporal relationships between the generated video frames. Secondly, it also allows us to condition the generation of the video on captions to effectively perform text-to-video generation (this would not have been easy if we trivially extended the image R-VAE to videos). Moreover, while it is possible to define attention as a spatio-temporal cuboid (3-D grid of Gaussian filters) across frame blocks, we instead define a separate individual attention mechanism for each frame of a video, each of which acts simultaneously across all frames, so as to learn latent spatio-temporal representations of video data. We found this proposed strategy of using a global R-VAE with a local frame-wise attention mechanism to be key in the results that we obtained. This approach introduces smoothness in the characteristics of objects across the frames in comparison to 3-D spatio-temporal approaches [25], which seem to have jitters in the generations, as highlighted in the experiments section (Section 4).

### 3.1 Sync-DRAW Architecture

Let  $X = \{x_1, x_2, \dots, x_N\}$  be a video comprising of frames  $x_i, i = 1, \dots, N$  where  $N$  is the number of frames in the video. Let the dimensions for every frame be denoted by  $A \times B$ . We generate  $X$  over  $T$  timesteps, where at each timestep  $t$ , canvases for all frames are generated in synchronization. The Sync-DRAW architecture comprises of: (i) a *read mechanism* which takes a region of interest (RoI) from each frame of the video; (ii) the *R-VAE*, which is responsible to learn a latent distribution for the videos from these RoIs;



**Figure 2: Sync-DRAW architecture. (a) The read-write mechanisms and the R-VAE; (b) Conditioning the R-VAE’s decoder on text features extracted from captions for text-to-video generation; (c) Canvases generated by Sync-DRAW over time (each row is a video at a particular time step, last row is the final generated video).**

and (iii) a *write mechanism* which generates a canvas focusing on a RoI (which can be different from the reading RoI) for each frame. During training, at each time step, a RoI from each frame of the original video is extracted using the read mechanism and is fed to the R-VAE to learn the latent distribution,  $Z$ , of the video data. Then, a random sample from this distribution,  $Z$ , is used to generate the frames of the video using the frame-wise write mechanisms. At test time, the read mechanism and the encoder of the R-VAE are removed; and at each time step, a sample is drawn from the learned  $Z$  distribution to generate a new video through the write mechanism. We now describe each of these components.

**3.1.1 Read Mechanism.** The read attention mechanism in Sync-DRAW is responsible for reading a region of interest (RoI) from each frame  $x_i$  at every timestep  $t$ . As mentioned earlier, the RoI is defined by a  $K \times K$  grid of Gaussian filters, where the grid’s center, grid’s stride and the Gaussian’s variance (which is the same across all grid points) are all learned while training the network. Varying these learned parameters allows us to vary the RoI read at every time step. In particular, due to the separability of the 2-D Gaussian filter, the RoI patch is read by dynamically computing 1-D Gaussian filters  $F_{i_p}^t$  and  $F_{i_q}^t$ , which correspond to a set of size  $K$  (a user-defined parameter for the grid size) of  $1 \times A$  and  $1 \times B$ -dimensional Gaussian filters respectively, for each frame  $x_i$  at timestep  $t$ , with  $i_p$  and  $i_q$  being the two spatial dimensions of frame  $x_i$ . (For convenience, we will refer to  $F_{i_p}^t$  and  $F_{i_q}^t$  as being of sizes  $K \times A$  and  $K \times B$  respectively.) These filters are evaluated together on a  $K \times K$  grid which defines the RoI of the frame to attend to. In order to compute these filters, we learn a set of read attention parameters for each frame:  $\tilde{g}_{i_p}^t, \tilde{g}_{i_q}^t, \sigma_i^t, \tilde{\delta}_i^t, \beta_i^t$  where  $i \in \{1, \dots, N\}$ .  $\tilde{g}_{i_p}^t$  and  $\tilde{g}_{i_q}^t$  correspond to the grid center coordinates along the two frame dimensions respectively; and  $\tilde{\delta}_i^t$  corresponds to the stride between the points in the grid. Evidently, we propose to learn a separate set

of these parameters for each frame in the video. In order to ensure that these parameters lie within the frame dimensions, they are slightly modified using:

$$g_{i_p}^t = \frac{A+1}{2}(\tilde{g}_{i_p}^t + 1); g_{i_q}^t = \frac{B+1}{2}(\tilde{g}_{i_q}^t + 1)$$

$$\delta_i^t = \frac{\max(A, B) - 1}{N - 1} \tilde{\delta}_i^t$$

These new parameters  $g_{i_p}^t$  and  $g_{i_q}^t$  are then used to calculate the centers for the horizontal and vertical filterbanks using:

$$\mu_{i_{p_u}}^t = g_{i_p}^t + (u - N/2 - 0.5)\delta_i^t \quad (1)$$

$$\mu_{i_{q_v}}^t = g_{i_q}^t + (v - N/2 - 0.5)\delta_i^t \quad (2)$$

for  $u, v \in \{1, \dots, K\}$ .  $\sigma_i^t$  serves as the standard deviation for all filterbanks at timestep  $t$  on  $x_i$ , the  $i^{th}$  frame.  $F_{i_p}^t$  and  $F_{i_q}^t$  are then obtained finally as follows:

$$F_{i_p}^t[u, a] = \frac{1}{Z_{i_p}} \exp\left(-\frac{(a - \mu_{i_{p_u}}^t)^2}{2(\sigma_i^t)^2}\right) \quad (3)$$

$$F_{i_q}^t[v, b] = \frac{1}{Z_{i_q}} \exp\left(-\frac{(b - \mu_{i_{q_v}}^t)^2}{2(\sigma_i^t)^2}\right) \quad (4)$$

where  $a \in \{1, \dots, A\}$ ,  $b \in \{1, \dots, B\}$ ,  $u, v \in \{1, \dots, K\}$  and  $Z_{i_p}, Z_{i_q}$  are the normalization constants.

