

---

# Game-Theoretic Models for Generative Learning

## PhD Thesis Proposal

---

### Abstract

Machine learning has achieved great success in object recognition and classification problems based on the rapid development of deep neural networks. However, there's still a big gap between computer and human intelligence. Discriminative learning attempts to infer knowledge from data by modeling the conditional distribution of class labels given data samples, while generative learning tries to estimate the full data distribution and synthesizes new samples. It can control the generated samples by changing its class domain or modifying its visual appearance. My thesis will focus on the latter problem and investigate generative learning from a game-theoretic perspective. We first formulate the problem as a distributionally robust game with payoff uncertainty, and then develop a robust optimization algorithm to solve the Nash equilibrium. Meanwhile, we will study the distance metrics that measure the similarity between distributions, which is a major issue in many generative learning problems. We then propose a conditional generative model to solve the style transfer problem in image and speech processing. The disentangled representations of domain-specific style information and domain-invariant content information are modeled by autoencoders and domain classifiers. The encoders, decoders and classifiers form a distributionally robust game with competitive and collaborative agent coalitions. Finally, we plan to test our approach on real-world applications including conditional image synthesis and emotional speech conversion.

## 1 Introduction

In the last decade, machine learning has achieved great success with the rapid development of deep neural networks. Recent algorithms beat humans on ImageNet Challenge [1], the largest benchmark for image recognition. On the other side, people use computers to mimic the creative power of humans. A stream of papers [2, 3, 4, 5, 6] claimed their ability to generate lifelike pictures, text and music that can fool the human evaluators as well as machine classifiers. I identify three levels of artificial intelligence: memorization, recognition and creativity. Computers are approaching humans on the first two levels, and now going to the third.

Apart from object recognition and classification, people want to learn the mechanism of data generation and synthesize new samples with desired properties. Discriminative learning tries to infer knowledge from data, while generative learning attempts to learn the full distribution of data and generate new samples. My research focuses on the latter problem.

In probability theory, discriminative models make predictions by learning a conditional distribution  $p(y | x)$ , and generative models synthesize new data by drawing samples from  $p(x)$ . Both distributions are estimated from a limited set of observations, which could be noisy and incomplete. The former problem is easier since  $y$  is usually in the low dimension space. Sometimes people only use a small portion of the data and ignore the other parts. For example in support vector machines (SVM) only points near the decision surface have influence on the classifier. The latter problem is harder because generative models need to estimate the full data distribution  $p(x)$ . Modeling the high-dimensional random variable  $x$  is difficult. As the number of configurations can grow exponentially with the number of dimensions, there are not enough training examples for each dimension. Another

42 concern is the computation challenge; many algorithms involve operations that grow exponentially  
43 with the number of dimensions.

44 Generative models are widely used in unsupervised learning. A major limitation of the current  
45 learning algorithms is that they rely on large amounts of well-labelled data to achieve good accuracy.  
46 However, it is not available in many industrial applications. For this reason, people adapt a pretrained  
47 model on a source dataset to a similar target dataset. The model may perform poorly as it is specialized  
48 to the source domain. Domain adaptation strategies can solve this problem. Let  $x_1 \in X_1$  be source  
49 domain data with associated labels  $y_1$ ,  $x_2 \in X_2$  be target domain data with unknown labels  $f(x_2)$ . A  
50 mapping from the source domain to the target domain can be established by the conditional generative  
51 model  $p(x_2 | x_1)$ . The classifier  $f$  in the target domain is learnt by the category knowledge from the  
52 source domain.

53 Apart from traditional generative models like Gaussian mixture model (GMM) and Naive Bayes,  
54 neural network models exhibit increasing importance in modeling high dimension data. It is amazing  
55 if machines can generate artistic work in painting, music and sculpture. There are three kinds of  
56 dominant approaches in deep generative learning: generative adversarial networks (GAN), variational  
57 auto-encoders (VAE) and autoregressive models. In general, training deep generative models is hard  
58 and time consuming. The high dimensional training data and complex objective structures lead to  
59 many problems in optimization, such as instability, saturation, and mode collapse. Moreover, some  
60 models fail to provide enough diversity in the generated examples or just memorize the training set.

61 In this research, we plan to develop a new game-theoretic framework for generative learning. We  
62 first address the problem of data synthesis. The goal is to produce lifelike artificial examples that  
63 are indistinguishable from the training data. This involves the problem of computing the statistical  
64 distance between two datasets: the original real sample set and the generated fake sample set. Several  
65 distance metrics will be investigated to measure the similarity between distributions. As an important  
66 part of this research, we will compare Wasserstein distance [7] with the most prevalent information-  
67 based metrics such as Kullback-Leibler [8] and Jensen-Shannon divergence [9] in generative learning.  
68 We plan to develop a practical method to approximately calculate the distance between distributions.  
69 Both theoretical analysis and experimental simulations will be provided.

70 We then propose a conditional generative model to solve the unsupervised style transfer problem in  
71 image and speech processing. The objective is to learn a translation model between domains that can  
72 modify the domain-specific style features and preserve the domain-invariant content information. The  
73 disentangled representations of style and content are modeled by autoencoders and domain classifiers.

74 We formulate the problem as a distributionally robust game with payoff uncertainty. The encoders,  
75 decoders and classifiers form competitive and collaborative agent coalitions in this game. They  
76 optimize towards their self objectives as well as the common interest with other players. Agents  
77 with similar objectives form a group and work together against the others in order to optimize their  
78 expected payoffs. The optimization process is not deterministic since the payoff function of each  
79 player depends on the actions of other players. The distributionally robust Nash equilibrium is  
80 achieved by solving a minimax optimization problem, in which each agent tries to maximize its  
81 worst-case payoff. An iterative learning algorithm is introduced to solve the robust Nash equilibrium.

82 We will first verify the theoretical results through simulations, and then apply our approach on  
83 real datasets of images and speech. The agents are implemented with deep convolution neural  
84 networks to capture the spatial and temporal correlations in high dimension data. We plan to work on  
85 several tasks such as object detection, vehicle tracking, image synthesis and voice conversion. This  
86 research contributes to the areas of distributionally robust game, deep generative learning, stochastic  
87 optimization and time-series data analysis.

88 The following chapters first describes the already completed work: the theory part of distributionally  
89 robust games (Chapter 2), the investigation of distance metrics (Chapter 3), learning algorithms  
90 and applications on image synthesis (Chapter 4). After that the progress on the current work of  
91 style transfer is described in Chapter 5. Finally, the thesis outline and research timeline are given in  
92 Chapter 6.

## 2 Distributionally Robust Games

One fundamental problem in generative learning is to estimate the full data distribution given finite noisy observations. Suppose the training samples  $x$  are drawn from a distribution  $m$ , and the generative model produces new data  $\tilde{x}$  by sampling from an estimated distribution  $\tilde{m}$ . The objective is to minimize the statistical distance between  $m$  and  $\tilde{m}$  so that the synthesized data is similar to the real ones. Since the comparison is made on two statistical objects instead of two individual sample points, [there is no deterministic objective to use](#). We formulate the problem as a distributionally robust game (DRG) with payoff uncertainty, and offer a corresponding distributionally robust equilibrium concept.

In this chapter, we first introduce the definition of distribution uncertainty set and formulate the distributionally robust game as a minimax optimization problem. Then, we define the distributionally robust equilibrium, analysis the complexity of the problem, and use triality theory to reduce model dimensionality via Lagrangian relaxation. Finally, we apply this model to solve an unsupervised generative learning problem.

### 2.1 Introduction

Games with payoff uncertainty refer to games where the outcome of a play is not known with certainty by the players. Such games are also called incomplete information games and can be formalized in different ways. Distribution-free models of incomplete information games, both with and without private information are examined in [10, 11]. There the players use a robust optimization approach to contend with the worst-case scenario payoff. The distribution free approach relaxes the well-known Bayesian game model approach by Harsanyi. The limitation of the distribution-free model is that the uncertainty set has to be carefully designed and in most cases such approach leads to too conservative and unrealistic scenarios. [We propose a distribution-dependent model in which the uncertainty set contains all possible distributions within a bounded region](#).

Distributionally robust game (DRG) is a game with incomplete information. Instead of assuming an underlying mean-field or exactly known probability distribution, one acts with an uncertainty set, which could be distributions chosen by other players. The set of distributions should be chosen to fit for the applications at hand. In robust best-response problems, the uncertain sets are represented by deterministic models. The opponent players have a bounded capability to change the uncertain parameters, and therefore affects the objective function the decision maker seeks to optimize. Each player has his own robust best response optimization problem to solve. Thus, the standard best response problem of agent  $j$ :  $\inf_{a_j \in \mathcal{A}_j} l_j(a_j, a_{-j}, \omega)$  becomes the minimax robust best response:

$$\inf_{a_j \in \mathcal{A}_j} \sup_{\omega \in \Omega} l_j(a_j, a_{-j}, \omega) \quad (1)$$

where  $a_j$  is the action to choose and  $l$  is the objective functional evaluated at uncertain state  $\omega$ . This approach to uncertainty has a long history in optimization, control and games [12, 13, 14]. A credible alternative to this set-based uncertainty is to use a stochastic model, in which the uncertain state  $\omega$  is a random variable with distribution  $m$ . If we assume the generating mean-field distribution  $m$  is known, it becomes the standard stochastic optimal control paradigm. If  $m$  is not known and the only known is a set of distributions lie in the neighborhood of  $m$ :  $\tilde{m} \in B_\rho(m)$ , the resulting best response to mean-field formulation is the so-called distributionally robust best response:

$$\inf_{a_j \in \mathcal{A}_j} \sup_{\tilde{m} \in B_\rho(m)} \mathbb{E}_{\omega \sim \tilde{m}} l_j(a_j, a_{-j}, \omega) \quad (2)$$

#### 2.1.1 Distribution Uncertainty Set

Let  $(\Omega, \mathcal{F}, m)$  be a probability space. Here  $m$  is a probability measure defined on  $(\Omega, \mathcal{F})$ . The distribution  $m$  of the state  $\omega$  is used to capture the probability of the different scenarios and of the corresponding performance function obtained under each of such scenarios for fixed action profile. We assume that the exact distribution of the state is not available in general. Therefore we propose an uncertainty/constraint set among all possible distributions with a divergence bounded by above by a scalar  $\rho$ . When  $f$ -divergence [15, 16, 17] is applied, such a constraint takes the form

$$B_\rho(m) = \{\tilde{m} \mid \int_{\Omega} d\tilde{m} = \tilde{m}(\Omega) = 1, D_f(\tilde{m} \parallel m) \leq \rho\}, \quad (3)$$

When the Wasserstein distance [7] is applied, the uncertainty set contains probability distributions within a Wasserstein ball.

$$B_\rho(m) = \{\tilde{m} \mid \int_{\Omega} d\tilde{m} = \tilde{m}(\Omega) = 1, W(m, \tilde{m}) \leq \rho\}, \quad (4)$$

**$f$ -Divergence** Let  $m$  and  $\tilde{m}$  be two probability measures over a space  $\Omega$  such that  $\tilde{m}$  is absolutely continuous with respect to  $m$ . Then, for a convex function  $f$ , the  $f$ -divergence of  $\tilde{m}$  from  $m$  is defined as

$$D_f(\tilde{m} \parallel m) \equiv \int_{\Omega} f\left(\frac{d\tilde{m}}{dm}\right) dm - f(1),$$

where  $\frac{d\tilde{m}}{dm}$  is the Radon-Nikodym derivative of the measure  $\tilde{m}$  with respect to the measure  $m$ .

