

DeepNAG: Deep Non-Adversarial Gesture Generation

Mehran Maghoubi
NVIDIA, University of Central Florida
mehran@cs.ucf.edu

Eugene M. Taranta II
University of Central Florida
etaranta@gmail.com

Joseph J. LaViola Jr.
University of Central Florida
jjl@cs.ucf.edu

ABSTRACT

Synthetic data generation to improve classification performance (data augmentation) is a well-studied problem. Recently, generative adversarial networks (GAN) have shown superior image data augmentation performance, but their suitability in gesture synthesis has received inadequate attention. Further, GANs prohibitively require simultaneous generator *and* discriminator network training. We tackle both issues in this work. We first discuss a novel, device-agnostic GAN model for gesture synthesis called *DeepGAN*. Thereafter, we formulate *DeepNAG* by introducing a new differentiable loss function based on dynamic time warping and the average Hausdorff distance, which allows us to train *DeepGAN*'s generator without requiring a discriminator. Through evaluations, we compare the utility of *DeepGAN* and *DeepNAG* against two alternative techniques for training five recognizers using data augmentation over six datasets. We further investigate the perceived quality of synthesized samples via an Amazon Mechanical Turk user study based on the $HYPE_{\infty}$ benchmark. We find that *DeepNAG* outperforms *DeepGAN* in accuracy, training time (up to 17× faster), and realism, thereby opening the door to a new line of research in generator network design and training for gesture synthesis. Our source code is available at <https://www.deepnag.com>.

CCS CONCEPTS

• Human-centered computing → Gestural input; • Computing methodologies → Neural networks.

KEYWORDS

deep neural networks, generative modeling, gesture generation, generative adversarial networks, dynamic time warping

ACM Reference Format:

Mehran Maghoubi, Eugene M. Taranta II, and Joseph J. LaViola Jr.. 2021. DeepNAG: Deep Non-Adversarial Gesture Generation. In *26th International Conference on Intelligent User Interfaces (IUI '21)*, April 14–17, 2021, College Station, TX, USA. ACM, New York, NY, USA, 11 pages. <https://doi.org/10.1145/3397481.3450675>

1 INTRODUCTION

Recently we observe that system designers are integrating gestures into almost every product with a user interface, igniting the need for accurate gesture recognizers [28, 38]. As these recognizers get

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 the author(s) 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.
IUI '21, April 14–17, 2021, College Station, TX, USA

© 2021 Copyright held by the owner/author(s). Publication rights licensed to ACM.
ACM ISBN 978-1-4503-8017-1/21/04...\$15.00
<https://doi.org/10.1145/3397481.3450675>

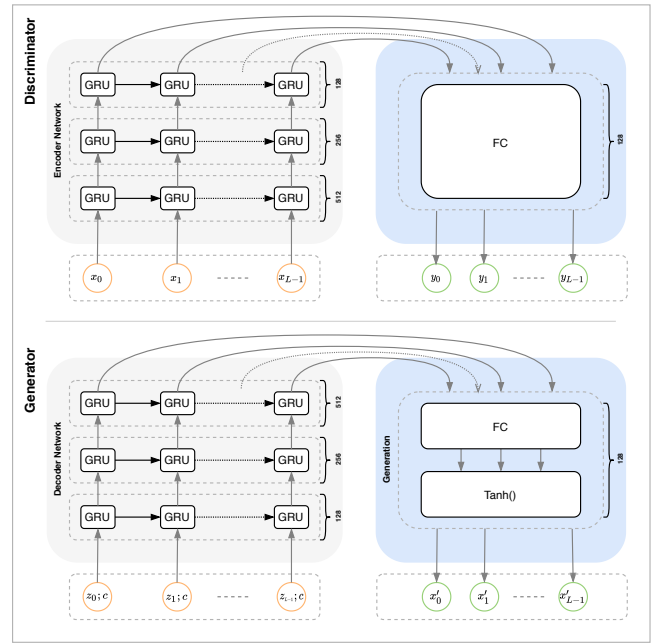


Figure 1: Our proposed model for gesture generation. *DeepGAN* consists of the discriminator and generator networks, whereas *DeepNAG* only consists of the generator network. The generator takes a *class-conditioned* random noise vector as the input and produces a gesture of the specified class. The discriminator (critic) takes the raw gesture points as the input and produces a set of features y_i used in the computation of the Wasserstein loss [14].

more sophisticated and accurate, so does their need for more data. While the size of publicly available datasets continues to grow, obtaining more task-specific data is not always easy which highlights the importance of synthetic data generation. Among such methods, generative adversarial networks (GAN)[12] have shown great promise in various problem domains [10, 32, 48], including handwriting and gesture generation [40, 50]. Typically, these networks consist of a generator and a discriminator. The discriminator aims to determine if a given example is real or fake, whereas the generator aims to fool the discriminator into confusing fake examples with the real ones. To our knowledge, such models have not received much attention towards *modality-agnostic* gesture generation wherein gestures are represented by a sequence of 2D or 3D spatio-temporal features typically produced by touch interfaces, Kinect or similar input devices. Additionally, GANs require the concurrent training of two networks, which makes the training procedure long and challenging.

This work focuses on modality-agnostic gesture generation and aims to address the challenges involved with training GANs. We start by discussing *DeepGAN*, our novel GAN approach for dynamic gesture generation, through which our deep recurrent gesture generator network is born. We thereafter discuss our unique solution for alleviating the difficulties associated with training GANs. Specifically, we formulate a novel and intuitive loss function for training our gesture generator. Our loss function, which is based on the dynamic time warping (DTW) algorithm [36], completely replaces *DeepGAN*'s discriminator network. This, transforms the significantly complex adversarial training procedure of a GAN to the much simpler *non-adversarial* training problem: train the generator by minimizing a loss function that directly maps the quality of the generated examples to their similarity to the real examples. We call this approach *DeepNAG* (see Figure 1). We evaluate both methods by using their generated gestures in data augmentation for improved gesture recognition across a variety of datasets of different sizes and modalities, as well as different gesture recognizers. We additionally conduct a user study to evaluate the human's perception of the realism of our generated results. Such evaluations, which have recently become common practice in the literature of generative modeling [22, 51], provide insight into the visual quality of generated samples.

Contributions. Our main contributions are (1) a novel recurrent GAN model for gesture generation that works across a variety of datasets and modalities, (2) a novel and intuitive loss function that completely replaces the discriminator of our GAN model, which not only simplifies and significantly speeds up the training process, but also yields a generator that produces high quality examples, (3) an evaluation of the improvements in gesture recognition accuracy when our generator is used for data augmentation.

2 RELATED WORK

Synthetic data generation is an effective approach in addressing data shortage, which in turn can improve recognition performance [11, 21, 42, 44]. Some data generation methods rely on perturbing existing samples to generate new ones. Taranta *et al.* [42] introduced GPSR which works by selecting random points along a given gesture's trajectory and scaling the between-point distances to create realistic gesture variations. Other perturbation models include the use of Perlin-noise [8] or the Sigma-Lognormal model [23, 24, 31]. These works differ from ours in that we do not rely on inputting existing gestures to produce new ones. Rather, our model generates new samples from random noise.

Generative models often involve the use of deep networks. One popular approach that predates GANs is *language modeling*, a probabilistic technique for sequence prediction which works well for handwriting generation [13]. More recently, GANs have gained popularity for such tasks [10]. Relevant examples include GestureGAN [40] a model for hand gesture-to-gesture translation. Given an image of a hand gesture and a target skeleton pose, GestureGAN produces a new hand image holding the target gesture. Yang *et al.* [47] presented a pose-guided human video generation method in which videos of a person performing a desired action are generated. Zhang's *et al.* [50] proposed a recurrent GAN for Chinese character generation which generates the temporal pen movements.

The problem domain of these works is different from ours. We focus on modality agnostic gesture generation to produce hand, full-body or 2D pen gestures. Our model learns the representation of a given gesture and produces new gestures in that same representation. Also, our generator can be trained without a discriminator, which sets us apart from the work of Zhang *et al.* [50].