The last parameter,  $\beta_i^t$ , allows the network to learn the temporal relationships between the patches read from the frames. It lets the model decide the level of importance to be given to each frame for generating the video at any time step. Before feeding to the R-VAE, the patch read from each frame of the input video is scaled by its respective  $\beta_i^t$ ,  $i \in \{1, \dots, N\}$  as follows:

$$read(x_i) = \beta_i^t (F_{i_p}^t) x_i (F_{i_q}^t)^T \quad (5)$$

where  $read(x_i)$  is of size  $K \times K$ . This learned set of parameters,  $\beta$ , allows Sync-DRAW to selectively focus on one (or a subset) of frames of the video at a given timestep if required.<sup>2</sup>

**3.1.2 R-VAE.** The R-VAE is responsible to generate the video sequence by sampling from the latent distribution learnt by the VAE at each time step. The core components of a standard R-VAE are the encoder LSTM,  $LSTM_{enc}$  (which outputs  $h_{enc}$ ), the latent representation,  $Z$ , and the decoder LSTM,  $LSTM_{dec}$  (which outputs  $h_{dec}$ ) [2]. The R-VAE runs for  $T$  time steps (user-defined) and generates a set of canvases  $C^t = \{c_1^t, c_2^t, \dots, c_N^t\}$  at every timestep  $t$ , where  $\sigma(C^t)$  is the video generated after  $t$  timesteps. Since the R-VAE works at the global video level, the canvases for all the frames in the video are generated in synchronization, which by our observation, preserve the inter-frame temporal relationships. At every  $t$ , we define a new quantity,  $\hat{X}^t$ , which represents the video-level error, as follows:

$$\hat{X}^t = X - \sigma(C^{t-1}) \quad (6)$$

where  $\sigma$  is the sigmoid function, and  $X$  is the input video. The LSTM-encoder at time  $t$  is then processed as follows:

$$h_{enc}^t = LSTM_{enc} \left( h_{enc}^{t-1}, [R^t, h_{dec}^{t-1}] \right) \quad (7)$$

where  $R^t = [r_1^t, r_2^t, \dots, r_N^t]$  with  $r_i^t = [read(x_i), read(\hat{x}_i^t)]$ .  $h_{enc}^t$  is then used to estimate the parameters of latent variable  $Z^t$ , which characterizes the approximate posterior  $Q(Z^t | h_{enc}^t)$  that captures the latent representation. This is done using the reparameterization technique described in [12]. A sample  $z^t \sim Q(Z^t | h_{enc}^t)$  is then obtained and fed to the LSTM-decoder,  $LSTM_{dec}$ , producing  $h_{dec}^t$ :

$$h_{dec}^t = LSTM_{dec} \left( h_{dec}^{t-1}, z^t \right) \quad (8)$$

$h_{dec}^t$  is used to compute the write attention mechanism parameters as described below. We note that the read attention parameters at time  $t$  are learned as a linear function of  $h_{dec}^{t-1}$  and the input  $X$ , whose weights are learned during training in the read block.

**3.1.3 Write Mechanism.** As part of the write mechanism, we need two things per frame: what to write, and where to write. The former is obtained using an attention window,  $w^t$ , which is computed as a linear function of  $h_{dec}^t$ , whose weights are learnt during training. The attention window is, once again, defined by an equivalently similar set of parameters as the read mechanism, which are now computed using  $h_{dec}^t$  (as against using  $h_{dec}^{t-1}$  and  $X$  for the read parameters). Next, similar to the read mechanism, a set of filterbanks  $\hat{F}_{i_p}^t$  and  $\hat{F}_{i_q}^t$  are computed using these write attention parameters, defining the RoI to be attended to while writing into each frame. Finally, the canvas  $c_i^t$  corresponding to every frame  $i$  at time step  $t$  is created as follows:

$$write(h_{dec}^t) = \frac{1}{\hat{\beta}_i^t} (\hat{F}_{i_p}^t)^T w_i^t (\hat{F}_{i_q}^t) \quad (9)$$

$$c_i^t = c_i^{t-1} + write(h_{dec}^t) \quad (10)$$

The canvasses are accumulated over timesteps  $T$ , denoted by  $c^T$  and the final video generated is computed as  $\sigma(c^T)$ . By using

<sup>2</sup> $\sigma$ ,  $\delta$ , and  $\beta$  are defined in the logarithm scale to ensure that values are always positive.

sigmoid function over  $c^T$ , we can consider the output as the emission probabilities making it equivalent to  $P(X|Z_1 \dots T)$ . In Figure 2, the final arrow is used to represent the connection between the technical implementation of the algorithm and the mathematical interpretation of the underlying model.