Recall that for a proper convex function  $f$ , the Legendre-Fenchel duality holds:  $(f^*)^* = f$  where

$$f^*(\xi) = \sup_x [\langle x, \xi \rangle - f(x)] = -\inf_x [f(x) - \langle x, \xi \rangle]. \quad (5)$$

## 2.1.2 Related Work

Strategic learning has proven to be a powerful approach in stochastic games. In particular, its algorithmic nature is well suited to accommodate parallel and distributed information exchange and processing as well as hardware realizability. However, almost all existing learning approaches work well only for specific classes of games such as concave-convex zero-sum games, convex potential games, and some S-modular game problems with unimodal objective functions. For more general classes of games, convergence of strategic learning dynamics is still an open issue. In addition, learning algorithms for finding fixed points or equilibria for general classes of games still present several challenges. For polymatrix games with finite action spaces, there have been great progress including Cournot adjustment, Brown-von Neumann-Nash dynamics, reinforcement learning, combined learning (see [18] and the references therein). For continuous action spaces, however, only few handful works are available. Evolutionary dynamics and revision protocols based actions has been proposed [19, 20, 21]. The exploration of continuous action takes too much time if the dynamics is based on individual action or measurable subsets [22, 23, 24]. Moreover, the convergence time of these existing strategic learning algorithms (when a convergence to a point or a limit cycle occurs) are unacceptably high even for potential games and it often requires strong assumptions such as bounded densities. The above mentioned prior works do not consider robust games setting. In [10, 11] a robust game framework is presented. The authors defined a distributed-free approach by considering the worst-case performance. However the choice of the uncertainty set remains an important part of the robust game modelling [25]. In this chapter we are interested in learning in distributionally robust games under  $f$ -divergence [15, 16, 17, 26].

## 2.2 Problem Formulation

In distributionally robust games, each agent  $j$  adjusts his action  $a_j \in \mathcal{A}_j$  to optimize the worst-case payoff functional  $\mathbb{E}_{\tilde{m}} l_j(a, \omega)$  subject to the constraint that the divergence  $D_f(\tilde{m} \parallel m) \leq \rho$ . This means that the worst-case performance is obtained under the assumption that a virtual player (nature) acts as a defender who modifies the distribution  $m$  into  $\tilde{m}$  with an effort capacity that should not exceed  $\rho > 0$ .

Throughout the paper we assume that the following conditions hold. The measure  $\tilde{m}$  is continuous with respect to  $m$ ; it is not a given profile, and it could be deformed or falsified by the defender. The payoff function  $l_j(\cdot, \omega)$  is proper and upper semi-continuous for  $m$ -almost all  $\omega \in \Omega$ . Either the domain  $\mathcal{A}_j$  is a non-empty compact set or  $\mathbb{E}_{\omega \sim \tilde{m}} l_j(a, \omega)$  is coercive.

**Definition 1** (Distributionally Robust Game). *The robust game  $\mathcal{G}(m)$  involves*

- The set of agents  $\mathcal{J} = \{1, 2, \dots, n\}$ ,  $n \geq 2$
- The action profile of each agent  $j$ ,  $\mathcal{A}_j$ ,  $j \in \mathcal{J}$
- The uncertainty set of probability distributions defined on  $m$ ,  $\Omega$  and  $\rho > 0$ :  $B_\rho(m)$
- The expected payoff function of agent  $j$ :  $\mathbb{E}_{\omega \sim \tilde{m}} l_j(a, \omega)$ , where  $\tilde{m}$  is an alternative distribution of  $m$  within some bounded distance

Then the robust stochastic optimization of agent  $j$  given the uncertain set and the action of other agents is formulated as a minimax problem

$$(P_j) : \inf_{a_j \in \mathcal{A}_j} \sup_{\tilde{m} \in B_\rho(m)} \mathbb{E}_{\omega \sim \tilde{m}} l_j(a, \omega). \quad (6)$$

With the above game in mind, we introduce the following solution concept for game  $\mathcal{G}(m)$ .

**Definition 2** (Distributionally Robust Nash Equilibrium). *Let  $a_j^*$  be the configuration of agent  $j$  and  $a_{-j}^* := (a_k^*)_{k \neq j}$ . A strategy profile  $a^* = (a_1^*, \dots, a_n^*)$  satisfying*

$$\sup_{\tilde{m} \in B_\rho(m)} \mathbb{E}_{\tilde{m}} l_j(a^*, \omega) \leq \sup_{\tilde{m} \in B_\rho(m)} \mathbb{E}_{\tilde{m}} l_j(a_j, a_{-j}^*, \omega)$$

for every  $a_j \in \mathcal{A}_j$  and every agent  $j \in \mathcal{J}$ , is said to be the distributionally robust pure Nash equilibrium of game  $\mathcal{G}(m)$ .

**Existence of distributionally robust Nash equilibria** In other words, reaching the robust Nash equilibrium means all agents achieve the optimal payoff in their worst-case scenario. As in classical game theory, sufficient condition for existence of robust equilibrium can be obtained from the standard fixed-point theory: if  $\mathcal{A}_j$  are nonempty compact convex sets and  $l_j$  are continuous functions such that for any fixed  $a_{-j}$ , the function  $a_j \mapsto l_j(a, \omega)$  is quasi-convex for each  $j$ , then there exists at least one distributionally robust *pure* Nash equilibrium. This result can be easily extended to the coupled-action constraint case for generalized robust Nash equilibria.

## 2.3 Minimax Robust Game

### 2.3.1 Infinite Dimension

Assume that  $a \mapsto \mathbb{E}_{\tilde{m}} l_j(a, \omega)$  is continuous for  $m$ -almost all  $\omega$ . Then, the functional  $F_j : \tilde{m} \mapsto \inf_{a_j} \mathbb{E}_{\tilde{m}} l_j(a, \omega)$  is Gateaux differentiable with derivative

$$F_{j,m}(\hat{m}) = \inf_{a_j \in \mathcal{A}_j^*(m)} \mathbb{E}_{\hat{m}} l_j(a, \omega),$$

where  $\mathcal{A}_j^*(m) = \arg \min_{a_j} \mathbb{E}_m l_j(a, \omega)$  is the best-response under  $m$ . This derivative in the space of square integrable measurable functions under  $m$  is of infinite dimensions, which does not facilitate the computation of the robust optimal strategy  $a_j^*, \tilde{m}^*$ . Below we propose an equivalent problem that reduces considerably the curse of dimensionality of the problem.

### 2.3.2 Triality Theory

We here streamline the basic idea of triality theory. For this purpose, consider uncoupled domains  $\mathcal{A}_j$ ,  $j \in \mathcal{J}$ . For a general real-valued loss function  $l_2$  on  $\mathcal{A}_1 \times \mathcal{A}_2$ , one has

$$\sup_{a_2 \in \mathcal{A}_2} \inf_{a_1 \in \mathcal{A}_1} l_2(a_1, a_2) \leq \inf_{a_1 \in \mathcal{A}_1} \sup_{a_2 \in \mathcal{A}_2} l_2(a_1, a_2)$$

and the difference is called the well-known duality gap.

$$\min_{a_1 \in \mathcal{A}_1} \max_{a_2 \in \mathcal{A}_2} l_2(a_1, a_2) - \max_{a_2 \in \mathcal{A}_2} \min_{a_1 \in \mathcal{A}_1} l_2(a_1, a_2)$$

It is widely known in duality theory from Sion's minimax Theorem [13] (which is an extension of von Neumann minimax Theorem) there is an equality, for example for convex-concave function, and the value is achieved by a saddle point in the case of non-empty convex compact domain. Similarly, for a general loss function  $l_3 : (a_1, a_2, a_3) \mapsto l_3(a_1, a_2, a_3)$ , one has

$$\begin{aligned} \inf_{a_3 \in \mathcal{A}_3} \sup_{a_2 \in \mathcal{A}_2} \inf_{a_1 \in \mathcal{A}_1} l_3(a_1, a_2, a_3) &\leq \inf_{a_1 \in \mathcal{A}_1, a_3 \in \mathcal{A}_3} \sup_{a_2 \in \mathcal{A}_2} l_3(a_1, a_2, a_3), \\ \sup_{a_3 \in \mathcal{A}_3} \inf_{a_2 \in \mathcal{A}_2} \sup_{a_1 \in \mathcal{A}_1} l_3(a_1, a_2, a_3) &\geq \sup_{a_1 \in \mathcal{A}_1, a_3 \in \mathcal{A}_3} \inf_{a_2 \in \mathcal{A}_2} l_3(a_1, a_2, a_3). \end{aligned}$$