Training generative models without a discriminator has also been explored in the literature. Yu *et al.* [48] and Guo *et al.* [15] proposed generating text sequences using reinforcement learning techniques. Lin *et al.* [27] presented the use of a ranking mechanism instead of a discriminator, and Li *et al.* [26] introduced an adversarial optimization procedure to train a text generator. All of these work focus on generating sequences of discrete tokens (*e.g.* text), whereas our goal is to generate real-valued and continuous multi-dimensional gesture sequences, which is highly challenging as data can take arbitrary values. Lastly, our loss formulation is different from [26] in that our formulation is inherently non-adversarial.

3 GESTURE GENERATION WITH DEEP RECURRENT NETWORKS

In this section we present our proposed deep learning approaches for gesture synthesis. We first discuss our initial GAN approach from which the *DeepNAG* generator is born. We then describe the intuition behind our loss function, followed by its formal definition.

3.1 Notations and Problem Definition

In this work, we represent gestures as a temporal sequence of input device samples (*e.g.* time-series of 3D joint positions or 2D touch coordinates), thus our data is spatio-temporal. At time step t , the gesture data is the column vector $x_t \in \mathbb{R}^N$, where N is the dimensionality of the feature vector. Thus, the entire temporal sequence of a single gesture sample is the matrix $\mathbf{x} \in \mathbb{R}^{N \times L}$, where L is the length of the sequence in time steps. For simplicity, and as typical in most gesture recognition work [41, 43], we spatially resample all gesture samples to the same length L^1 as described in [45]. We denote the vector trajectory path of a gesture \mathbf{x} with $\vec{\mathbf{x}} = \{(\overrightarrow{x_i - x_{i-1}}), \forall x_{i>0} \in \mathbf{x}\}$. Lastly, we use $|A|$ to denote the cardinality of a point set A , thus $|\mathbf{x}| = L$ and $|\vec{\mathbf{x}}| = L - 1$.

We define gesture generation as producing synthetic examples $\mathbf{x}' = \{x'_0, x'_1, \dots, x'_{L-1}\}$ over a dataset of gestures \mathbb{D} such these samples mirror data-specific properties of the samples in \mathbb{D} , as if these examples were seemingly sampled from \mathbb{D} . Formally, if $p_{\mathbb{D}}$ is dataset's distribution such that $\mathbf{x} \in \mathbb{D} \implies \mathbf{x} \sim p_{\mathbb{D}}$, our goal is to synthesize \mathbf{x}' where $\mathbf{x}' \notin \mathbb{D}$ but $\mathbf{x}' \sim p_{\mathbb{D}}$. We aim to achieve this using a deep recurrent network G which maps a *class-conditioned* latent vector $\mathbf{z}_{\hat{c}}$ to a synthetic example $\mathbf{x}' = G(\mathbf{z}_{\hat{c}}; \theta_G)$, where θ_G are the trainable parameters of G . Henceforth, we use $G_{\theta}(\mathbf{z}_{\hat{c}})$ in place of $G(\mathbf{z}_{\hat{c}}; \theta_G)$, and use \mathcal{L} to denote an objective function to minimize (a training loss function).

3.2 Gesture Generation with GANs

Our initial approach for gesture generation uses the well-known GAN training setting comprised of a generator and a discriminator.

¹We use $L = 64$

We call this recurrent model *DeepGAN*, which we designed incrementally, informed by the latest developments in deep learning. Early on, the simplicity and the recognition power of the recently proposed DeepGRU model [28] inspired us to adopt it as our discriminator. This encoder-style model has shown promising results in various recognition tasks [6, 38]. Through experiments guided by an ablation study on DeepGRU [28], and with the goal of managing design complexity, we settled for the simpler *uDeepGRU* [6] variant as our discriminator. As for our generator, we conducted experiments across different datasets with generators consisting of both LSTM and GRU units, as well as a varying number of recurrent layers. We observed stabler training, less overfitting and more plausible outputs with a decoder-style network resembling the *flipped* version of our discriminator. A possible explanation for this could be that this choice potentially benefits from the balance between the two *D* and *G* networks. Figure 1 depicts the architecture of DeepGAN, which we believe is easy to understand and straightforward to implement in any modern deep learning framework. A common design for RNN-based generators is the use of the *tanh()* activation function in the last layer, which is what we use as well.

To generate a gesture sample \mathbf{x}' of class c , the class-conditioned latent vector $\mathbf{z}_{\hat{c}}$ is fed to *G* where $\mathbf{z}_{\hat{c}}$ is defined as Equation 1. Note that each time step $z_i \in \mathbf{z}_{\hat{c}}$ is sampled independently from the standard normal distribution², and class-conditioning is done by appending the *one-hot* representation of c to each time step which avoids ignoring the conditioning through forgetting [10].

$$\mathbf{z}_{\hat{c}} = \left\{ [z_i; \hat{c}] : \forall z_i \sim \mathcal{N}(0, 1), \hat{c} = \text{one-hot}(c) \right\} \quad (1)$$

We experimented with different loss functions to train DeepGAN. Even though training with the classic adversarial loss [12] yielded plausible results, we observed improved sample quality and better convergence with the improved version [14] of the Wasserstein loss (WGAN) [2], which is what we settled on using for our evaluations. Figure 1 depicts DeepGAN’s architecture.

3.3 Non-Adversarial Gesture Generation

Although DeepGAN shows promising results (see Section 4), it demonstrated a few shortcomings early on. Most importantly, the need for training two networks simultaneously increases the training burden: changes in one network may adversely affect the other and most hyperparameters need to be tuned twice. Moreover, training times are long and we observed that the model required tens of thousands of generator iterations to converge. Aiming to reconcile these challenges, we present our loss function that completely replaces the discriminator. We start by providing an intuition for training a sequence generator without a discriminator, then proceed with the formal definition of our loss function.

Intuition. The key to simplify the network design is the answer to a fundamental question: can we possibly train the single network generator network *G*? A generator network aims to learn the distribution of the underlying dataset \mathbb{D} , so that new examples can be sampled from the distribution, which is typically done with the help of a discriminator (or a critic, in the case of WGAN-based models).

²The dimensionality of the latent space was fixed to 32 dimensions.

To simplify the gesture generation procedure, we pose the problem in a slightly different way: let us train a generator that aims to produce gestures that are *similar* to their real counterparts. The fundamental question then becomes, how to train a generator that increases the similarity between the generated and the real gestures? The answer is surprisingly simple: by reducing the *dissimilarity* between the two! Conveniently, a well-studied dissimilarity measure for time-series (as well as gestures [43]) is *dynamic time warping* (DTW) [36]. A differentiable formulation of DTW called soft DTW (sDTW) was recently proposed by Cuturi and Blondel [7].

DTW is a dynamic programming algorithm that was originally proposed for speech recognition [36]. It is a dissimilarity measure of two time-series that can be used to find their optimal alignment for various time-series analysis tasks [33]. Given two time-series $\mathbf{X} = \{x_1, x_2, \dots, x_n\}$ and $\mathbf{Y} = \{y_1, y_2, \dots, y_m\}$, a cost matrix Δ of size $n \times m$ is built. Each element Δ_{ij} is the matching cost of x_i to y_j , computed via the following recursion:

$$\Delta_{ij} = f(x_i, y_j) + \min \left\{ \Delta_{i-1, j}, \Delta_{i, j-1}, \Delta_{i-1, j-1} \right\} \quad (2)$$

where $f(x_i, y_j)$ is a problem-specific cost (distance) function. Although Euclidean distance (ED) is widely used, Taranta *et al.* [41, 43] demonstrated the superiority of using the cosine similarity measure (COS) for gesture recognition problems. Once Δ is fully computed, the value Δ_{nm} is the dissimilarity measure of \mathbf{X} and \mathbf{Y} , and the path through the matrix that yields Δ_{nm} is the optimal alignment between the two time-series. Cuturi and Blondel’s sDTW formulation [7] replaces the $\min\{\}$ operator with a differentiable (soft) minimum defined as:

$$\min^{\gamma > 0} \left\{ \Delta_1, \Delta_2, \dots, \Delta_n \right\} = -\gamma \log \sum_{i=1}^n e^{-\Delta_i/\gamma} \quad (3)$$

where γ controls the smoothness (smaller γ yields a closer approximation of classic DTW). Cuturi and Blondel show that the resulting Δ_{nm} would be the expected value of dissimilarity between \mathbf{X} and \mathbf{Y} , over every possible alignment between them weighted by their probability under the Gibbs distribution [7]. Further, they formulate the derivative of sDTW using backpropagation through a computation graph. A detailed explanation is available in [7].