**3.1.4 Sync-DRAW with captions.** To generate videos from captions, we use a pre-trained model of skip-thought vectors,  $SK_{pre}$  [13], and pass the caption,  $Y$ , through this model to generate a latent representation for the caption,  $S$ .

$$S = SK_{pre}(Y) \quad (11)$$

We then condition  $LSTM_{enc}$  and  $LSTM_{dec}$  on  $S$  [20] to modify equations 7 and 8, and generate corresponding hidden representations as:

$$h_{enc}^t = LSTM_{enc} \left( h_{enc}^{t-1}, [R^t, h_{dec}^{t-1}, S] \right) \quad (12)$$

$$h_{dec}^t = LSTM_{dec} \left( h_{dec}^{t-1}, [z^t, S] \right) \quad (13)$$

At test time, when the R-VAE’s encoder and read mechanisms are removed, a caption  $Y$  is provided as input to the R-VAE’s decoder to condition the generation on the given caption,  $Y$ .

**3.1.5 Loss function.** The loss function for Sync-DRAW is composed of two types of losses (similar to a standard VAE [12]), both of which are computed at the video level. The first is the reconstruction loss,  $L_X$ , computed as the binary pixel-wise cross-entropy loss between the original video  $X$  and the generated video,  $\sigma(C^T)$ . The second is the KL-divergence loss,  $L_Z$ , defined between a latent prior  $P(Z^t)$  and  $Q(Z^t | h_{enc}^t) \sim \mathcal{N}(\mu^t, (\sigma^t)^2)$  and summed over all  $T$  timesteps. We assume prior  $P(Z^t)$  as a standard normal distribution and thus,  $L_Z$  is given by:

$$L_Z = \frac{1}{2} \left( \sum_{t=1}^T \mu_t^2 + \sigma_t^2 - \log \sigma_t^2 \right) - T/2 \quad (14)$$

The final loss is the sum of the two losses,  $L_X$  and  $L_Z$ .

**3.1.6 Testing.** During the testing phase, the encoder is removed from Sync-DRAW and due to the reparameterization trick of the VAE [12],  $z$  is sampled from  $\mathcal{N}(0, I)$ , scaled and shifted by the learned  $\mu^t$  and  $(\sigma^t)^2$ , to generate the final videos.

## 4 EXPERIMENTS AND RESULTS

We studied the performance of Sync-DRAW on the following datasets with varying complexity: (i) Single-Digit Bouncing MNIST (which has been used in similar earlier efforts [9, 22, 26]), (ii) Two-digit Bouncing MNIST; and (iii) KTH Human Action Dataset [19]. Figure 3 shows sample videos from these datasets. Considering this is the first work in text-to-video and there is no dataset yet for such a work, we chose these datasets since they provide different complexities, and also allow for adding captions<sup>3</sup>. We manually created text captions for each of these datasets (described later) to demonstrate Sync-DRAW’s ability to generate videos from captions. We varied the number of timesteps  $T$  and  $K$  for read and write attention parameters across the experiments based on the size of the

<sup>3</sup>The codes, the captioned datasets and other relevant materials are available at <https://github.com/Singularity42/Sync-DRAW>

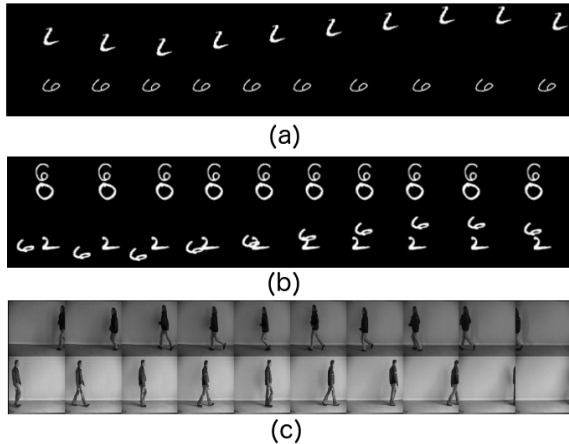


Figure 3: A few sample videos from each of the dataset used for evaluation. (a) Single-Digit Bouncing MNIST. (b) Two-Digit Bouncing MNIST. (c) KTH Dataset.

frame and complexity of the dataset (grayscale or RGB). Stochastic Gradient Descent with Adam [10] was used for training with initial learning rate as  $10^{-3}$ ,  $\beta_1$  as 0.5 and  $\beta_2$  as 0.999. Additionally, to avoid gradient from exploding, a threshold of 10.0 was used.

#### 4.1 Baseline Methodology

For performance comparison and to highlight the significance of having a separate set of attention parameters for each frame of the video, we designed another methodology to extend [6] to generate videos, by modeling attention parameters as spatio-temporal cuboids and adding a set of filterbanks  $F_Z^t$  of size  $K \times N$  to operate over the temporal dimension (in addition to the two filterbanks in the spatial dimensions). We use this methodology as our baseline for comparison. We tried on several settings of parameters and present the best results for the baseline methodology in this paper.