**Proposition 1** (Triality). *Let  $(a_1, a_2, a_3) \mapsto l_3(a_1, a_2, a_3) \in \mathbb{R}$  be a loss function  $l_3$  defined on  $\prod_{i=1}^3 \mathcal{A}_i$ . Then, the following inequalities hold:*

$$\begin{aligned} &\sup_{a_2 \in \mathcal{A}_2} \inf_{a_1 \in \mathcal{A}_1, a_3 \in \mathcal{A}_3} l_3(a_1, a_2, a_3) \\ &\leq \inf_{a_3 \in \mathcal{A}_3} \sup_{a_2 \in \mathcal{A}_2} \inf_{a_1 \in \mathcal{A}_1} l_3(a_1, a_2, a_3) \\ &\leq \inf_{a_1 \in \mathcal{A}_1, a_3 \in \mathcal{A}_3} \sup_{a_2 \in \mathcal{A}_2} l_3(a_1, a_2, a_3), \end{aligned} \quad (7)$$

209 and similarly

$$\begin{aligned}
& \sup_{a_1 \in \mathcal{A}_1, a_3 \in \mathcal{A}_3} \inf_{a_2 \in \mathcal{A}_2} l_3(a_1, a_2, a_3) \\
& \leq \sup_{a_3 \in \mathcal{A}_3} \inf_{a_2 \in \mathcal{A}_2} \sup_{a_1 \in \mathcal{A}_1} l_3(a_1, a_2, a_3) \\
& \leq \inf_{a_2 \in \mathcal{A}_2} \sup_{a_1 \in \mathcal{A}_1, a_3 \in \mathcal{A}_3} l_3(a_1, a_2, a_3).
\end{aligned} \tag{8}$$

210 We use the above inequalities in the Lagrangean relaxation of the minimax robust game.

### 211 2.3.3 Dimension Reduction

212 In order to reduce the curse of dimensionality of the problem, we use the triality theory and convert  
 213 the original infinite dimensional problem  $(P_j)$  into a finite dimensional problem.

214 **Proposition 2.** *The (distributionally) robust best-response problem of player  $j$  is equivalent to the*  
 215 *finite dimensional stochastic optimization problem when  $\mathcal{A}_j$  are of finite dimensions:*

$$(P_j^*) \begin{cases} \inf_{a_j \in \mathcal{A}_j, \lambda \geq 0, \mu \in \mathbb{R}} l_j^*(a, \lambda, \mu, m) \\ l_j^*(a, \lambda, \mu, m) = \lambda(\rho + f(1)) + \mu + \int \lambda f^*\left(\frac{l_j - \mu}{\lambda}\right) dm = \mathbb{E}_m h_i, \end{cases} \tag{9}$$

216 where  $h_i$  is the integrand cost  $\lambda(\rho + f(1)) + \mu + \lambda f^*\left(\frac{l_j - \mu}{\lambda}\right)$ .

217 We have converted the infinite dimensional problem  $(P_j)$  into a finite dimensional problem  $(P_j^*)$ .  
 218 The above calculations culminate in the following result:

219 **Proposition 3.** *If  $a, \lambda^*(a), \mu^*(a)$  is a solution of  $(P_j^*)$  then the optimal likelihood  $L^*$  is such that*  
 220  $\int_{\omega} L^* dm = 1$ ,  $f'(L^*) = \frac{l_j - \mu^*}{\lambda^*}$ . *This means that  $a_j$  and  $d\tilde{m}^* = L^* dm$  provide a solution of the*  
 221 *original problem  $(P_j)$ .*

## 222 2.4 Case Study: Unsupervised Generative Learning

223 In unsupervised learning, people want to learn from the distribution of training data without labeled  
 224 information. Denote  $m$  as the real data distribution, and  $\tilde{m}$  as the distribution of synthesized outputs.  
 225 Learning generative model means minimizing the statistical distance between these two distributions  
 226  $\inf_m D_f(m \parallel \tilde{m})$ . However,  $m$  is not a given profile, so it needs to be estimated from the training  
 227 samples. We use Legendre-Fenchel duality to transform it into a minimax optimization problem

$$\begin{aligned}
\inf_m D_f(m \parallel \tilde{m}) &= \inf_m \int_{\Omega} f^{**}\left(\frac{dm}{d\tilde{m}}\right) d\tilde{m} - f(1) \\
&= \inf_m \int_{\Omega} [\sup_{\xi} \langle \xi, \left(\frac{dm}{d\tilde{m}}\right) \rangle - f^*(\xi)] d\tilde{m} - f(1) \\
&= \inf_m \int_{\Omega} \sup_{\xi} \langle \xi, dm \rangle - f^*(\xi) d\tilde{m} - f(1) \\
&\geq \inf_m \sup_{\xi} \int_{\Omega} \mathbb{E}_m \xi(\omega) - \mathbb{E}_{\tilde{m}} f^*(\xi(\omega)) - f(1)
\end{aligned} \tag{10}$$

228 The optimal structure for  $\xi$  is  $\xi^*(\omega) = f'\left(\frac{dm}{d\tilde{m}}\right)(\omega)$ . We assume that the optimal  $\xi^*$  can be parametrized  
 229 by and represented as  $\xi^*(\omega) = h_f(V(a_1, \omega))$  for some functions  $h_f$  and  $V$ , where  $V(a_1, \omega) : A_1 \times \Omega \rightarrow \mathbb{R}$ , and  $h_f : \mathbb{R} \rightarrow \text{Dom}(f^*)$ ,  $\text{Dom}(f^*) = \{x \mid f^*(x) \in \mathbb{R}\}$ . The measure  $\tilde{m}$  is  
 230 re-parametrized with  $\tilde{m}_{a_2}$ . Then,  
 231

$$\inf_m D_f(m \parallel \tilde{m}) = \inf_{a_1} \sup_{a_2} \mathbb{E}_{m_{a_1}} h_f(V) - \mathbb{E}_{\tilde{m}_{a_2}} f^*(h_f(V)) - f(1) \tag{11}$$

232 Therefore, the optimization problem  $\inf_m D_f(m \parallel \tilde{m})$  is formulated as a zero-sum game with two  
 233 adversarial players. The defender acts as a discriminator to find the maximum gap between  $m$  and  
 234  $\tilde{m}$ , while the attacker acts as a generator that tries to synthesize indistinguishable samples to fool  
 235 the discriminator. Since the payoff function is not deterministic and dependent on the other side,  
 236 each agent optimizes its worst-case performance. It forms a distributionally robust game and the  
 237 corresponding optimization problem is as follow:

$$(P) \begin{cases} \inf_{a_1 \in A_1} \sup_{a_2 \in A_2} L(a_1, a_2), \\ L(a_1, a_2) = \mathbb{E}_{m_{a_1}} h_f(V(a_1, \omega)) + \mathbb{E}_{\tilde{m}_{a_2}} (-f^*)(h_f(V(a_1, \omega))) - f(1) \end{cases} \tag{12}$$

238 where  $A_1$  denotes the decision space of the attacker and  $A_2$  is the decision space of the defender.

**Example 1.** Let  $f(x) = x \log x - (1+x) \log(1+x)$ . Then the derivative of the function  $f$  is  $f'(x) = 1 + \log x - 1 - \log(1+x) = \log(1 - \frac{1}{1+x})$  and  $f''(x) = \frac{1}{x(1+x)} > 0$ . As  $f$  is a convex function, we can apply the Legendre transform

$$f^*(\xi) = \sup_x [x\xi - f(x)] = -\log(1 - e^\xi)$$

The parameterized objective function is now

$$L(a_1, a_2) = \mathbb{E}_{m_{a_1}} h_f(V) + \mathbb{E}_{\tilde{m}_{a_2}} \log(1 - e^{h_f(V)}) - f(1)$$

If  $h_f = \log$ , then

$$L(a_1, a_2) = \mathbb{E}_{m_{a_1}} \log(V) + \mathbb{E}_{\tilde{m}_{a_2}} \log(1 - V) - f(1).$$

In order to provide training data with enough diversity, a particular sampling process was considered in [27]. Real data is randomly sampled from the training set  $\omega \sim m$ ; fake data is first sampled from a noise distribution  $z$  then fed into the generative model  $G(z) \sim \tilde{m}$ . The discriminator takes in both real and synthetic examples, and outputs two probabilities  $V = D(\omega)$  and  $D(G(z))$  for  $a_1, a_2$  respectively. In this game, two opponent agents optimize on different objective functions. It is a zero-sum game with no cooperation. The attacker  $G(z)$  tries to fool the defender  $D$  by minimizing  $\log(1 - D(G(z)))$ , which means incorrectly predicting a synthetic example as real. The defender optimizes the discriminator by maximizing the sum  $\log(D(\omega)) + \log(1 - D(G(z)))$ , which means correctly predicting both the real and fake examples. The generative model is trained by solving the Nash equilibrium in a distributional robust game. The formulation is robust in the sense that the attacker tries to sneak into the real data distribution while facing the supervision of the defender.

Based on the ideas above, we formulate the generative learning problem into the DRG framework. As depicted in Figure 1, there are two groups of players in this game. The attackers train the generative model  $G_{\theta_a}(z)$  to produce fake samples  $\tilde{x}_i$  that are similar to the real ones, where  $z$  is a low dimension random variable feed to the generator and  $\theta_a$  is the model parameter. The defenders explore the neighborhood of  $m$  and slightly change the real data to produce hard samples  $x'_i$  which have the maximum distance from the model distribution. The attackers again refine the model to fit for those hard samples. The loss function is defined by the discrepancy  $D(\tilde{m}, m')$ , where  $m'$  is the hard sample distribution chosen by the defender. Since the real distribution  $m$  is unknown and we only have an observation dataset  $\{x_1, \dots, x_N\} \subset \mathbb{R}^d$ , the optimization is performed by iteratively updating the generative model as well as the hard sample distribution  $m'$ , which is an approximation of  $m$  within some bounded uncertain set  $B_\rho(m)$ .