Loss formulation. To increase the similarity of a real example \mathbf{x} and a fake example \mathbf{x}' , one could naively decide to learn θ by minimizing $\mathcal{L} = \text{sDTW}(G_{\theta}(\mathbf{z}_{\hat{c}}), \mathbf{x}; f)$ where f is the cost function of Equation 2. Unfortunately, this formulation merely states that generated samples \mathbf{x}' should be as similar to the real samples \mathbf{x} as possible, ignoring inter-class variations. A trivial solution to this formulation, that was easily achievable in our tests, is the per-class centroid sample.

To overcome this issue and account for inter-class variability, we formulate the loss function as a coverage measure (point set similarity) of a set of fake and real examples. For this, we propose using the *Hausdorff distance*, a well-studied measure for point set similarity [9] that can be easily implemented and computed. Conveniently, the average Hausdorff distance (denoted as $d_{\mathcal{H}}$) between two point sets \mathbf{A} and \mathbf{B} is differentiable [34] and is defined as:

$$d_{\mathcal{H}}^f(\mathbf{A}, \mathbf{B}) = \frac{1}{|\mathbf{A}|} \sum_{a \in \mathbf{A}} \min_{b \in \mathbf{B}} d(a, b; f) + \frac{1}{|\mathbf{B}|} \sum_{b \in \mathbf{B}} \min_{a \in \mathbf{A}} d(b, a; f) \quad (4)$$

where $d(a, b; f)$ is the distance (dissimilarity) between two points a and b parameterized by f . We use $d(a, b; f) = \text{sDTW}(a, b; f)$, and will discuss the choice of f (sDTW’s cost function) shortly. Finally, we propose the following as DeepNAG’s loss function to minimize. To our knowledge, this is the very first formulation of a single loss measure to train a deep recurrent neural network for generating synthetic gesture sequences:

$$\mathcal{L}_f(\mathbf{X}', \mathbf{X}) = \underbrace{d_{\mathcal{H}}^f(\mathbf{x}'_1, \mathbf{x}_1)}_{\text{Similarity term}} + \underbrace{\left| d_{\mathcal{H}}^f(\mathbf{x}'_1, \mathbf{x}'_2) - d_{\mathcal{H}}^f(\mathbf{x}_1, \mathbf{x}_2) \right|}_{\text{Variation term}} \quad (5)$$

where $\mathbf{X}' = \{\mathbf{x}'_1, \mathbf{x}'_2\}$ (two generated samples), and $\mathbf{X} = \{\mathbf{x}_1, \mathbf{x}_2\}$ (two real samples), and both sample sets belong to the same gesture class. Intuitively, Equation 5 expresses that training G should aim to increase the similarity of fake and real examples³ (similarity term), while maintaining the similarity balance between two batches of fake and real examples (variation term). The former term ensures real and generated samples are similar, while the latter term ensures generated examples maintain proper overall inter-class variations, effectively avoiding pitfalls such as mode-collapse typically encountered in GANs.

We briefly experimented with non-elastic dissimilarity functions as alternatives to sDTW in our formulation. This was motivated by the fact that our training samples were all resampled to the same length. Our experiments revealed a clear benefit in using sDTW over Euclidean distance in terms of sample generation quality. This observation is in line with prior work [43], and we posit that this is due to the elastic nature of gesture articulations.

Practical notes. During training, Equation 5 is computed over a batch of training samples, thus the number of gesture samples involved in the computation of either loss term depends on the training batch size. Also, we only used the derivative $\partial \mathcal{L}_f / \partial \mathbf{x}'_1$ to optimize the model since training was faster and more stable. This implies that the computation of $d_{\mathcal{H}}^f(\mathbf{x}'_2, \mathbf{x}_2)$ as an additional similarity term in Equation 5 is unnecessary, simply because its gradient w.r.t. \mathbf{x}'_1 is zero.

When computing $\text{sDTW}(a, b)$ we ensure *class-awareness*: the value is only computed if samples a and b belong to the same gesture class. As for the choice of sDTW’s internal cost function f , we started with ED, but the benefits of using COS quickly became apparent to us: minimizing \mathcal{L}_{ED} yielded high-quality results but convergence was slow. Conversely, minimizing \mathcal{L}_{COS} led to much faster convergence with sometimes noisier results. In the end, we settled for minimizing both and leave a thorough study on the effects of each cost function to future work. Lastly, recall our use of fixed-length sequences ($L = 64$) where the points in the sequence are equidistant. To enforce the production of such sequences by G

we add an additional term $\mathcal{L}_{\text{Resample}}$ to our objective. Putting everything together, the following is the loss function that we minimize for our experiments:

$$\begin{aligned} \mathcal{L}_{\text{DeepNAG}}(\mathbf{X}', \mathbf{X}) &= \mathcal{L}_{\text{ED}}(\mathbf{X}', \mathbf{X}) + \mathcal{L}_{\text{COS}}(\mathbf{X}', \mathbf{X}) + \alpha \cdot \mathcal{L}_{\text{Resample}}(\mathbf{x}'_1) \\ \mathcal{L}_{\text{Resample}}(\mathbf{x}') &= \frac{1}{|\mathbf{x}'|} \sum_{\forall \vec{x}'_i \in \mathbf{x}'} \left(\|\vec{x}'_i\| - \tilde{L}(\mathbf{x}') \right)^2, \quad \tilde{L}(\mathbf{x}) = \frac{1}{|\mathbf{x}|} \sum_{\forall \vec{x}_i \in \mathbf{x}} \|\vec{x}_i\| \end{aligned} \quad (6)$$

where $\tilde{L}(\mathbf{x})$ is the length of each $\vec{x}_i \in \mathbf{x}$ after \mathbf{x} is resampled to L equidistant points. Thus, $\mathcal{L}_{\text{Resample}}$ simply enforces that points in \mathbf{x}' be equidistant with α as its regularizer. Note that minimizing \mathcal{L}_{ED} alone (Equation 5) yields good-quality results in most cases. Yet, we achieved faster convergence and better data augmentation performance using Equation 3.3.

As we discuss shortly, a generator trained with our loss function demonstrates promising results compared to the same generator trained in a GAN setting. In addition to these two training configurations, we considered training our generator in a variational autoencoder (VAE) [19] setting. Although we successfully generated gesture sequences using our proposed VAE, the quality and the variety of our produced samples were sub-par compared to either DeepGAN or DeepNAG. Refer to Appendix A for more details.

4 EVALUATION

We evaluate DeepGAN and DeepNAG from two aspects. First, we conduct experiments to determine the efficacy of either model in data augmentation tasks to address data shortage. We then evaluate the perceived realism of our synthetic gestures through a user study on Amazon Mechanical Turk based on a recently introduced benchmark [51].

4.1 Data Augmentation Performance

Our experiment design for this study is as follows. Given a dataset of gestures collected from multiple participants, we simulate small training sets by splitting the data into training (50%), validation (20%) and test (30%) sets. Our experiments are all *subject-independent*; i.e. the data of each participant only appears in one of these sets. This is a more challenging and realistic evaluation protocol, as it ensures that during training, the recognizer never sees any data from the participant that it will be evaluated on during the validation and testing phases.

We begin by training a gesture recognizer on the training set, and use the validation set for model selection. We evaluate the best performing model on the test set and record its recognition error (baseline). Next, we augment the training set with a selected data generation method and repeat the experiment: train the recognizer with this new training set, use the validation data for model selection, and evaluate the best model on the test set. We record the recognizer’s recognition error again, which will be the error after augmenting the training set. Comparing this number with the baseline benchmarks the synthetic data generation method. In total we perform 150 experiments: we train five gesture recognizers on six different datasets to evaluate four synthetic data generation methods against the baseline.