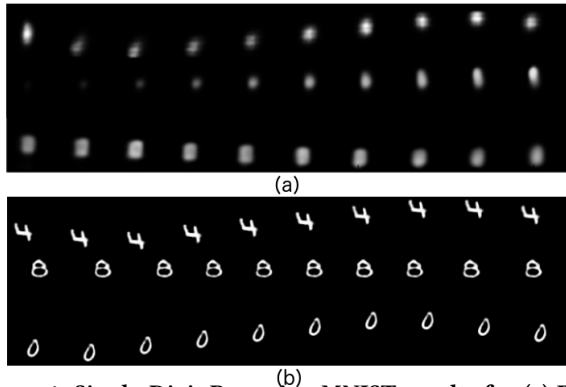


Figure 4: Single-Digit Bouncing MNIST results for (a) Baseline method and (b) Sync-DRAW

#### 4.2 Results on Single-Digit Bouncing MNIST

As in earlier work [9, 22, 26], we generated the Single-Digit Bouncing MNIST dataset by having the MNIST handwritten digits move over time across the frames of the sequence (Figure 3a). Each video contains 10 frames each of size  $64 \times 64$  with a single  $28 \times 28$  digit either moving *left and right* or *up and down*. The initial position

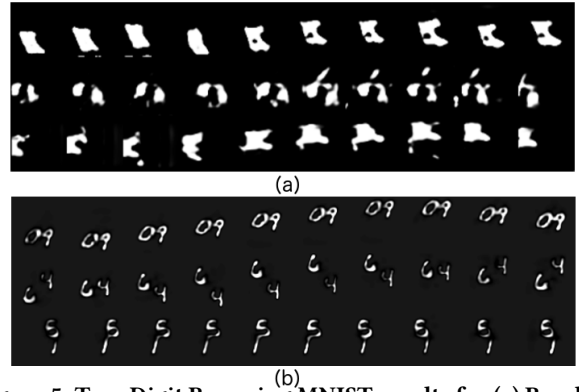


Figure 5: Two-Digit Bouncing MNIST results for (a) Baseline method and (b) Sync-DRAW

for the digit in each video is chosen randomly in order to increase variation among the samples. The training set contains 12,000 such videos. The results of Sync-DRAW on Single Digit MNIST are shown in Figures 2(c) and 4(b). The figures illustrate the usefulness of the proposed methodology in gracefully generating the final video as a sequence of canvases. When compared to the baseline methodology results, as shown in Figure 4(a), the quality of the generated digits is clearly superior with the digits having well-defined structure and boundaries. An interesting observation in Figure 2(c) is that while each frame has its own attention mechanism in the proposed framework, it can be seen that the same region of the digit is being attended to at every timestep  $t$ , even though they are present at different pixel locations in each frame, thus validating the proposed approach to video generation.

#### 4.3 Results on Two-Digit Bouncing MNIST

Extending the Bouncing MNIST dataset, we generated the two-digit Bouncing MNIST dataset where two digits move independent of one another, either *up and down* or *left and right*, as shown in Figure 3(b). The dataset consists of a training set of 12,000 videos, with each video containing 10 frames of size  $64 \times 64$  with two  $28 \times 28$  digits. In order to ensure variability, the initial positions for both the digits were chosen randomly. In case of an overlap, the intensities are added, clipping the sum if it goes beyond 1. The results of applying Sync-DRAW to this two-digit Bouncing MNIST dataset are shown in Figures 5(b) and 6(a). Once again, the baseline method performs rather poorly as shown in Figure 5(a). Figure 6(a) also shows that Sync-DRAW attends to both the digits at every time step simultaneously. In the initial time steps, a single structure is formed which then breaks down to give the two digits in the subsequent time steps. We believe that this is the reason that though the generated digits have a well defined structure and a boundary, they are still not as clear as when compared to results on Single-Digit Bouncing MNIST. We infer that there is a need for a “stimulus” for attention mechanism to know that there are two digits that need to be attended to. This claim is substantiated in subsequent sections where we include alignment with captions in Sync-DRAW, giving the attention mechanism this required stimulus.

#### 4.4 Results on KTH dataset

We evaluated the performance of Sync-DRAW on the KTH Human Action Database [19] to benchmark the approach on a more real

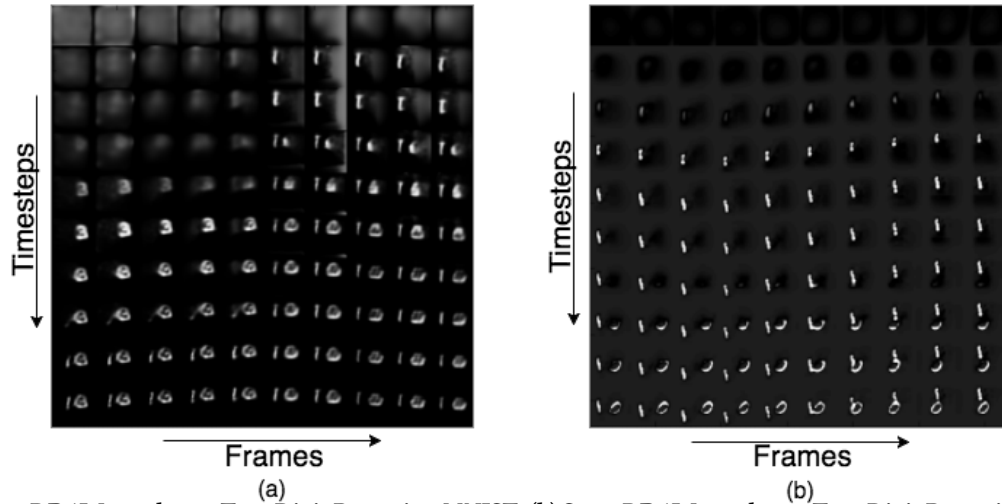


Figure 6: (a) Sync-DRAW results on Two-Digit Bouncing MNIST; (b) Sync-DRAW results on Two-Digit Bouncing MNIST videos when captions are included. Clearly, captions improve generation of ‘objectness’ of the digits.