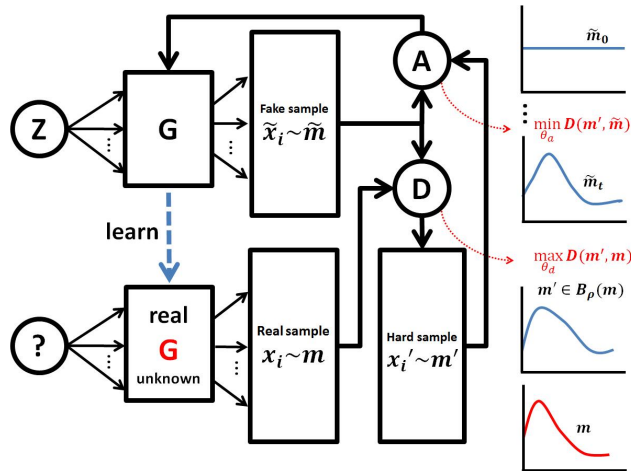


Figure 1: Distributionally robust game (DRG) framework for generative learning

As displayed in the right column of Figure 16, initially, the fake samples are drawn from an arbitrary distribution  $\tilde{m}_0$ , e.g., a uniform distribution. In each iteration, the attackers refine it closer to the hard sample distribution  $m'$ , while the defenders keep looking for the worst-case approximation  $m'$  within

269 the uncertain set  $B_\rho(m)$ . Once the algorithm converges, i.e.  $D(\tilde{m}, m') \leq \epsilon$ , we can ensure that the  
 270 discrepancy is bounded by a small value if  $D(\cdot, \cdot)$  satisfies triangle inequality:

$$D(\tilde{m}, m) \leq D(\tilde{m}, m') + D(m', m) \leq \epsilon + \rho \quad (13)$$

271 If  $m'$  is exactly the worst distribution in  $B_\rho(m)$ , we can say that  $D(\tilde{m}, m) \leq |\rho - \epsilon|$  (see Figure 2).  
 272 Therefore, the learning task is completed and the fake samples drawn from  $\tilde{m}$  will be indistinguishable  
 273 from the real ones.

274 Next chapter will discuss the properties of Wasserstein distance as a metric for  $D(\cdot, \cdot)$  and compare it  
 275 with the popular used KL divergence.



## 3 Wasserstein Metric

### 3.1 Introduction

The uncertain set involves a term to measure the similarity between two distributions. There are various ways to define the divergence between two sets. A straightforward way is to sum up the point-wise loss on those two sets, such as  $L^p$  distances used in ridge regression ( $L^2$ -norm) and Lasso ( $L^1$ -norm), KL divergence and its symmetric alternative Jensen-Shannon (JS) divergence. These loss functions are decomposable and widely adopted in both discriminative and generative tasks. Since the evaluation can be conducted on individual parts, they provide convenience for incremental learning and it's easier to develop efficient algorithms. However, they do not take into account the interactions of the individual points within a set. Non-decomposable losses such as F-measure, total variation and Wasserstein distance capture the entire structure of data and provide better topologies for optimization, at the cost of additional computational burden in loss evaluation.

#### 3.1.1 Optimal Transportation Problem

The optimal transport cost measures the least energy required to move all the mass in the initial distribution  $f_0$  to match the target distribution  $f_1$ .

$$C(f_0, f_1) = \inf_{\pi \in \Pi(f_0, f_1)} \int_{\mathcal{X} \times \mathcal{Y}} c(x, y) d\pi(x, y) \quad (14)$$

where  $c(x, y)$  is the ground cost for moving one unit of mass from  $x$  to  $y$ , and  $\Pi(f_0, f_1)$  denotes the set of all measures on  $\mathcal{X} \times \mathcal{Y}$  with marginal distributions  $f_0, f_1$ , i.e., the collection of all possible transport plans [7].

Wasserstein distance is a specific kind of optimal transport cost in which  $c(x, y)$  is a distance function. The  $p^{th}$  Wasserstein distance ( $p \geq 1$ ) is defined on a completely separable metric space  $(\mathcal{X}, d)$ :

$$W_p(f_0, f_1) := \left( \inf_{\pi \in \Pi(f_0, f_1)} \int_{\mathcal{X} \times \mathcal{X}} d(x, y)^p d\pi(x, y) \right)^{\frac{1}{p}} \quad (15)$$

Specifically, when  $p = 1$ ,  $W_1$  is called the Kantorovich-Rubinstein distance or Earth-Mover (EM) distance. It has duality as a supremum over all 1-Lipschitz functions  $\psi$ :

$$W_1(f_0, f_1) = \sup_{\|\psi\|_{Lip} \leq 1} \mathbb{E}_{f_0}[\psi(x)] - \mathbb{E}_{f_1}[\psi(x)] \quad (16)$$

Moreover, if  $f_0, f_1$  are densities defined on  $\mathbb{R}$ ,  $F_0, F_1$  denotes their cumulative distribution functions (CDF). In this one dimensional case  $f_0, f_1$  have inverse CDFs  $F_0^{-1}$  and  $F_1^{-1}$ . Then the Wasserstein distance can be computed by the integration of their CDF difference.

$$W_1(f_0, f_1) = \int_{\mathbb{R}} |F_0(x) - F_1(x)| dx \quad (17)$$

### 3.2 From KL divergence to Wasserstein Metric

Although KL divergence and its generalized version  $f$ -divergence are very popular in generative learning literatures, using Wasserstein metric  $D(m, \tilde{m}) = W_p(m, \tilde{m})$  has at least three advantages. First, Wasserstein metric is a true distance: the properties of positivity, symmetry and triangular inequality are fulfilled. Thanks to triangular inequality (Figure 2), the maximum discrepancy  $D(m, \tilde{m})$  is bounded by  $\rho + \epsilon$  when the optimizer approaches the distributionally robust equilibrium.

Second, the Wasserstein space has a finer topology where the loss changes smoothly with respect to the model parameters, thus effective gradients are always available during optimization. The information based divergence like KL doesn't recognize the spatial relationship between random variables.  $D_{KL}(m, \tilde{m}) = \int m(x) \log(\frac{m(x)}{\tilde{m}(x)}) dx$  is invariant of reversible transformations on  $x = (x_1, x_2, \dots)^T$  because  $m(x)dx$  removes the dimensional information. This property is illustrated in Figure 3. Therefore defining the uncertain set in Wasserstein space is more reasonable than using KL. It ensures the hard samples drawn from the Wasserstein neighborhood  $B_\rho(m)$  will not deviate too far from the real distribution.



### 3.4 Case Study: A Toy Example

This example shows learning the optimal transportation between two density distributions. The initial and target distributions are given by two Gaussians (Figure 5).

$$\begin{aligned} m_0 &= \mathcal{N}(\mu = (0.2, 0.3), \sigma = 0.1) \\ m_1 &= \mathcal{N}(\mu = (0.6, 0.7), \sigma = 0.07) \end{aligned} \tag{18}$$



Figure 5: Initial and target distribution  $m_0, m_1$ .

Initially, the intermediate states are set as linear interpolations, where a unimodal distribution is changed to bimodal. This is not desirable because it doesn't keep the mode. (think of an object being split into two parts while moving and then put together, as the lower part of Figure 4). It has been proved in [29] that optimizing with Wasserstein loss can preserve variance of the intermediate distributions, thus the object moves from one position to another without changing its shape. Figure 6 shows that after 1000 iterations, the Gaussian peak is preserved in the intermediate distributions, and the loss decreases to the optimal transport cost (Figure 7).

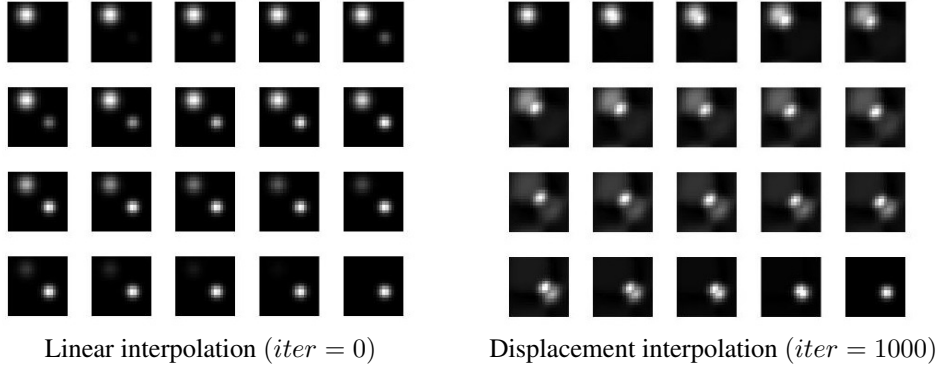


Figure 6: Learning the optimal transportation between two normal distributions.

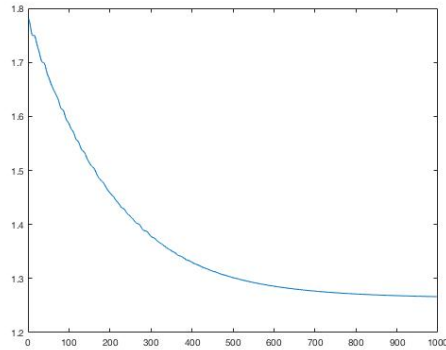


Figure 7: [Transportation cost during optimization.](#)

## 4 Learning Algorithms

In this section we provide learning procedure to solve the distributionally robust Nash equilibria, and develop practical implementation algorithms to train deep generative models.