³In other words *decrease their dissimilarity*

Datasets. We selected six datasets among the ones frequently studied in the literature. They vary in size and span across gesture modalities and input devices: JK2017 (Kinect) [43] (14 full-body fighting gestures of 20 participants with Kinect v2), JK2017 (Leap Motion) [43] (eight hand-gestures of 20 participants with Leap Motion), UT-Kinect [46] (ten full-body daily activities of ten participants with Kinect v1), MSR Action 3D [25] (20 full-body actions of ten participants with Kinect v1), SBU Kinect Interactions [49] (8 two-person interaction of seven participants with Kinect v1) and \$1-GDS [45] (16 2D pen gestures of ten participants).

Recognizers. We selected five gesture recognizers: support vector machine (SVM), random forest, naïve Bayes, DeepGRU [28] and Jackknife [43]. These represent classic machine learning algorithms, deep learning as well as rapid prototyping [43] approaches, which are common choices for gesture recognizers. The first three methods require explicit feature extraction for which we use the Rubine [35] feature set extended to 3D gestures [37]. Jackknife [43] is a 1-nearest neighbor DTW-based template matching recognizer.

Data generation methods. We compare four data generation methods against the baseline: random Gaussian noises, GPSR [42], DeepGAN and DeepNAG. The effectiveness of GPSR for 3D gestures has been previously demonstrated [6, 28].

Implementation. We implemented DeepGAN and DeepNAG with the PyTorch [29] framework which we have publicly released⁴. Additionally, our implementation requirements yielded multiple other standalone projects, which we have made public in the hope of benefiting the deep learning community. Inspired by [52], we implemented a CUDA version of sDTW with a PyTorch interface using Numba [20]. Our novel implementation runs more than 100× faster than any other publicly available implementation that we know of⁵. Additionally, we implemented fast GRU units using PyTorch’s just-in-time (JIT) compilation features⁶ to allow computing their higher-order derivatives, a feature that is missing in PyTorch⁷. Such derivatives are required to implement the improved WGAN loss [14] for GRUs.

Hyperparameters. All hyperparameters were tuned across different datasets, but the same set of parameters were used for every experiment. Both DeepGAN and DeepNAG were trained on the 50% split training set, and shared most hyperparameter settings. We use the Adam [18] solver ($\beta_1 = 0.5$, $\beta_2 = 0.9$), with a learning rate of 10^{-4} and a mini-batch size of 64. DeepGAN-specific hyperparameters were chosen from [14], as they performed the best in our validation runs. Other parameters were chosen via cross-validation as follows. For DeepNAG we used $\gamma = 0.1$ and $\alpha = 10^3$. GPSR parameters were set to $r = 2$, $\sigma = 0.25$ and the magnitude of random noise was set to 2% of the bounding box of each feature.

Results and discussion. Figure 2 depicts the results of our experiments⁸. In many cases the use of some form of data augmentation

decreases the recognition error, indicating that our 50% split to simulate small training sets is working as expected. To better contrast the generation methods we employ a scoring scheme that quantifies whether the use of a given augmentation method is both *warranted* and *effective*. Data augmentation is only warranted if a recognizer trained with the additional data outperforms the baseline. Additionally, a method is effective only if it outperforms random noise. We start with a score of zero for a given generator. In each experiment set, we increment this score by one if the method outperforms all other methods in addition to the baseline and random noise. Ties are only counted if the method outperforms both random noise and the baseline, and we use the cumulative score for comparison.

Table 1 presents the computed score aggregates over each dataset and recognizer. We observe that across both aggregate groups, DeepNAG outperforms other methods by a large margin, suggesting its suitability for data augmentation regardless of the choice of dataset or recognizer. In a few cases, DeepNAG reduced the recognition error to zero. Compared to GPSR, these results are notable as DeepNAG generates new examples purely from random noise. Conversely, GPSR perturbs existing examples to generate new ones. This process leaves some characteristics of the original gesture (e.g. bounding box size) largely unchanged, which benefits recognizers that rely on such features.

Figure 2 also shows cases wherein data augmentation seems harmful. In particular, we observe increased errors in almost all cases where data generation is used with multi-actor gestures (Figure 2e). This suggests that our generators may not be suitable for generating multi-actor gestures, which we confirmed by visual inspection. In some cases, both DeepGAN and DeepNAG confuse the main and the secondary actors, yielding malformed gestures. We intend to study the generation of such gestures in future work. We additionally inspected some of the generated examples of Figure 2f wherein our generators increased recognition errors. Most synthetic examples were visually fine which suggests that the use of domain adaptation techniques may be helpful [1, 30, 53]. We plan to explore this in the future.

During visual inspection, we did not observe any mode-collapse issues with DeepNAG. We observed healthy variations across all gesture classes and datasets with minimal amounts of degenerate samples (except for the few cases noted above, typically less than 1% of the generation samples were degenerate). Lastly, factors such as ease of training and training times compel the use of DeepNAG over DeepGAN as the former offers a significant reduction in training times. Training DeepGAN on a Tesla V100 GPU takes between 3-5 days depending on the dataset size, whereas DeepNAG takes around 3-7 hours under the same conditions, a speedup of 12–17×.

Overall, our results indicate that DeepNAG outperforms DeepGAN on data augmentation tasks, regardless of the choice of dataset or recognizer. These results are notable, as the generator model in both DeepGAN and DeepNAG is exactly the same. In other words, the generator which is trained using our novel loss function outperforms the same generator trained in a GAN setting with the improved Wasserstein loss. Additionally, our loss function trains the generator in a much shorter period of time.

⁴<https://github.com/Maghoumi/DeepNAG>

⁵<https://github.com/Maghoumi/pytorch-softdtw-cuda>

⁶<https://github.com/Maghoumi/JitGRU>

⁷To our knowledge, the cuDNN framework is missing this feature. Thus, at the time of this writing one cannot compute higher-order derivatives for GPU-based GRUs in any deep learning framework that relies on cuDNN.

⁸A video demo of generated gestures is available at <https://www.deepnag.com>

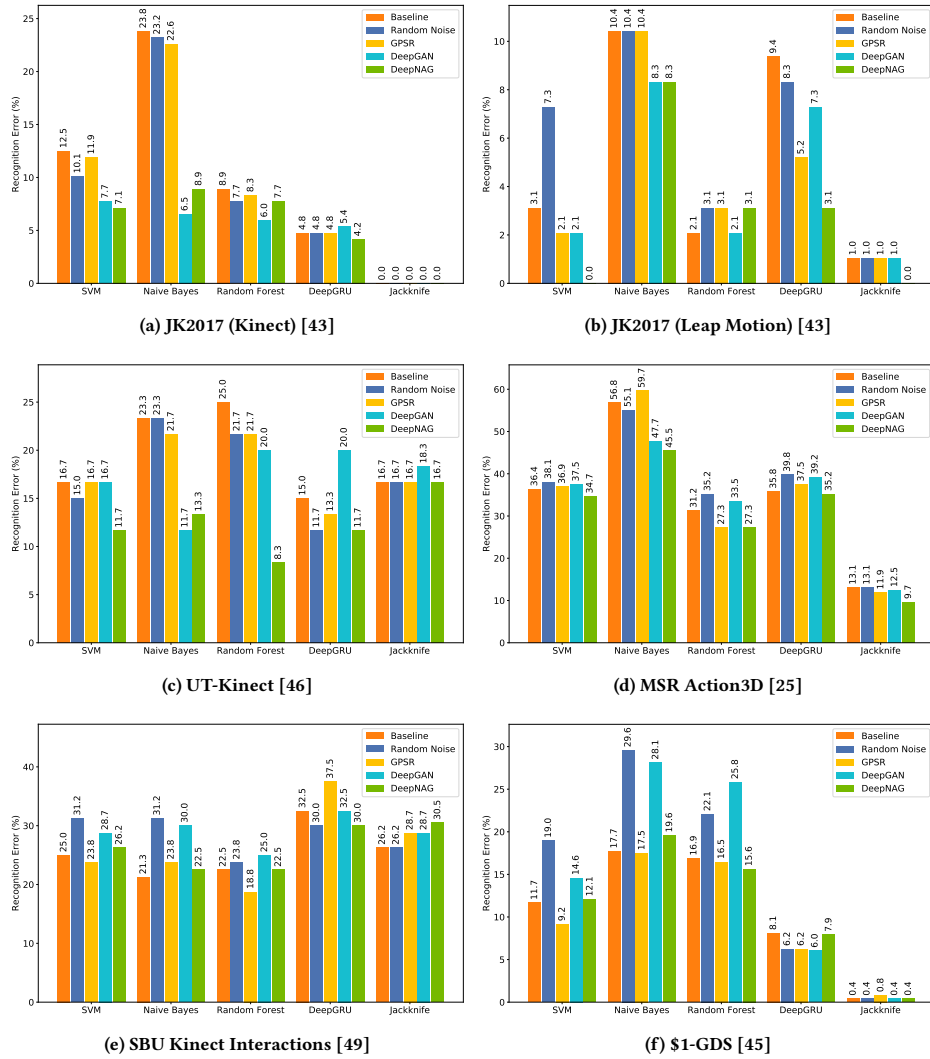


Figure 2: Results of evaluation across six datasets (best viewed in color).