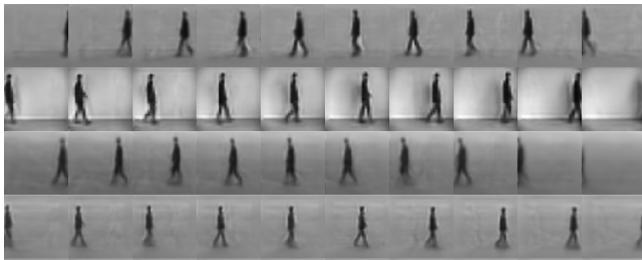


Figure 7: Videos generated by Sync-DRAW for KTH Human Action Dataset without captions.

dataset. We chose KTH dataset over other video datasets such as UCF as the videos from KTH were accompanied with meta-data which helped us to create meaningful captions for the videos and experiment Sync-DRAW with captions on this dataset. The KTH dataset consists of over 2000 videos of 25 subjects performing six different actions - walking, running, jogging, hand-clapping, hand-waving and boxing. We resized the frames to  $120 \times 120$  and performed generation for 10 frames. Given the complexity of the scene and increased frame size, we increased  $K$  to 12 and  $T$  to 32 to obtain videos with better resolution.

Figure 7 shows some test videos generated for the KTH dataset. The clarity of the person in the generated scenes suggests that Sync-DRAW is able to perform equally well on a real-world dataset.

#### 4.5 Sync-DRAW with captions

As discussed in Section 3.1.4, the proposed Sync-DRAW methodology can be used to generate videos from captions.

**4.5.1 Single-Digit Bouncing MNIST.** For every video in the bouncing MNIST dataset, a sentence caption describing the video was included in the dataset. For Single-Digit Bouncing MNIST, the concomitant caption was of the form ‘digit 0 is moving left and right’ or ‘digit 9 is moving up and down’. Hence, for the single-digit dataset, we have 20 different combinations of captions. In order to challenge Sync-DRAW, we split our dataset in such a way that for all the digits, the training and the test sets contained different motions

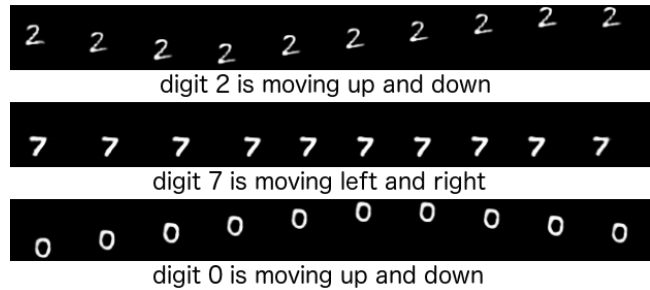


Figure 8: Sync-DRAW generates videos from just captions on the Single-Digit Bouncing MNIST: Results above were automatically generated by the trained model at test time.

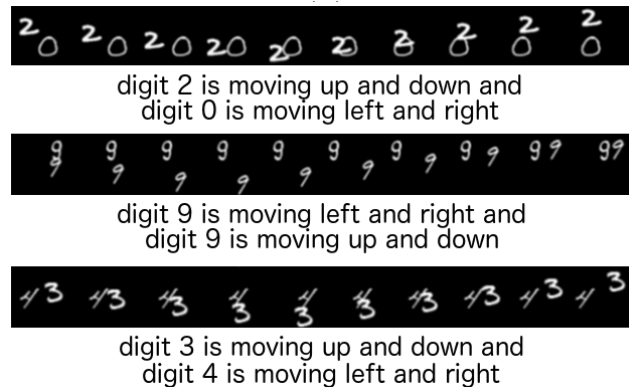
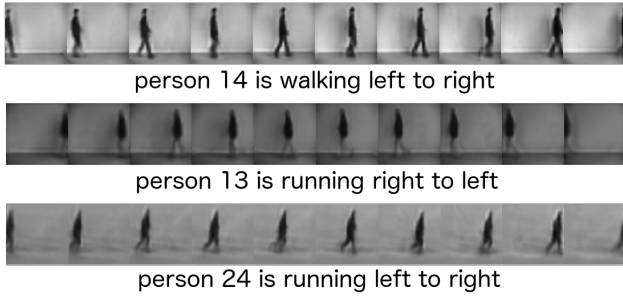


Figure 9: Sync-DRAW generates videos from just captions on the Two-Digit Bouncing MNIST: Results above were automatically generated by the trained model at test time.

for the same digit, i.e. if a digit occurs with the motion involving *up and down* in the training set, the caption for the same digit with the motion *left and right* (which is not used for training) is used in the testing phase. Figure 8 shows some of the videos generated from captions present in the test set. It can be observed that even though the caption was not included in the training phase, Sync-DRAW



**Figure 10: Sync-DRAW generates videos from just captions on KTH dataset.**

is able to capture the implicit alignment between the caption, the digit and the movement fairly well.