### 4.1 Distributionally Robust Optimization

In DRG, the attacker and defender work against each other to find the robust Nash equilibrium by solving a minimax optimization problem

$$(P_j) : \inf_{a_j \in \mathcal{A}_j} \sup_{m' \in B_\rho(m)} \mathbb{E}_{m'} l_j(a, \omega')$$

Since  $B_\rho(m)$  is a subset of Lebesgue space (the set of integrable measurable functions under  $m$ ), the original problem  $(P_j)$  has infinite dimensions, which does not facilitate the computation of robust optimal strategies. It has been proved in section 2.3.3 that  $(P_j)$  can be reduced to a finite dimensional stochastic optimization problem when  $\omega' \mapsto l_j(a, \omega')$  is upper semi-continuous and  $(\Omega, d)$  is a Polish space. We introduce a Lagrangian function for constraint (4),

$$\tilde{l}_j(a, \lambda, m, m') = \int_{\omega'} l_j(a, \omega') dm' + \lambda(\rho - W(m, m')) \quad (19)$$

the original problem  $(P_j)$  becomes

$$(\tilde{P}_j) : \inf_{a_j \in \mathcal{A}_j, \lambda \geq 0} \sup_{m' \in B_\rho(m)} \tilde{l}_j(a, \lambda, m, m') \quad (20)$$

In robust game  $\mathcal{G}(m)$ , the defenders search for the worst hard sample distribution  $m'$  in the Wasserstein neighborhood of  $m$  to maximize its loss against the model  $\tilde{m}$ . According to the definition of Wasserstein metric with ground distance  $d(\cdot, \cdot)$ ,

$$\begin{aligned} \sup_{m'} \tilde{l}_j &= \lambda \rho + \sup_{m'} \int_{\omega'} [l_j(a, \omega')] dm' - \lambda W(m, m') \\ &= \lambda \rho + \int_{\omega'} \sup_{\omega'} [l_j(a, \omega') - \lambda d(\omega, \omega')] dm \end{aligned} \quad (21)$$

Define the integrand cost as

$$h_j(a, \lambda, \omega) = \lambda \rho + \sup_{\omega'} [l_j(a, \omega') - \lambda d(\omega, \omega')], \quad (22)$$

then  $(\tilde{P}_j)$  becomes a finite dimension problem on  $\mathcal{A}_j \times \mathbb{R}_+ \times \Omega$  if  $\mathcal{A}_j$  and  $\Omega$  have finite dimensions

$$(\tilde{P}_j^*) \inf_{a_j \in \mathcal{A}_j, \lambda \geq 0} \mathbb{E}_m h_j(a, \lambda, \omega) \quad (23)$$

Since  $m$  is an unknown distribution observed by the noisy unsupervised dataset  $x_1, \dots, x_N$ , it is challenging to compute the expected payoff  $\mathbb{E}_{m'} l_j(a, \omega')$ ,  $\mathbb{E}_m h_j(a, \lambda, \omega)$  and their partial derivatives. We need a stochastic learning algorithm to estimate the empirical gradients for the Wasserstein metric.

For a single player, the stochastic state  $\omega_j$  leads to error

$$\varepsilon_j = \nabla_{a, \lambda} h_j(a, \lambda, \omega_j) - \nabla_{a, \lambda} \mathbb{E}_m h_j(a, \lambda, \omega)$$

The variance of  $\varepsilon_j$  is high and not vanishing. To handle this, we introduce a swarm of players  $\omega_j \sim m$ ,  $j \in \mathcal{J}$ , then the error term becomes

$$\varepsilon = \frac{1}{|\mathcal{J}|} \sum_j \nabla_{a, \lambda} h_j(a, \lambda, \omega_j) - \nabla_{a, \lambda} \mathbb{E}_m h_j(a, \lambda, \omega)$$

It has zero mean and standard deviation as

$$\sqrt{\mathbb{E}[\varepsilon^2]} = \frac{1}{|\mathcal{J}|} \sqrt{\text{var}[\nabla_{a, \lambda} h_j(a, \lambda, \cdot)]}$$

For realized  $\omega \leftarrow \{x_1, \dots, x_N\}$ , the expected payoff for  $N$  players is  $\frac{1}{N} \sum_{j=1}^N h_j(a, \lambda, \omega_j)$ , and the optimal strategy is

$$(a^*, \lambda^*) \in \arg \min_{a, \lambda} \sum_{j=1}^N h_j(a, \lambda, \omega_j)$$

This provides an accurate robust equilibrium payoff when  $N$  is very large.

## 4.2 Numerical Investigation

To illustrate the stochastic learning algorithm we consider specific robust games with finite number of players. Each player acts as if he is facing a group of opponents whose randomized control actions are limited to a Wasserstein ball, and tries to optimize the worst case payoff. The random variable  $\omega$  is distributed over  $m$  and we assume it has finite  $p$  moments. We choose  $|\mathcal{J}| = 2$ ,  $p = 2$ ,  $d(\omega, \omega') = \|\omega - \omega'\|_2^2$  and a convex payoff function  $l_j(a, \omega')$  defined on  $\mathbb{R}^2 \times \mathbb{R}^2$

$$l_j(a, \omega') = \|\omega' - a\|_2^2 = (\omega'_1 - a_1)^2 + (\omega'_2 - a_2)^2 \quad (24)$$

The optimal defender state  $\omega'^*$  is computed through the Moreau-Yosida regularization, and the attacker's action pushes it closer to the destination  $\omega$  as shown in Figure 8.

$$\begin{aligned} \sup_{m'} \tilde{l}_j &= \lambda \rho + \int_{\omega \in \Omega} \phi_j(a, \lambda, \omega) dm \\ \phi_j(a, \lambda, \omega) &= \sup_{\omega' \in \mathbb{R}^2} [l_j(a, \omega') - \lambda d(\omega, \omega')] \\ &= \sup_{\omega' \in \mathbb{R}^2} (\|\omega' - a\|_2^2 - \lambda \|\omega' - \omega\|_2^2) \end{aligned} \quad (25)$$

$$\omega'^* = \omega + \frac{\omega - a}{\lambda - 1}, \quad (\lambda > 1) \quad (26)$$

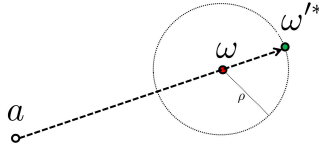


Figure 8: Action pushes the particle toward  $\omega$  (which is unknown) given  $\omega'^*$

Then  $d(\omega, \omega'^*) = \|\frac{\omega - a}{\lambda - 1}\|_2^2$ , leads to the worst-case loss

$$l_j(a, \omega'^*) = \|\omega'^* - a\|_2^2 = \frac{\lambda^2}{(\lambda - 1)^2} \|\omega - a\|_2^2 \quad (27)$$

The Moreau-Yosida regularization on  $m'$  realized at  $\omega'^*$  is

$$\begin{aligned} \phi_j(a, \lambda, \omega) &= l_j(a, \omega'^*) - \lambda d(\omega, \omega'^*) \\ &= \frac{\lambda}{\lambda - 1} \|\omega - a\|_2^2 \end{aligned} \quad (28)$$

The integrand cost function  $h_j = \lambda \rho^2 + \frac{\lambda}{\lambda - 1} \|\omega - a\|_2^2$ . Thus, problem  $(\tilde{P}_j^*)$  becomes

$$\inf_{a, \lambda} \mathbb{E}_m h_j = \inf_{a, \lambda} \int_{\omega} \lambda \rho^2 + \frac{\lambda}{\lambda - 1} \|\omega - a\|_2^2 dm, \quad (\lambda > 1) \quad (29)$$

Given  $N$  observations, the stochastic robust loss is

$$\begin{aligned} l_N^* &= \frac{1}{N} \sum_{j=1}^N h_j(a, \lambda, \omega_j) \\ &= \lambda \rho^2 + \frac{\lambda}{N(\lambda - 1)} \sum_{j=1}^N \|\omega_j - a\|_2^2 \end{aligned}$$

We set  $\rho = 1$  and  $m$  is a dirac distribution where  $\omega_j \equiv 1$ . Figure 9 plots the trajectories of strategies during learning.

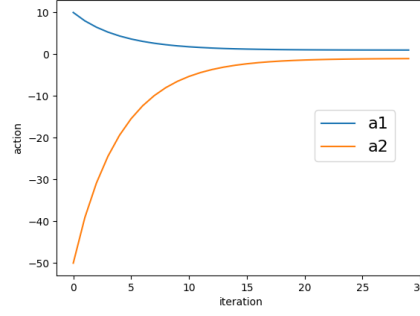


Figure 9: The optimal strategies converges to  $(a_1^*, a_2^*) = (1, -1)$

### 381 4.3 Train a Deep Generative Model

382 Image generative models such as VAE [30], GAN [31], WGAN [32] have shown great success in  
 383 recent years. VAE trains an encoder network and a decoder network by minimizing the reconstruction  
 384 loss, i.e., the negative log-likelihood with a regularizer. It tends to produce blurring images due to the  
 385 additional noise terms in their model. GAN trains a generator network and a discriminator network  
 386 by solving a minimax problem based on the KL-divergence. The model is unstable (Figure 13) due to  
 387 the discontinuity of the information-based loss functions, and the generator is vulnerable to saturation  
 388 as the discriminator getting better. [33, 34] gives some empirical solutions to these problems, e.g.,  
 389 keeping balance in training generator and discriminator networks, designing a customized network  
 390 structure. WGAN [32] defines a GAN model by an efficient approximation (equation 16) of the Earth  
 391 Mover distance. During training, it simply crops all weights of the discriminator network to maintain  
 392 the Lipschitz constraint.

393 In our framework, generative learning is formulated as a distributionally robust game with two  
 394 competitive groups of players, whose actions are defined on the parameter space  $\theta = (\theta_a, \theta_d) \in \Theta$ .  
 395 In stochastic settings,  $\omega, \omega'$  and  $\tilde{\omega}$  are instantiated to sample vectors  $\{x_1, x_2, \dots\}, \{x'_1, x'_2, \dots\}$  and  
 396  $\{\tilde{x}_1, \tilde{x}_2, \dots\}$ . The attacker produces indistinguishable artificial samples  $\tilde{x}_i = G_{\theta_a}(z)$  to minimize the  
 397 discrepancy  $\inf_{\theta} D(\tilde{m}, m')$ , where  $\tilde{x}_i \sim \tilde{m}$ . Meanwhile, the defender produce adversarial samples  
 398  $x'_i = G_{\theta_d}(x_i)$ , which are substitutes of the real ones, to maximize the loss  $\sup_{m' \in B_\rho(m)} D(\tilde{m}, m')$ ,  
 399  $x_i \sim m$ , where  $x'_i \sim m'$ .