Table 1: Generator scores aggregated over dataset and recognizer.

Dataset	Generator Score			Recognizer	Generator Score		
	GPSR	DeepGAN	DeepNAG		GPSR	DeepGAN	DeepNAG
JK2017 (Kinect) [43]	0	2	2	SVM	2	0	4
JK2017 (LeapMotion) [43]	0	1	4	Naïve Bayes	1	3	2
UT-Kinect [46]	0	1	2	Random Forest	2	1	3
MSR Action3D [25]	1	0	3	DeepGRU [28]	0	1	3
SBU Kinect	2	0	0	Jackknife [43]	0	0	2
\$1-GDS [45]	2	1	1				
Total Score	5	5	12		5	5	14

4.2 User Study

Qualitative evaluation of generative models through user studies has become a common practice in the literature [22, 51]. We therefore turn our focus to the comparison of DeepGAN and DeepNAG based on the perceived quality of the generated samples using human evaluators based on the HYPE_∞ [51] benchmark. This benchmark defines the gold standard for evaluating generative realism on crowd-sourcing platforms.

To compare different generative models, HYPE_∞ defines an experiment with 30 participants: each participant only sees the results from one of the models. For a given generative model, every participant is shown a total of 100 samples comprised of 50 fake and 50 real samples. Given each sample, participants are asked to indicate whether they think that sample is real or computer-generated. Participants have an infinite amount of time to make this binary choice. Afterwards, the percentage of the samples that were judged incorrectly is computed for every participant. Obtained values are averages over 30 participants and the final result is reported as the HYPE_∞ score for the understudied generative model. Zhou *et al.* [51] showed that this protocol ensures repeatability and maintains the separability between different generative models and can be used as a reliable measure of the generative model’s quality.

Using this protocol, we conduct our user study on Amazon Mechanical Turk for a given dataset \mathbb{D} and a generator G . We first train the G on \mathbb{D} to convergence. Using the trained model, we then sample as many fake samples as the real samples in \mathbb{D} . The 100 samples needed to show a given participant are selected by randomly choosing 50 fake and 50 real samples from the pool of all available samples. We recruit 30 participants for every combination of \mathbb{D} and G . Participants begin by studying the purpose of the study and answering demographic questions as detailed in Table 2. We then show them each of the 100 samples and ask them to indicate whether they think a given sample is produced by human or computer-generated (see Figure 3). Every gesture sequence is drawn in the form of a looping *gif* animation with a framerate of 32. Between each animation loop, we display a countdown with a duration of 0.25 seconds. This was inspired by [51] and was done to avoid confusing participants who may be unaware that they are watching an infinitely-looping animation. Note that the order of the samples shown to every participant is random and varies among different participants. This was done to avoid any learning effects in the experiment.

Participants are given an infinite amount of time to respond to each question. Similar to the HYPE_∞ benchmark, we reveal the correct answer to the participant upon submitting a response to every question. Participants can participate in our study only once, which ensures unique responses across all experiment conditions. At the end of the study, we reveal the overall accuracy of the participant in our task and pay them \$2 for their time [51].

When posting our study on the Mechanical Turk platform, we created a list of criteria to ensure the selection of a pool of high-quality workers. We refined and validated these criteria through trials on Mechanical Turk prior to starting our actual study. First, participants must have an approval rating of at least 97% to participate in our study to filter out low-quality workers. Our next participation requirement is that workers must have completed at

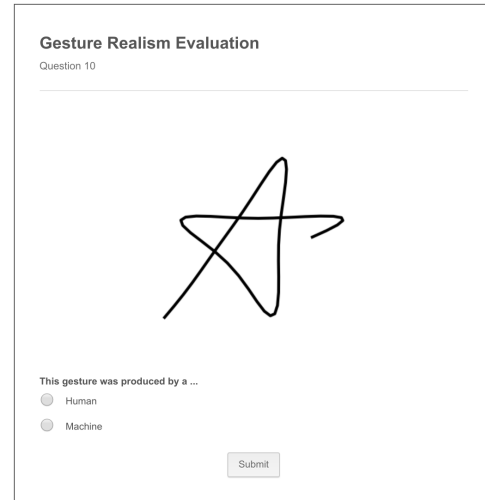


Figure 3: The interface of our user study application. Participants are shown the gesture animation and are asked to select either “human” or “machine”. Once “submit” is clicked, the correct answer is revealed.

Table 2: Pre-study questionnaire. Except for age, the remaining questions are multiple-choice.

-
- | | |
|----|---|
| Q1 | What is your gender? |
| Q2 | What is the highest level of education that you have completed? |
| Q3 | What is your age? |
-

least 5000 studies. This criterion filters out participants who may have high approval ratings because they recently joined the platform. Lastly, participants must be *Mechanical Turk Masters* to be eligible to participate in our study. Amazon uses proprietary criteria to grant top-performing workers this qualification. Although the exact criteria is not publicly disclosed, Amazon claims they continuously monitor the performance of master workers across different user studies on the platform to ensure consistent performance⁹.

Our study consisted in evaluating each of DeepGAN and DeepNAG on three datasets covering different gesture modalities: Kinect (JK-2017 [43]), Leap Motion (JK-2017 [43]) and Pen gestures (\$1-GDS [45]). Thus our *generator* factor has two levels and our *dataset* factor has three levels, yielding a total of six experiments.

Results and discussion. In total, we recruited 180 participants (30 participants for each of our six experiments) with an average age of 41 years ($\sigma=10.8$). Figure 4 depicts the demographics of our participants. Across all tasks, participants spent an average of 12.3 minutes ($\sigma=3.8$), and each question was answered in 7.4 seconds on average ($\sigma=2.3$). Considering a payment of \$2 per study, our participants were compensated well above the minimum wage specified by the United States federal guidelines (\$7.25 per hour at the time of this writing).

⁹Details available at <https://www.mturk.com/worker/help>

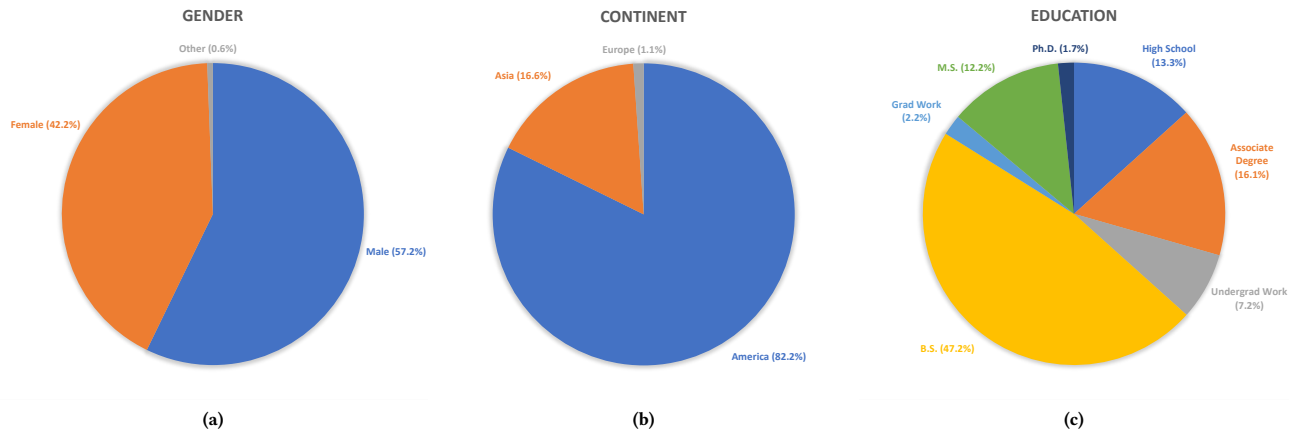


Figure 4: Demographics of our user study participants (best viewed in color).