**4.5.2 Two-Digit Bouncing MNIST.** We conducted similar experiments for the Two-Digit Bouncing MNIST, where the captions included the respective motion information of both the digits, for example *'digit 0 is moving up and down and digit 1 is moving left and right'*. Figure 9 shows the results on this dataset, which are fairly clear and good on visual inspection (we show quantitative results later in Section 4.6). Interestingly, in Figure 6(b), we notice that when Sync-DRAW is conditioned on captions, the quality of the digits is automatically enhanced, as compared to the results in Section 4.3 (Figure 6(a)). These results give the indication that in the absence of captions (or additional stimuli), the attention mechanism in Sync-DRAW focuses on a small patch of a frame at a time, but possibly ignores the presence of different objects in the scene and visualizes the whole frame as one entity. However, by introducing captions, the attention mechanism receives the much needed "stimulus" to differentiate between the different objects and thereby cluster their generation, resulting in videos with better resolution. This is in concurrence with the very idea of an attention mechanism, which when guided by a stimulus, learns the spatio-temporal relationships in the video in a significantly better manner.

**4.5.3 KTH Dataset.** As mentioned in Section 4.4, we were able to generate descriptive captions for the videos in the KTH dataset by using the metadata which included the person and the corresponding action. We carried out our experiments on videos for walking and running as it further helped to deterministically introduce the notion of direction. Some examples of the captions are *'person 1 is walking right to left'* and *'person 3 is running left to right'*. Figure 10 shows some of the videos generated by Sync-DRAW using just the captions for KTH dataset. The generated videos clearly demonstrate Sync-DRAW's ability to learn the underlying representation of even real-world videos and create high quality videos from text.

## 4.6 Quantitative analysis

**4.6.1 Reconstruction Loss.** While quantitative analysis is difficult in image/video generation methods (due to lack of ground truth for the generations), we analyzed the quality of the videos generated by Sync-DRAW by computing the Negative Log Likelihood (NLL) of the generated samples (as recommended by Theis et al. in [23] for evaluating generative models). We compared the NLL values at convergence against the baseline approach, and a setup where, instead of a hierarchical approach, we allowed both

the spatial and temporal relationships to be learned at the global level (called Global in table below). The results are presented in Table 1. From the losses, we can clearly infer that Sync-DRAW performs significantly better than the other two approaches.

Experiments	One Digit MNIST	One Digit MNIST*	Two Digit MNIST	Two Digit MNIST*
Sync-DRAW	340.39	327.11	639.71	524.41
Global	500.01	512.46	899.06	860.54
Baseline	3561.23	3478.26	5240.65	5167.94

**Table 1: Negative Log Likelihoods at convergence for Sync-DRAW, Baseline and Global methods. \* = with captions (Lower is better).**

**4.6.2 Psychophysical Analysis.** In order to quantify the visual quality of the generated videos, we further performed a psychophysical analysis where we asked a group of human subjects to rate the videos generated by Sync-DRAW for different datasets. 37 subjects (each ignorant of this work) were given a set of 10 generated videos from each experiment we conducted with Sync-DRAW, and were asked to rate the perceived visual quality of a video on a scale of 1 – 10 (with 10 being the highest score, and also corresponded to the score for original videos from the respective datasets shown to the subjects). The statistics of this analysis are presented in Tables 2 and 3. It can be observed that the average rating for Single-Digit Bouncing MNIST (with and without captions), Two-Digit Bouncing MNIST (with captions) and KTH (with captions) is very high which correlates well with the qualitative results. The scores for Two-Digit Bouncing MNIST and KTH without captions were slightly lower (although not unacceptable) due to the reduction in clarity. Generations on these datasets, however, improved with the availability of a stimulus in the form of captions (as explained before in Figure 6).

## 4.7 Qualitative comparison with VGAN [25]

To compare our results with the recently proposed VGAN method in [25], which is the only earlier effort for video generation from scratch, we ran a comprehensive set of experiments. We first studied the performance of VGAN's [25] code on the Bouncing MNIST dataset. In order to align with VGAN's requirements, we modified the Bouncing MNIST dataset to comprise of videos containing 32

Datasets	One Digit MNIST	Two Digit MNIST	KTH Dataset
Avg. Video Quality Score	9.37 ± 0.63	7.86 ± 0.75	7.27 ± 0.80

**Table 2: Average score ± Standard Deviation given to the videos generated by Sync-DRAW without captions for different datasets. Baseline: low score (2.7 ± 0.53) for all datasets.**

Datasets	One Digit MNIST	Two Digit MNIST	KTH Dataset
Avg. Video Quality Score	9.43 ± 0.60	8.97 ± 0.83	9.10 ± 0.77

**Table 3: Average score ± Std Deviation given to videos generated by Sync-DRAW with captions for different datasets.**

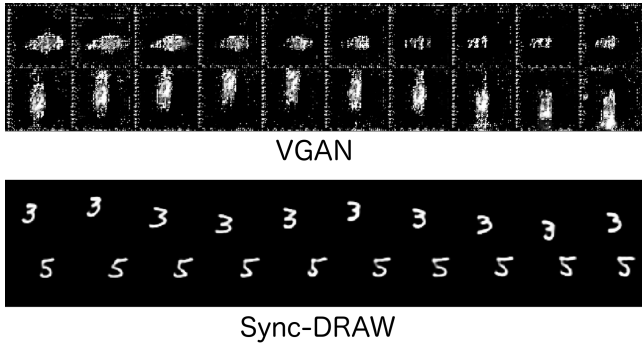


Figure 11: Qualitative comparison of video generation on MNIST dataset between VGAN and Sync-DRAW. First 10 frames generated by both approaches are shown for brevity.