400 With Moreau-Yosida regularization, the defenders work on the following maximization problem to  
 401 generate the optimal adversarial samples in Wasserstein ball  $B_\rho(m)$ ,

$$\theta_d^* \in \arg \max_{\theta_d} l(\theta_a, \omega') - \lambda d(\omega, \omega')$$

402 and the attackers work on the minimization problem to find the best generative parameters  $\theta_a^*$

$$\theta_a^* \in \arg \min_{\theta_a, \lambda} \lambda \rho + l(\tilde{x}, x'^*) - \lambda d(x, x'^*)$$

403 Given enough observations  $\{x_1, x_2, \dots\}$  from the unknown real distribution  $m$ , a similar distribution  
 404  $\tilde{m}$  can be learned by solving the distributionally robust Nash equilibrium. New samples generated  
 405 from  $\tilde{x}_i \sim \tilde{m}$  should be indistinguishable from the real ones. The DRG algorithm is summarized in  
 406 Algorithm 1.

### 407 4.4 Case Study: Unsupervised Learning for Clustering

408 We apply our algorithm to learn a three-dimensional data distribution for the Fisher's Iris dataset [35].  
 409 This dataset contains 150 samples, each having five attributes: petal length, petal width, sepal length,  
 410 sepal width and class. We remove the class label and using PCA to extract three most prominent  
 411 features. The observed data samples are plotted in Figure 10, each color represents a class of flower.  
 412 The generative model is set as the affine transformations of three unit balls  $z_1, z_2, z_3$ . Initially, the  
 413 attacker parameters  $W_1, W_2, W_3$  are set as identity matrices and  $b_1 = b_2 = b_3 = 0$ , so all fake

---

**Algorithm 1** DRG with Wasserstein metric

---

**Input:** real data  $(x_i)_{i=1}^N$ , batch size  $n$ , initial attacker parameters  $\theta_{a0}$ , Lagrangian multiplier  $\lambda_0$ , initial defender parameters  $\theta_{d0}$ , number of defender updates per attacker loop  $n_d$ , Wasserstein ball radius  $\rho$ , learning rate  $\eta_a, \eta_d$ , low-dimension random noise  $z \sim \zeta$

**Output:**  $\theta_a, \theta_d, \lambda$

**while**  $\theta_a$  has not converged **do**

**for**  $t = 1, 2, \dots, n_d$  **do**

    Sample  $(x_i)_{i=1}^n \sim m$  from real dataset

    Sample  $(\tilde{x}_i)_{i=1}^n \sim \tilde{m}$  from generator  $G_{\theta_a}(z)$

$y_i \leftarrow E_{\theta_d}(x_i), \quad \tilde{y}_i \leftarrow E_{\theta_d}(\tilde{x}_i)$

    Modify to adversarial samples  $y'_i \leftarrow G_{\theta_d}(y_i)$

$g_d \leftarrow \nabla_{\theta_d} l(\tilde{y}_1^n, y'_1^n) - \lambda d(y_1^n, y'_1^n)$

$\theta_d \leftarrow \theta_d + \eta_d \text{RMSPProp}(g_d)$

**end for**

  Sample  $(x_i)_{i=1}^n \sim m$  from real dataset

  Sample  $(\tilde{x}_i)_{i=1}^n \sim \tilde{m}$  from generator  $G_{\theta_a}(z)$

$y_i \leftarrow E_{\theta_d}(x_i), \quad \tilde{y}_i \leftarrow E_{\theta_d}(\tilde{x}_i)$

  Modify to adversarial samples  $y'_i \leftarrow G_{\theta_d}(y_i)$

$g_{a,\lambda} \leftarrow \nabla_{\theta_a, \lambda} \lambda \rho + l(\tilde{y}_1^n, y'_1^n) - \lambda d(y_1^n, y'_1^n)$

$\theta_a \leftarrow \theta_a - \eta_a \text{RMSPProp}(g_{a,\lambda})$

$\lambda \leftarrow \lambda - \eta_a \text{RMSPProp}(g_{a,\lambda})$

**end while**

---

414 samples are located in a unit ball shown in Figure 10.

$$\begin{aligned}\tilde{x} &= G_{\theta_a}(z) \in \{W_1 z_1 + b_1, W_2 z_2 + b_2, W_3 z_3 + b_3\} \\ x' &= G_{\theta_d}(x) = W_0 x + b_0\end{aligned}$$

415 Hard samples are generated from the real ones by another affine transformation with parameters  
416  $W_0, b_0$ . Since the dimension is low, we directly work with the raw data vectors. We set  $n_d = 20, \rho =$   
417  $0.1, \lambda_0 = 10, \eta_a = 0.1, \eta_d = 0.01$ , and the training cost for attackers and defenders are displayed in  
418 Figure 11. In each defender loop, the hard samples are refined to maximize the Wasserstein loss, and  
419 then the attackers minimize the worst-case cost. After 150 iterations, the algorithm converged at the  
420 distributionally robust Nash equilibrium. The generated samples (black dots in Figure 10) successfully  
421 covered the region of real samples, which demonstrated the effectiveness of our algorithm.

422 Next section will show the application of DRG on image synthesis, in which the attackers and  
423 defenders are realized with deep neural networks.

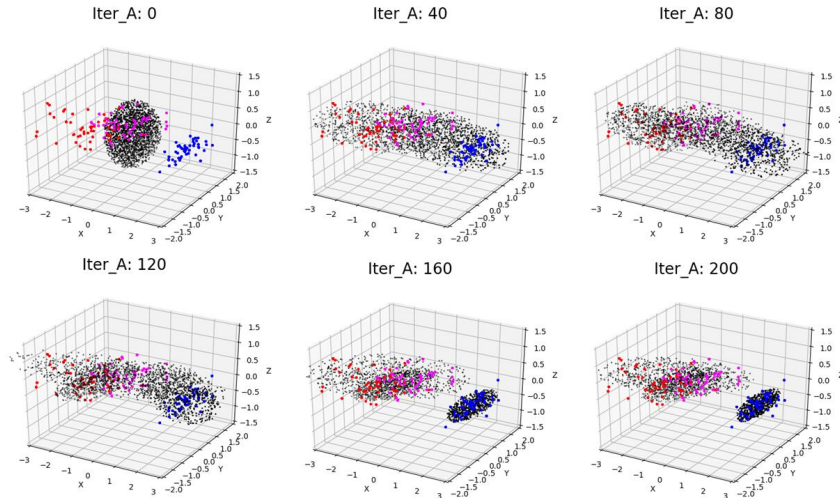


Figure 10: Generative learning for Iris Flower dataset

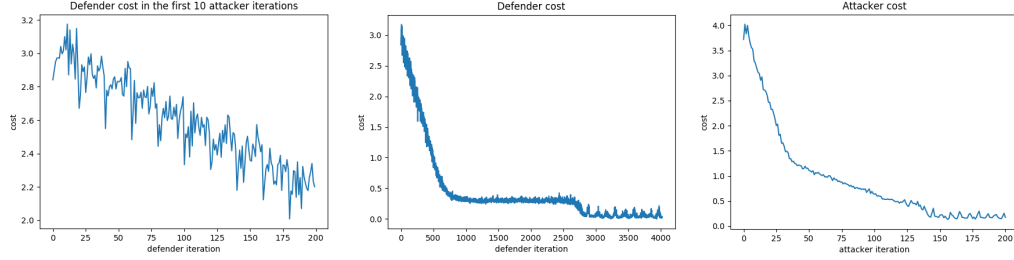


Figure 11: Attacker and defender cost during training

#### 4.5 Case Study: Generative Modeling for Image Synthesis

**Dataset** We apply our DRG algorithm on the CelebA [36] dataset to generate artificial human faces. The training set has 202K cropped face images with size  $64 \times 64$ , therefore each real sample  $x_i \sim m$  has 12288 dimensions. The artificial samples are generated from low-dimensional noise vectors  $z \sim \zeta$ , where  $\zeta$  is a random normal distribution.

**Network Structure** In this paper, the generative network  $x = G_{\theta_a}(z)$  follows the DCGAN [37] architecture. We design  $y' = G_{\theta_d}(y)$  as a single layer neural network to perform modification. For the encoding network  $y = E_{\theta_d}(x)$ , we use one CNN-ReLU layer followed by 3 CNN-BatchNorm layers and a fully connected layer to produce code vectors. Both networks have about 5 million training parameters.

**Loss Functions** In DRG algorithm, the Wasserstein distance  $l(\tilde{x}^n, x'^n)$  is implemented by Sinkhorn-Knopp’s algorithm [38], and the ground cost  $d(x^n, x'^n) = \frac{1}{n} \sum_{i=1}^n \|x_i - x'_i\|_2^2$ . Instead of directly computing the L2-norm on raw data vectors, Algorithm 1 uses an encoder network  $y = E_{\theta_d}(x)$  to learn useful features.

**Hyperparameters** The encoder maps the original data into a 100-dimension feature space, which matches the dimension of the random noise  $z$ . In all experiments, the cost based on Wasserstein metric is normalized to  $[0, 1]$ , where the supremum indicates the cost between images that are all black and all white. The hyperparameters listed in Algorithm 1 are chosen by validation and listed in table 1; others are set as the default values in their references. For training we choose the RMSProp optimizer [39] because it doesn’t involve a momentum term. Empirically, we found momentum-related optimizers may deteriorate the training. The reason is, in robust games the payoff function is dynamic and changes every time the other players take actions. Since the structure of the objective surface is not stationary, it’s meaningless to follow the velocity of the previous optimization steps.

Table 1: Hyper parameters

parameters	$n$	$\rho$	$\lambda_0$	$n_d$	$\eta_a, \eta_d$	$\theta_{a0}, \theta_{d0}, \eta$
values	64	0.1	10	1	0.00005	random normal

**Evaluation** Experimental results are demonstrated in Figure 12, in which the last line shows the most similar samples in the real dataset. The training curve for DRG is plotted in Figure 14. It means the Wasserstein loss is highly related to the sample quality. By optimizing the worst-case loss function, the DRG model converges very quickly to the real data distribution and successfully produce sharp and meaningful images. In experiments we found that the original GAN generator [37] suffers from unpredictable quality deterioration at iteration  $5.3K$ ,  $7.8K$ ,  $10.2K$  (Figure 13), etc, while our algorithm keeps improving the sample quality. This problem is caused by the discontinuity of the KL-divergence.