Table 3: Amazon Mechanical Turk user study results. Reported values are percentages (averaged over 30 participants). The top performing model on each dataset (with statistical significance) is boldfaced.

Dataset	Generator	HYPE _∞	Std.	Fake Errors	Real Errors
JK2017 (Kinect) [43]	DeepNAG	48.1	8.8	53.9	42.3
	DeepGAN	38.4	12.9	44.3	32.5
JK2017 (LeapMotion) [43]	DeepNAG	51.0	4.3	56.1	45.8
	DeepGAN	22.7	11.4	23.9	21.4
\$1-GDS [45]	DeepNAG	50.0	6.7	56.3	43.7
	DeepGAN	44.4	8.3	49.5	39.3

Table 3 summarizes our user study results. In all experiments, we observe higher HYPE_∞ scores for DeepNAG compared to DeepGAN. Unpaired t-tests confirm the difference is significant in every case: $t(58)=3.3$, $p=0.001$ (JK2017-Kinect [43]), $t(58)=12.4$, $p<0.001$ (JK2017-Leap Motion [43]) and $t(58)=2.8$, $p=0.006$ (\$1-GDS [45]).

Focusing on the results with JK2017 (Leap Motion) [43] dataset, we observe a large HYPE_∞ score gap between the two generators. Notably, DeepNAG achieves hyper-realism on this dataset: its fake samples look more realistic to humans than the real ones. These results correlate well with those in Section 4.1: on the Leap Motion dataset, DeepNAG significantly outperformed DeepGAN in reducing the recognition error (Table 1) and in some cases, DeepNAG reduced the recognition error to zero (Figure 2).

Similar to [51], we report a breakdown of the error on the real and fake samples. We observe higher fake errors with DeepNAG in all cases. Inline with Zhou *et al.*'s observation [51], real and fake errors track each other. This indicates participants become more confused when fake samples are particularly hard to distinguish from the real ones.

In summary, our study shows that it is harder for evaluators to distinguish DeepNAG's synthetic samples from the real samples compared to those produced by DeepGAN. This trend holds regardless of the dataset and gesture modality. Additionally, DeepNAG not only outperformed DeepGAN in every experiment, but it also achieved hyper-realism on the Leap Motion dataset. These results

correlate well with our data augmentation performance evaluations in Section 4.1 and are notable considering that DeepNAG and DeepGAN both use the same underlying generator.

5 CONCLUSION

We discussed modality-agnostic gesture generation with recurrent neural networks. We first presented DeepGAN, our GAN model for synthetic gesture generation across various datasets and gesture modalities. To reduce the training complexity, we formulated a novel loss function based on the dynamic time warping (DTW) algorithm and the average Hausdorff distance. Our loss function obviated the need for a separate discriminator network, and led to 12–17× faster training. We called this approach DeepNAG and evaluated it from two aspects. Our first evaluations focused on the use of either model towards data augmentation for improved gesture recognition. In these evaluations, DeepNAG outperformed DeepGAN, along with other synthetic gesture generators across various datasets and recognizers. Next, we evaluated the perceived quality of the synthetic samples produced by DeepGAN and DeepNAG using human evaluators. Our user study, which was based on the HYPE_∞ benchmark and was conducted using Amazon Mechanical Turk, demonstrated that DeepNAG consistently outperformed DeepGAN in terms of the realism of the synthetic samples. Users confused DeepNAG's samples with the real ones more frequently,

and on one of our studied datasets, DeepNAG achieved hyper-realism by obtaining a HYPE $_{\infty}$ score of 51%.

In the future, we plan to more deeply explore the generation of multi-actor gestures, as well as the use of domain adaptation techniques to further improve data augmentation performance. Lastly, we aim to explore the application of our loss function in problem domains besides gestures such as time-series generation.

ACKNOWLEDGMENTS

We thank Amir Gholaminejad, Rahul Sukthankar and Corey Pittman for the helpful discussions. We also thank the anonymous reviewers for their time and valuable feedback.