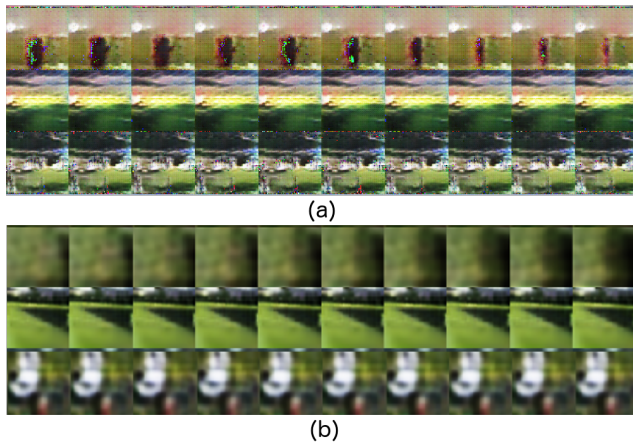


Figure 12: Qualitative comparison of video generation on dataset in [25] between VGAN and Sync-DRAW. (a) depicts the videos generated by VGAN model we trained. (b) depicts the videos generated by Sync-DRAW. The first 10 frames generated by both approaches are shown here for brevity.

frames in RGB, and then trained VGAN on this modified Bouncing MNIST dataset. To ensure an equitable comparison, we ran Sync-DRAW on the same dataset as well (note that our earlier results had different settings on this dataset). The results of the experiments are shown in Figure 11. Next, we ran Sync-DRAW on the dataset made available in VGAN’s work [25]. We changed Sync-DRAW’s parameters to support generation of videos containing 32 frames. The results can be seen in Figure 12. Lastly, in order to have a fair comparison, we ran both the codes on a dataset that they were never run on before, the UCF-101 dataset [21]. Each video in the training dataset from UCF-101 contained 32 colored frames.<sup>4</sup> These results can be seen in Figure 13.

From our experiments, we notice that there are significant differences between Sync-DRAW and VGAN results. While VGAN gives an overall semblance of sharper generations, on closer look, the generations are noisy with poor objectness. (While [25] reported

<sup>4</sup>VGAN [25] uses UCF-101 to show their network’s capability to learn representations for action classification but has shown no generation result on UCF-101

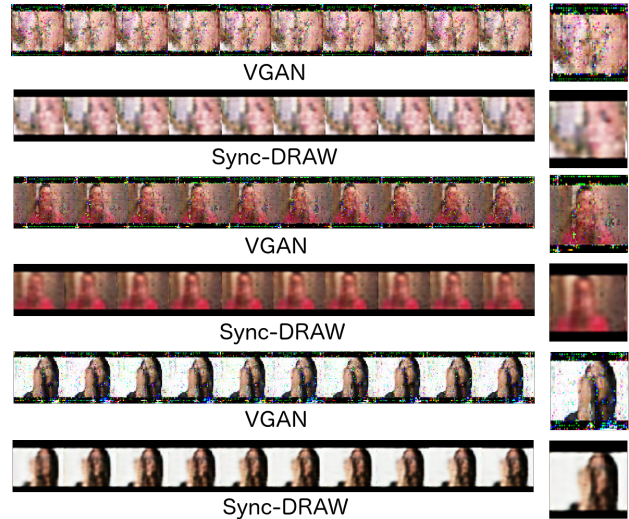


Figure 13: Qualitative comparison of video generation on UCF101 dataset between VGAN and Sync-DRAW. The first 10 frames generated by both approaches are shown here for brevity. On the right, a sample frame of each video is shown magnified for better comparison.

results that look better, we found this to be true on closer observation for many results there too). On the other hand, Sync-DRAW’s generations have a smoothness to them and are fairly noise-free in comparison, although they are a bit blurry. The results of VGAN on Bouncing MNIST highlight this issue of poor objectness further. Similar observations can also be made on the results generated from training on VGAN [25]’s own dataset. The results on UCF-101 serve as a fair comparison between the two methods because of the lack of any bias in the dataset used to train the two models. We also conducted a psychophysical analysis on the videos generated by VGAN, and found the scores on VGAN to be comparable to that of Sync-DRAW. Sharpening the generations of Sync-DRAW on real-world datasets, while maintaining its overall smoothness, is an important direction of our future work.

## 5 CONCLUSION AND FUTURE WORK

This paper presents Sync-DRAW, a new approach to automatically generate videos using a Variational Auto-Encoder and a Recurrent Attention Mechanism combined together. We demonstrate Sync-DRAW’s capability to generate increasingly complex videos starting from Single-Digit Bounding MNIST to more complex videos from KTH and UCF101 datasets. We show that our approach gives adequate focus to the objects via the attention mechanism and generates videos that maintain the structural integrity of objects. We also demonstrated Sync-DRAW’s ability to generate videos using just captions, which is the first of its kind. As evident from the results, our approach is highly promising in the domain of video generation. Our ongoing/future efforts are aimed at improving the quality of the generated videos by reducing blurriness and generating videos from even more complex captions.