The evaluation of generative models is itself a research topic. [40] figured out that different evaluation metrics favor different models. For example, a high log-likelihood doesn’t mean good visual quality, and vice versa. Therefore, the metric used in training and evaluation should fit for the specific



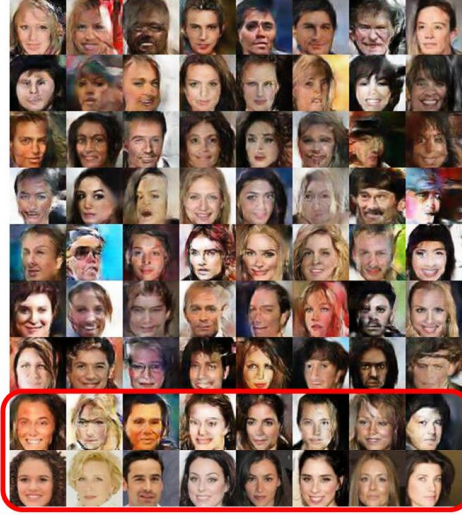


Figure 12: DRG results on CelebA, attacker iteration = 300K

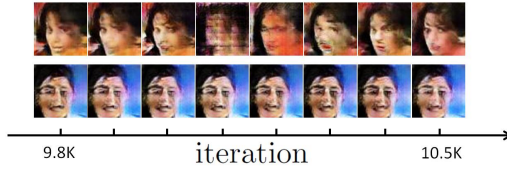


Figure 13: Stability of the generated models. Upper: DCGAN, Bottom: DRG

458 application. In our case, the learned fake data distribution should be as close as to the real one. So  
 459 we measure the discrepancy between these distributions using Wasserstein metric. We compare our  
 460 algorithm with DCGAN [37] and WGAN [32], and report the quantitative results in Table 2.

461 The computation complexity per attacker iteration is linear  $O(n)$  with respect to the batch size. We  
 462 use a Titan Xp to train the model and plot the computation time in Figure 15. When  $n = 64$ , it takes  
 463 0.2 seconds for an attacker update. Our algorithm has smaller constant factor than WGAN.

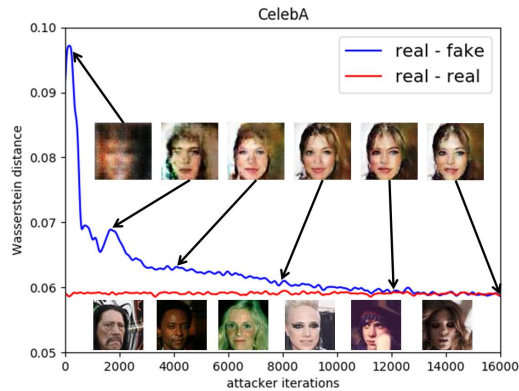


Figure 14: Training curve for DRG algorithm with Wasserstein metric. The loss goes down as generated samples getting better, and converges to the Wasserstein distance between two real data sets. The curves are smoothed for visualization purpose.

Table 2: Performance evaluation

$W(m, \tilde{m}) (\times 10^{-5})$	1K samples	10K samples
real - real	12.9	1.74
real - DRG	<b>22.6</b>	<b>15.9</b>
real - DCGAN	37.3	16.4
real - WGAN	31.0	17.2

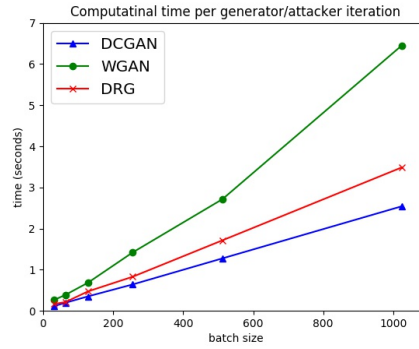


Figure 15: Computation time with respect to batch size.

## 5 Style Transfer as a Minimax Game

### 5.1 Introduction

Style transfer originally means rendering an image in different styles. This meaning can be extended to other kinds of data like music and speech. More generally, it refers to mapping data from one domain to another while keeping its semantic (underlying) content / the domain-invariant knowledge. For example, transfer photographs to artistic paintings, convert one person’s voice to another, or translate music to imitate different instruments. Another case is transfer learning. It adapts a pre-trained model in source domain to classify samples in target domain where labeled data is limited.

Let  $x_1 \in \mathbb{X}_1$  be samples in the source domain and  $x_2 \in \mathbb{X}_2$  be samples in the target domain. The goal of style transfer is to learn a mapping function  $\mathcal{T} : \mathbb{X}_1 \rightarrow \mathbb{X}_2$  such that the generated output  $x'_{2 \leftarrow 1} = \mathcal{T}(x_1)$  is indistinguishable from the real samples drawn from the target domain. The optimal mapping  $\mathcal{T}^*$  transforms  $x_1$  to  $x'_{2 \leftarrow 1}$  such that  $x'_{2 \leftarrow 1} \stackrel{d}{=} x_2$ . Semantic content should be preserved during the transformation.

We propose a game-theoretic approach to learn the mapping. The domain-invariant content information and domain-specific style information are decomposed by disentangled representation learning. For high dimensional data like image and speech waveform, we employ autoencoders to independently model the high-level semantic content and the low-level style information. The learning problem is formulated as a distributionally robust game with cooperative agents and payoff uncertainty. In this game, several groups of players run with different objectives. The intergroup competition and intragroup collaboration enable the players to learn from each other and optimize their worst-case performance.

This work has a wide range of applications. In visual and performing arts, it’s inspiring to automatically generate artificial paintings with user-specified style or play synthetic music with desired timbre and musical instrument. In informatics, it’s useful to transform speaker identity by modifying his voice to sound like another person. It is also possible to learn and mimic animal’s vocalization and study the feedback on the artificially generated sound. For case study, we apply our approach in two scenarios: image style transfer and emotional voice conversion.

### 5.2 Motivation

In machine learning, discriminative models predict labels from data by learning a conditional distribution  $p(y | x)$ , while generative models synthesize new data with desired labels by drawing samples from estimated distribution  $p(x | y)$ . From a perspective of probabilistic modelling, style transfer learns two conditional distributions  $p(x_1 | x_2)$  and  $p(x_2 | x_1)$ . When paired data is available, it is easy to infer from the joint distribution  $p(x_1, x_2)$ . For nonparallel data, the problem is ill-posed because the joint solution is not unique given two marginal distributions  $p(x_1), p(x_2)$ .

To solve this problem, additional constraints are required. Some researchers proposed to keep a particular part of the data unchanged, e.g., pixel intensity, gradient or object boundaries [41, 42]; others suggested to preserve some properties of the data, such as semantic features or class labels [43].

Zhu et al. [44] proposed a very straightforward constraint called cycle-consistency. It assumes that if a sample is translated from source domain to target domain and then translated back, it should be unchanged. Choi et al. [45] generalized it to perform multiple-domain translation using a single generative model. However, domain transfer is not a one-to-one mapping, but many-to-many. In some cases, the cycle-consistency constraint is too strong to provide enough diversity in the translated outputs.

Based on a similar idea, Liu et al. [46] developed the UNIT framework by making a fully shared latent space assumption, in which corresponding images across domains can be mapped to a same latent code in shared-latent space. This assumption implies the cycle-consistency constraint. Xun et al. [47] extended it to a partially shared latent space assumption, where each example is generated from a shared content code and a domain-specific style code. Images are translated across domains by replacing the style code.

Some approaches [45, 48, 49] assume there exists a transformation  $\mathcal{T}$  such that the source and target data can be matched in a new representation  $p(\mathcal{T}(x_1)) = p(\mathcal{T}(x_2))$ . The types of transformation includes projections, affine transform, and non-linear mapping defined by neural networks. The objective is to minimize the gap between two transformed distributions. Several metrics used to compare distributions have been discussed in chapter 3. One requirement is that  $p(x_1)$  and  $p(x_2)$  share a common support, otherwise the optimization process is instable due to vanishing gradient issues.

Another idea is to directly assume the mapping  $\mathcal{T} : x_2 = \mathcal{T}(x_1)$  to be the optimal transport between  $p(x_1)$  and  $p(x_2)$ . The optimal solution  $\mathcal{T}^*$  minimizes the global transportation cost (Wasserstein distance) between the source and target distributions. As a by-product, minimizing the transportation cost gives the alignment of corresponding samples as well as their labels. It can be used to solve the domain adaptation problem [50] by finding the joint distribution optimal transport between  $p(x_1, y_1)$  and  $p(x_2, f(x_2))$ .  $f$  is the classifier in the target domain, which can be inferred from the aligned data-label pairs.

While many approaches have been proposed, learning the mapping between different domains is still an open problem. Our work is motivated to support the development of more robust and comprehensible generative models for style transfer.

### 5.3 Related Work

There are several topics closely related to this work, but with different definitions.

**Deep Generative Modeling** The original objective of this topic is to generate new samples from scratch by learning complicated data distributions in an unsupervised way. At test time, it takes random noise as input and outputs realistic samples. In some cases, it can also take in conditional information to produce user-specified output. There are three main frameworks of deep generative modeling: generative adversarial networks (GANs) [31], variational auto-encoders (VAEs) [30], and auto-regressive models [51].

GANs build the generative model on the top of a discriminative network to force the output to be indistinguishable from the real samples. This model works pretty well for generating images with impressive visual quality [37] and high resolution [2]. Variations under this framework include conditional GAN [52] that generate samples conditioned on class labels, LAPGAN [53] that generates images in a coarse-to-fine fashion, WGAN-GP [54] that enables stable training of GANs without hyperparameter tuning.

VAEs use an encoder-decoder framework to model data in a latent space and optimize the reconstruction loss plus a regularizer. The generative process has two steps of sampling: first draw latent variables from  $p(z)$  and then draw datapoints from the conditional distribution  $p(x | z)$ . At test time, the encoder part is discarded and the decoder takes random noise as input to generate new samples. However, the reconstructed samples are blurry. This is because the VAE decoder assumes  $p(x | z)$  to be an isotropic Gaussian, which leads to the use of L2 loss. To remedy this, VAE-GAN [55] suggests learning the loss through a GAN discriminator.