REFERENCES

- [1] 2018. CyCADA: Cycle Consistent Adversarial Domain Adaptation. In *International Conference on Machine Learning (ICML)*.
- [2] Martin Arjovsky, Soumith Chintala, and Léon Bottou. 2017. Wasserstein Generative Adversarial Networks. In *Proceedings of the 34th International Conference on Machine Learning, ICML 2017, Sydney, Australia*.
- [3] Dana H. Ballard. 1987. Modular Learning in Neural Networks. In *Proceedings of the Sixth National Conference on Artificial Intelligence - Volume 1* (Seattle, Washington) (AAAI'87). AAAI Press, 279–284.
- [4] Samuel R. Bowman, Luke Vilnis, Oriol Vinyals, Andrew M. Dai, Rafal Jozefowicz, and Samy Bengio. 2015. Generating Sentences from a Continuous Space. arXiv:1511.06349 [cs.LG]
- [5] Christopher P. Burgess, Irina Higgins, Arka Pal, Loic Matthey, Nick Watters, Guillaume Desjardins, and Alexander Lerchner. 2018. Understanding disentangling in β -VAE. arXiv:1804.03599 [stat.ML]
- [6] F. M. Caputo, S. Burato, G. Pavan, T. Voillemin, H. Wannous, J. P. Vandeborste, M. Maghoumi, E. M. Taranta II, A. Razzmoo, J. J. LaViola Jr., F. Manganaro, S. Pini, G. Borghi, R. Vezzani, R. Cucchiara, H. Nguyen, M. T. Tran, and A. Giachetti. 2019. Online Gesture Recognition. In *Eurographics Workshop on 3D Object Retrieval*.
- [7] Marco Cuturi and Mathieu Blondel. 2017. Soft-DTW: a differentiable loss function for time-series. In *Proceedings of the 34th International Conference on Machine Learning - Volume 70*. JMLR. org, 894–903.
- [8] Kenny Davila, Stephanie Ludi, and Richard Zanibbi. 2014. Using off-line features and synthetic data for on-line handwritten math symbol recognition. In *Frontiers in Handwriting Recognition (ICFHR), 2014 14th International Conference on*. IEEE, 323–328.
- [9] M. . Dubuisson and A. K. Jain. 1994. A modified Hausdorff distance for object matching. In *Proceedings of 12th International Conference on Pattern Recognition*, Vol. 1. 566–568 vol.1. <https://doi.org/10.1109/ICPR.1994.576361>
- [10] Cristobal Esteban, Stephanie L. Hyland, and Gunnar Ratsch. 2017. Real-valued (Medical) Time Series Generation with Recurrent Conditional GANs. arXiv:1706.02633 [stat.ML]
- [11] Andreas Fischer, Muriel Visani, Van Cuong Kieu, and Ching Y. Suen. 2013. Generation of Learning Samples for Historical Handwriting Recognition Using Image Degradation. In *Proceedings of the 2Nd International Workshop on Historical Document Imaging and Processing* (Washington, District of Columbia, USA) (HIP '13). ACM, New York, NY, USA, 73–79.
- [12] 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.
- [13] Alex Graves. 2013. Generating Sequences With Recurrent Neural Networks. CoRR abs/1308.0850 (2013). arXiv:1308.0850
- [14] Ishaan Gulrajani, Faruk Ahmed, Martin Arjovsky, Vincent Dumoulin, and Aaron Courville. 2017. Improved Training of Wasserstein GANs. In *Proceedings of the 31st International Conference on Neural Information Processing Systems NIPS'17* (Long Beach, California, USA).
- [15] Jiaxian Guo, Sidi Lu, Han Cai, Weinan Zhang, Yong Yu, and Jun Wang. 2017. Long Text Generation via Adversarial Training with Leaked Information. arXiv preprint arXiv:1709.08624 (2017).
- [16] David Ha and Douglas Eck. 2017. A Neural Representation of Sketch Drawings. *ArXiv e-prints* (April 2017). arXiv:1704.03477 [cs.NE]
- [17] G. E. Hinton and R. R. Salakhutdinov. 2006. Reducing the Dimensionality of Data with Neural Networks. *Science* 313, 5786 (2006), 504–507. <https://doi.org/10.1126/science.1127647>
- [18] Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980 (2014).
- [19] Diederik P Kingma and Max Welling. 2013. Auto-Encoding Variational Bayes. arXiv:1312.6114 [stat.ML]
- [20] Siu Kwan Lam, Antoine Pitrou, and Stanley Seibert. 2015. Numba: A LLVM-Based Python JIT Compiler. In *Proceedings of the Second Workshop on the LLVM Compiler Infrastructure in HPC (LLVM '15)*. ACM, Article 7, 6 pages.
- [21] Do-Hoon Lee and Hwan-Gue Cho. 1998. A new synthesizing method for handwriting Korean scripts. *International Journal of Pattern Recognition and Artificial Intelligence* 12, 01 (1998), 45–61.
- [22] Luis A. Leiva. 2017. Large-Scale User Perception of Synthetic Stroke Gestures. In *Proceedings of the 2017 Conference on Designing Interactive Systems* (Edinburgh, United Kingdom). Association for Computing Machinery, 1135–1140.
- [23] Luis A. Leiva, Daniel Martín-Albo, and Réjean Plamondon. 2015. Gestures À Go Go: Authoring Synthetic Human-Like Stroke Gestures Using the Kinematic Theory of Rapid Movements. *ACM Trans. Intell. Syst. Technol.* 7, 2 (Nov. 2015), 15:1–15:29.
- [24] Luis A. Leiva, Daniel Martín-Albo, and Réjean Plamondon. 2017. The Kinematic Theory Produces Human-Like Stroke Gestures. *Interacting with Computers* 29, 4 (July 2017), 552–565.
- [25] W. Li, Z. Zhang, and Z. Liu. 2010. Action recognition based on a bag of 3D points. In *2010 IEEE Computer Society Conference on Computer Vision and Pattern Recognition - Workshops*.
- [26] Zhongliang Li, Tian Xia, Xingyu Lou, Kaihe Xu, Shaojun Wang, and Jing Xiao. 2019. Adversarial Discrete Sequence Generation without Explicit Neural Networks as Discriminators. In *Proceedings of Machine Learning Research (Proceedings of Machine Learning Research, Vol. 89)*, Kamalika Chaudhuri and Masashi Sugiyama (Eds.). PMLR, 3089–3098.
- [27] Kevin Lin, Dianqi Li, Xiaodong He, Zhengyou Zhang, and Ming-Ting Sun. 2017. Adversarial Ranking for Language Generation. In *Proceedings of the 31st International Conference on Neural Information Processing Systems* (Long Beach, California, USA) (NIPS'17). 11.
- [28] Mehran Maghoumi and Joseph J. LaViola. 2019. DeepGRU: Deep Gesture Recognition Utility. In *Advances in Visual Computing*. Springer International Publishing, 16–31.
- [29] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. 2017. Automatic differentiation in PyTorch. In *NIPS-W*.
- [30] Luis Perez and Jason Wang. 2017. The Effectiveness of Data Augmentation in Image Classification using Deep Learning. arXiv:1712.04621 [cs.CV]
- [31] Réjean Plamondon and Moussa Djioua. 2006. A multi-level representation paradigm for handwriting stroke generation. *Human movement science* 25, 4 (2006), 586–607.
- [32] Alec Radford, Luke Metz, and Soumith Chintala. 2015. Unsupervised Representation Learning with Deep Convolutional Generative Adversarial Networks. arXiv:1511.06434 [cs.LG]
- [33] Thanawin Rakthanmanon, Bilson Campana, Abdullah Mueen, Gustavo Batista, Brandon Westover, Qiang Zhu, Jesin Zakaria, and Eamonn Keogh. 2012. Searching and Mining Trillions of Time Series Subsequences Under Dynamic Time Warping. In *Proceedings of the 18th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD '12)* (Beijing, China).
- [34] Javier Ribera, David Güera, Yuhao Chen, and Edward J. Delp. 2019. Locating Objects Without Bounding Boxes. *Proceedings of the Computer Vision and Pattern Recognition (CVPR)* (June 2019). Long Beach, CA.
- [35] Dean Rubine. 1991. Specifying Gestures by Example. In *Proceedings of the 18th Annual Conference on Computer Graphics and Interactive Techniques (SIGGRAPH '91)*. 329–337.
- [36] H. Sakoe and S. Chiba. 1978. Dynamic programming algorithm optimization for spoken word recognition. *IEEE Transactions on Acoustics, Speech, and Signal Processing* 26, 1 (February 1978), 43–49. <https://doi.org/10.1109/TASSP.1978.1163055>
- [37] Jia Sheng. 2003. A study of adaboost in 3d gesture recognition. *Department of Computer Science, University of Toronto* (2003).
- [38] S. Shin and W. Kim. 2020. Skeleton-Based Dynamic Hand Gesture Recognition Using a Part-Based GRU-RNN for Gesture-Based Interface. *IEEE Access* 8 (2020), 50236–50243.
- [39] Kihyuk Sohn, Honglak Lee, and Xinchen Yan. 2015. Learning Structured Output Representation using Deep Conditional Generative Models. In *Advances in Neural Information Processing Systems* 28, C. Cortes, N. D. Lawrence, D. D. Lee, M. Sugiyama, and R. Garnett (Eds.). Curran Associates, Inc., 3483–3491.
- [40] Hao Tang, Wei Wang, Dan Xu, Yan Yan, and Nicu Sebe. 2018. GestureGAN for Hand Gesture-to-Gesture Translation in the Wild. In *Proceedings of the 26th ACM International Conference on Multimedia* (Seoul, Republic of Korea) (MM '18). ACM, 774–782.
- [41] Eugene M. Taranta, II and Joseph J. LaViola, Jr. 2015. Penny Pincher: A Blazing Fast, Highly Accurate β -family Recognizer. In *Proceedings of the 41st Graphics Interface Conference (GI '15)* (Halifax, Nova Scotia, Canada).
- [42] Eugene M Taranta II, Mehran Maghoumi, Corey R Pittman, and Joseph J LaViola Jr. 2016. A rapid prototyping approach to synthetic data generation for improved 2D gesture recognition. In *Proceedings of the 29th Annual Symposium on User Interface Software and Technology*. ACM, 873–885.