## REFERENCES

- [1] Yoshua Bengio. 2009. Learning deep architectures for AI. *Foundations and trends® in Machine Learning* 2, 1 (2009), 1–127.
- [2] Otto Fabius and Joost R van Amersfoort. 2014. Variational recurrent auto-encoders. *arXiv preprint arXiv:1412.6581* (2014).
- [3] Leon A Gatys, Alexander S Ecker, and Matthias Bethge. 2015. A neural algorithm of artistic style. *arXiv preprint arXiv:1508.06576* (2015).
- [4] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. 2014. Generative adversarial nets. In *Advances in Neural Information Processing Systems*. 2672–2680.
- [5] Alex Graves, Greg Wayne, and Ivo Danihelka. 2014. Neural Turing machines. *arXiv preprint arXiv:1410.5401* (2014).
- [6] Karol Gregor, Ivo Danihelka, Alex Graves, Danilo Jimenez Rezende, and Daan Wierstra. 2015. DRAW: A recurrent neural network for image generation. *arXiv preprint arXiv:1502.04623* (2015).
- [7] Geoffrey E Hinton, Simon Osindero, and Yee-Whye Teh. 2006. A fast learning algorithm for deep belief nets. *Neural computation* 18, 7 (2006), 1527–1554.
- [8] Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. *Neural computation* 9, 8 (1997), 1735–1780.
- [9] Samira Ebrahimi Kahou, Vincent Michalski, and Roland Memisevic. 2015. RATM: Recurrent Attentive Tracking Model. *arXiv preprint arXiv:1510.08660* (2015).
- [10] Diederik Kingma and Jimmy Ba. 2015. Adam: A method for stochastic optimization. In *International Conference on Learning Representations (ICLR)*.
- [11] Diederik P Kingma and Max Welling. 2013. Auto-encoding variational bayes. *International Conference on Learning Representations* (2013).
- [12] Diederik P Kingma and Max Welling. 2013. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114* (2013).
- [13] Ryan Kiros, Yukun Zhu, Ruslan R Salakhutdinov, Richard Zemel, Raquel Urtasun, Antonio Torralba, and Sanja Fidler. 2015. Skip-thought vectors. In *Advances in neural information processing systems*. 3294–3302.
- [14] Hugo Larochelle and Geoffrey E Hinton. 2010. Learning to combine foveal glimpses with a third-order Boltzmann machine. In *Advances in neural information processing systems*. 1243–1251.
- [15] Elman Mansimov, Emilio Parisotto, Jimmy Lei Ba, and Ruslan Salakhutdinov. 2016. Generating images from captions with attention. *International Conference on Learning Representations* (2016).
- [16] Alec Radford, Luke Metz, and Soumith Chintala. 2015. Unsupervised representation learning with deep convolutional generative adversarial networks. *International Conference on Learning Representations* (2015).
- [17] Scott Reed, Zeynep Akata, Xinchun Yan, Lajanugen Logeswaran, Bernt Schiele, and Honglak Lee. 2016. Generative adversarial text to image synthesis. In *Proceedings of The 33rd International Conference on Machine Learning*, Vol. 3.
- [18] Ruslan Salakhutdinov and Geoffrey E Hinton. 2009. Deep Boltzmann Machines.. In *AISTATS*, Vol. 1. 3.
- [19] Christian Schuldt, Ivan Laptev, and Barbara Caputo. 2004. Recognizing human actions: A local SVM approach. In *Pattern Recognition, 2004. ICPR 2004. Proceedings of the 17th International Conference on*, Vol. 3. IEEE, 32–36.
- [20] Kihyuk Sohn, Honglak Lee, and Xinchun Yan. 2015. Learning structured output representation using deep conditional generative models. In *Advances in Neural Information Processing Systems*. 3483–3491.
- [21] Khurram Soomro, Amir Roshan Zamir, and Mubarak Shah. 2012. UCF101: A dataset of 101 human actions classes from videos in the wild. *arXiv preprint arXiv:1212.0402* (2012).
- [22] Nitish Srivastava, Elman Mansimov, and Ruslan Salakhutdinov. 2015. Unsupervised learning of video representations using LSTMs. In *International Conference on Machine Learning (ICML)*.
- [23] Lucas Theis, Aäron van den Oord, and Matthias Bethge. 2015. A note on the evaluation of generative models. *arXiv preprint arXiv:1511.01844* (2015).
- [24] Aaron van den Oord, Nal Kalchbrenner, and Koray Kavukcuoglu. 2016. Pixel Recurrent Neural Networks. *arXiv preprint arXiv:1601.06759* (2016).
- [25] Carl Vondrick, Hamed Pirsiavash, and Antonio Torralba. 2016. Generating videos with scene dynamics. In *Advances In Neural Information Processing Systems*. 613–621.
- [26] SHI Xingjian, Zhourong Chen, Hao Wang, Dit-Yan Yeung, Wai-kin Wong, and Wang-chun Woo. 2015. Convolutional LSTM network: A machine learning approach for precipitation nowcasting. In *Advances in Neural Information Processing Systems*. 802–810.
- [27] Kelvin Xu, Jimmy Ba, Ryan Kiros, Kyunghyun Cho, Aaron Courville, Ruslan Salakhutdinov, Rich Zemel, and Yoshua Bengio. 2015. Show, Attend and Tell: Neural Image Caption Generation with Visual Attention. In *Proceedings of the 32nd International Conference on Machine Learning (ICML-15)*. 2048–2057.