- [43] Eugene M. Taranta II, Amirreza Samiei, Mehran Maghoughi, Pooya Khaloo, Corey R. Pittman, and Joseph J. LaViola Jr. 2017. Jackknife: A Reliable Recognizer with Few Samples and Many Modalities. In *Proceedings of the 2017 CHI Conference on Human Factors in Computing Systems (CHI '17)* (Denver, Colorado, USA), 12.
- [44] Tamás Varga, Daniel Kilchhofer, and Horst Bunke. 2005. Template-based synthetic handwriting generation for the training of recognition systems. In *Proceedings of the 12th Conference of the International Graphonomics Society*. 206–211.
- [45] Jacob O Wobbrock, Andrew D Wilson, and Yang Li. 2007. Gestures Without Libraries, Toolkits or Training: A \$1 Recognizer for User Interface Prototypes. In *Proceedings of the 20th Annual ACM Symposium on User Interface Software and Technology (UIST '07)*.
- [46] L. Xia, C.C. Chen, and JK Aggarwal. 2012. View invariant human action recognition using histograms of 3D joints. In *Computer Vision and Pattern Recognition Workshops (CVPRW), 2012 IEEE Computer Society Conference on*. IEEE, 20–27.
- [47] Ceyuan Yang, Zhe Wang, Xinge Zhu, Chen Huang, Jianping Shi, and Dahua Lin. 2018. Pose Guided Human Video Generation. In *Computer Vision – ECCV 2018*, Vittorio Ferrari, Martial Hebert, Cristian Sminchisescu, and Yair Weiss (Eds.). Springer International Publishing, Cham, 204–219.
- [48] Lantao Yu, Weinan Zhang, Jun Wang, and Yong Yu. 2017. SeqGAN: Sequence Generative Adversarial Nets with Policy Gradient. In *Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence* (San Francisco, California, USA) (AAAI'17). AAAI Press, 2852–2858.
- [49] Kiwon Yun, Jean Honorio, Debaleena Chattopadhyay, Tamara L. Berg, and Dimitris Samaras. 2012. Two-person Interaction Detection Using Body-Pose Features and Multiple Instance Learning. In *Computer Vision and Pattern Recognition Workshops (CVPRW), 2012 IEEE Computer Society Conference on*. IEEE.
- [50] X. Zhang, F. Yin, Y. Zhang, C. Liu, and Y. Bengio. 2018. Drawing and Recognizing Chinese Characters with Recurrent Neural Network. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 40, 4 (April 2018), 849–862. <https://doi.org/10.1109/TPAMI.2017.2695539>
- [51] Sharon Zhou, Mitchell Gordon, Ranjay Krishna, Austin Narcomey, Li F Fei-Fei, and Michael Bernstein. 2019. Hype: A benchmark for human eye perceptual evaluation of generative models. In *Advances in Neural Information Processing Systems*. 3449–3461.
- [52] H. Zhu, Z. Gu, H. Zhao, K. Chen, C. Li, and L. He. 2018. Developing a pattern discovery method in time series data and its GPU acceleration. *Big Data Mining and Analytics* 1, 4 (2018), 266–283.
- [53] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. 2017. Unpaired Image-to-Image Translation using Cycle-Consistent Adversarial Networks. In *Computer Vision (ICCV), 2017 IEEE International Conference on*.

A APPENDIX: GESTURE GENERATION WITH VARIATIONAL AUTOENCODERS

Thus far, we have shown that training our proposed RNN-based gesture generator using our novel loss function outperforms the same generator that is trained with the improved Wasserstein loss in a GAN training setting. In this section we investigate whether our generator can be trained in a variational autoencoder setting.

A.1 Background

Variational autoencoders (VAE) [19] are a class of autoencoders [3, 17] designed for generative modeling. Similar to autoencoders, VAEs consist of encoder and decoder networks. However, the goal of VAEs is to model the distribution of the input data by learning a latent representation thereof. As such, the encoder network maps the input data \mathbf{x} to a probability distribution (latent space) while the decoder network aims to reconstruct the original data from a vector \mathbf{z} in that latent space. Once training concludes, the decoder network can be used to generate synthetic samples. The loss function for VAEs consists of reconstruction as well as regularization terms. The reconstruction term ensures that the reconstructed data closely resembles the input data. The regularization term ensures that the learned distribution of the latent space is as close to some known distribution as possible (typically the standard normal distribution). Assuming that ϕ and θ denote the trainable parameters

of the encoder and decoder respectively, the following is the loss function that is minimized [5]:

$$\mathcal{L}(\theta, \phi; \mathbf{x}, \mathbf{z}) = \underbrace{-\mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{z})]}_{\text{reconstruction}} + \underbrace{D_{KL}(q_{\phi}(\mathbf{z}|\mathbf{x}) \parallel p(\mathbf{z}))}_{\text{regularization}} \quad (7)$$

where $D_{KL}(\parallel)$ is the Kullback-Leibler (KL) divergence between two probability distributions. Equation 7 simply aims to minimize the reconstruction error as well as the KL divergence between the learned latent space and the standard normal distribution ($p(\mathbf{z}) = \mathcal{N}(0, 1)$). Note that in this formulation, class labels are not considered, which means one cannot control what sample class is produced for a given \mathbf{z} . Sohn *et al.* [39] proposed conditional VAEs in which a conditioning criteria is applied to the input data \mathbf{x} as well as the latent vector \mathbf{z} similar to conditional GANs as described in Section 3.2.

A.2 Model Architecture and Objective Function

We iteratively designed our VAE’s overall architecture. Given our goal of training the generator of DeepGAN/DeepNAG in a VAE framework, we reused the aforesaid generator as the decoder in our VAE network. As for the encoder, we started with using the uDeepGRU model as the encoder. This way, our overall VAE network closely resembled that of DeepGAN’s. After running some preliminary experiments, we observed that the choice of the encoder architecture did not result in perceptible difference in the model’s performance. In fact, adding or removing layers in either the encoder or the decoder made little difference in the produced results, inline with what Bowman *et al.* [4] observed. We ultimately decided to carry on with an architecture similar to Figure 1.

We now discuss our proposed training objective function which can be used to train our RNN-based generator in a VAE framework. As previously mentioned, the VAE objective function typically consist of reconstruction and regularization terms. We can conveniently reuse the regularization term of Equation 7, as it simply ensures that the learned latent space follows the standard normal distribution. The reconstruction term, however, is domain-specific. Recall that the this term ensures that the output of the decoder (generator) closely resembles the input data. Conveniently, a differentiable measure that can be used for this purpose is sDTW. Thus, we propose the following loss function as the objective to minimize during the training our gesture generating VAE:

$$\mathcal{L}(\theta, \phi; \mathbf{x}, \mathbf{z}_c) = \text{sDTW}(\mathbf{x}, G_{\theta}(\mathbf{z}_c); f) + D_{KL}(q_{\phi}(\mathbf{z}_c|\mathbf{x}) \parallel p(\mathbf{z}_c)) \quad (8)$$

where f is sDTW’s internal cost function and \mathbf{z}_c is the latent vector conditioned on the class label c . In simple terms, we define the reconstruction error as the sDTW dissimilarity between the input data and the output of the generator. As mentioned in Section 3.3, Cuturi and Blondel show that the computed sDTW value would be the expected value of dissimilarity between two time-series, over every possible alignment between them weighted by their probability under the Gibbs distribution [7]. This closely resembles the original reconstruction term of VAEs (Equation 7). Though VAEs have previously been used to generate 2D sketches [16], to

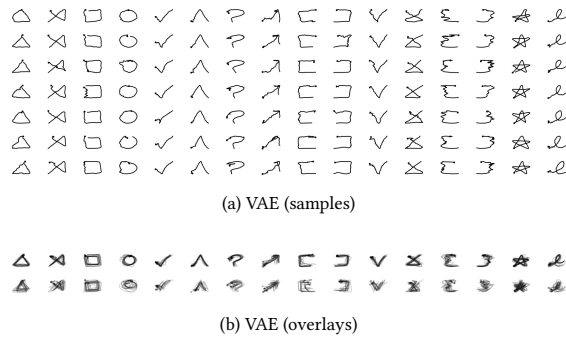


Figure 5: Synthetic gestures produced by our VAE-based generator trained on \$1-GDS [45] dataset. Note that most samples are noisy and lack visual fidelity. We further show overlaid rendering of synthetic samples from our VAE model (b – top), and real samples (b – bottom). Each overlay consists of 16 samples per class. Note the lack of variety in the synthetic results compared to the real samples.

our knowledge our proposed loss function is novel, and we are the first to formulate such function specifically for gesture generation of various modalities using VAEs.

A.3 Results

As previously mentioned, changes in the architecture of our model made little difference in the produced results. Although training plots showed a steady decrease of the loss value and training converged, the generated results lacked visual quality and diversity. Most gesture trajectories were rather noisy. Again, we observed this trend regardless of the architecture of encoder or decoder networks, various hyperparameter settings, the choice of f (we tried ED and COS) or even the gesture dataset. We additionally experimented with alternative reconstruction terms. Specifically, we experimented with the mean squared error (MSE) of both the Euclidean distance as well as the cosine similarity of gesture paths (\vec{x}) between the input and reconstructed samples. These alternate formulations performed worse than our sDTW-based reconstruction term.

Some samples produced by our VAE model when trained on the \$1-GDS dataset [45], along with overlaid samples for each of the real and synthetic data are depicted in Figure 5. These results show that the produced samples lack sufficient diversity when compared to the real samples. Given the visual quality of the results, we hypothesize that our proposed VAE framework is not suitable for training our generator to produce good synthetic samples.