Auto-regressive model is quite different from the above two. It aims at modeling time-varying processes by assuming that the value of a time series depends on its previous values and a stochastic term. For a sequential data sample  $x = (x_1, x_2, \dots, x_T)$ , the joint distribution  $p(x)$  is factorised as a product of conditional distributions

$$p(x) = \prod_{t=1}^T p(x_t | x_1, \dots, x_{t-1}) \quad (30)$$

This idea is quite straightforward for modeling audio sequence [56], but it also works for images. In PixelRNN [51], each image is written as a sequence, in which pixels are taken row by row from the image. The two-dimensional spatial autocorrelation of pixels is modeled by one-dimensional temporal correlations. Since the generation process is sequential, it requires a lot of GPU memory and computation time (200K updates over 32 GPUs) even after some modifications [57].

**Image Style Transfer** There are two types of style transfer problems: example-based style transfer where the style comes from one image, and domain-based style transfer where the style is learnt

from a collection of images in a specific domain. The former problem originates from nonphoto-realistic rendering (NPR) [58] in computer graphic, and has the similar meaning of realistic image manipulation. The goal is to edit image in a user-specified way and keep it as realistic as possible. Practical issues include texture synthesis and transfer [59], photo manipulation of shape and color [60], photorealistic image stylization [61], etc. In general, the output should be similar to the input in high-level structures and varies in low-level details such as color and texture.

Recently, Gatys et al. [62] claimed the image content and style information are separable in Convolutional Neural Network representations. They introduced a method [29] to separate and recombine content and style of natural images by matching feature correlations (Gram matrix) in different convolutional layers. However, their synthesis process is slow (an hour for a 512\*512 image). Moreover, the style from a single image is ambiguous and may not capture the general theme of an entire domain of images.

The second problem, also named as image-to-image translation, learns a mapping to transfer images from one domain to another. For example, super-resolution [63] maps low-dimentional images to high-dimention, colorization [64] maps gray images to color; other cases include day to night, dog to cat, yong to old, summer to winter, photographs to paintings, aerial photos to maps [43, 61, 65, 66, 67, 68, 69]. The mapping can be learnt in a supervised or unsupervised manner. In supervised settings [70, 71, 72], corresponding image pairs across domains are available for training. In unsupervised settings [2,4,8,9,10], there's no paired data and the training set only contains independent set of images for each domain. Our work is under the unsupervised settings because it is more applicable, and the training data is almost free and unlimited.

**Domain Adaptation** Most recognition algorithms are developed and evaluated on the same data distributions, e.g, the public datasets ImageNet, MS-COCO, CIFAR-10, MNIST. In real applications, people often confront performance degradation when apply a classifier trained on a source domain to a target domain.

In unsupervised domain adaptation, source domain has labeled data  $\mathcal{D}_s = \{x_i^s, y_i^s\}_{i=1}^{n_s}$  while target domain contains data without labels  $\mathcal{D}_t = \{x_i^t\}_{i=1}^{n_t}$ . The goal is to learn a classifier  $f : x_i^t \mapsto y_i^t$  for the unseen target samples by exploring the knowledge learnt from the source domain. Domain adaptation algorithms attempts to transfer knowledge across domains by solving the domain shift problem, i.e., the data-label distributions  $p(x^s, y^s)$  and  $p(x^t, y^t)$  are different.

There are many approaches to address this issue. One is to extract transferable features that are invariant across domains [49, 73], or learn representative hash codes [74] to find a common latent space where the classifier can be used without considering the data's origin. Another trend is to learn the transformation between domains [75] to align the source and target datapoints through barycentric mapping, and train a classifier on the transferred source data. Courty [76] and Damodaran [50] proposed to look for a transformation that matches the data-label joint distributions  $p(x^s, y^s)$  in source domain to its equivalent version  $p(x^t, y^t)$  in target domain. The predictive function  $f$  is learnt by minimizing the optimal transport loss between the distributions  $p(x^s, y^s)$  and  $p(x^t, f(x^t))$ . As a by-product, minimizing the optimal transport cost is equivalent to mapping a source domain sample to a target domain sample with similar semantic content, and this is the domain transfer problem.

**Voice Conversion** Voice conversion (VC) aims to change a speaker's voice to make it sounds like spoken by another person. It is a special case of voice transformation (VT), whose goal is to modify human speech without changing its content. VC transforms speaker identity by replacing speaker-dependent components of the signal while maintaining the linguistic information. Speech quality and speaker similarity are two important factors to evaluate a VC system. There are a bunch of VC applications, such as movie dubbing, personalized TTS (Text To Speech) systems, speaker accent or emotion transformation, speaking-aid devices, call quality enhancement, etc.

Most VC frameworks involve three steps: feature extraction, feature conversion, waveform generation. In speech analysis, waveform signals are encoded into feature representations that are easy to control and modify. Spectral envelope, mel-cepstrum, fundamental frequency ( $f_0$ ), formant frequencies and bandwidths are the most widely used features to represent speech in short-time segments. To capture contextual information across frames, implicit methods such as hidden Markov models (HMMs), Long Short-Term Memory (LSTM) and recurrent neural networks (RNNs) [77] were developed.

The main work in VC is to transform the source feature sequences to target feature sequences that capture the speaker identity. Most traditional VC systems perform frame-by-frame mapping under the assumption that speech segments are independent from each other. Some recent models such as HMM and RNN incorporate speech dynamics implicitly. There are four typical approaches to learn the mapping function: codebook mapping (e.g., Vector quantization (VQ) [78]), mixed linear mappings (e.g., Gaussian mixture model (GMM) [79]), neural network mapping (e.g., RBM, DNN, RNN [80]), and exemplar-based mapping (e.g., non-negative matrix factorization (NMF) [81]). Beyond these, an autoregressive neural network model called WaveNet [56] was proposed. It can directly learn the mapping based on raw audio and generate speech waveforms conditioning on the speaker identity.

There are various assumptions in speech analysis and waveform generation. Source-filter models assume speech to be generated by excitation signals passing through a vocal tract, and encode speech waveforms as acoustic features that represent sound source and vocal tract independently. However, the original phase information will lose under this assumption. At conversion time, the converted target features are passed through a vocoder based on the source filter model to reconstruct the waveform. Quality degradation may happen due to the inaccurate assumption. Iterative phase reconstruction algorithm Griffin-Lim [82] was adopted to alleviate this issue. Harmonic plus noise models (HNM) [83] assume speech to be a combination of a noise component and a harmonic component, i.e., sinusoidal waves with frequencies relevant to pitch. Speech is parameterized by the fundamental frequency  $f_0$  and a spectrum which consists of a lower band of harmonic and a higher band of noise. Other assumptions include stationary speech signal, frame-by-frame mapping, time-invariant linear filter, etc. Recently, Tamamori et. al [84] proposed a speaker-dependent WaveNet vocoder that does not require explicit modeling of excitation signals and those assumptions.

In terms of conversion conditions, VC can be categorized into parallel and non-parallel, text-dependent and text-independent systems [85]. In parallel systems, the training corpus consists of paired recordings from the source and target speakers with same linguistic contents. The shared acoustic features can be used to train the mapping model. To get parallel feature sequences of equal length, a time-alignment step must be included to remove the temporal differences in the recordings, for example, the dynamic time warping (DTW) [78] algorithm. Phoneme transcriptions are also useful for time alignment. Non-parallel system does not require sentences with the same linguistic contents. It is much more useful and practical because non-parallel speech data is easier to collect and therefore can get larger training sets. There are several ways to learn the mapping without paired data: (1) use unit selection [86] to choose matched linguistic feature pairs; (2) build pseudo parallel sentences on extra automatic speech recognition (ASR) modules [87]; (3) extract speaker-independent features in shared latent space [88]; (4) use unpaired image-to-image translation approaches [46][47, 70].

Parallel, text-dependent systems are supposed to have better performance. However, parallel utterance pairs are difficult to get. Most parallel VC systems require time alignment to extract parallel source-target features. The misalignment in automatic time alignment algorithms often leads to degradation in speech quality, while manual correction is arduous. Recently, the winner of VC Challenge 2018 [89] showed their algorithm can achieve similar results in both parallel and non-parallel settings. It first uses a lot of external speech data with phonetic transcriptions to train a speaker-independent content-posterior-feature extractor, followed by a speaker-dependent LSTM-RNN to predict fundamental frequency  $f_0$  and STRAIGHT spectral features [90], and then reconstruct the waveforms with a speaker-dependent WaveNet vocoder [84]. Moreover, Kaneko et.al [91] and Fang et. al [83] claimed their nonparallel, text-independent VC algorithms based on CycleGAN [44] perform comparable to or better than the state-of-the-art parallel approaches.

## 5.4 Method

We propose a general approach for style transfer. Let  $x_i \in \mathbb{X}_i$  be a datapoint sampled from domain  $i$ . The goal is to learn a conversion model  $p(x_j | x_i)$  that maps  $x_i$  to domain  $j$  ( $j \neq i$ ) by changing its style and preserving the original content. For unpaired training data, we have marginal distributions  $p(x_i), p(x_j)$  instead of the joint  $p(x_i, x_j)$ , so it requires additional constraints to determine  $p(x_j | x_i)$ . Inspired by disentangled representation learning in [29], we assume that each example  $x_i \in \mathbb{X}_i$  can be decomposed into a content code  $c \in \mathcal{C}$  that encodes domain-invariant information and a style code  $s_i \in \mathcal{S}_i$  that encodes domain-dependent information.  $\mathcal{C}$  is shared across domains and contains the information we want to preserve.  $\mathcal{S}_i$  is domain-specific and contains the information we want to change. In conversion stage, we extract the content code of  $x_i$  and recombine it with a style code

randomly sampled from  $\mathcal{S}_j$ . A generative adversarial network (GAN) [31] is added to ensure that the converted samples are indistinguishable from the real ones.

Figure 16 shows the autoencoder model of style transfer with a partially shared latent space. Any pair of corresponding datapoints  $(x_i, x_j)$  is assumed to have a shared latent code  $c \in \mathcal{C}$  and domain-specific style codes  $s_i \in \mathcal{S}_i$ ,  $s_j \in \mathcal{S}_j$ . The generative model of domain  $i$  is an autoencoder that consists of a deterministic decoder  $x_i = G_i(c_i, s_i)$  and two encoders  $c_i = E_i^c(x_i)$ ,  $s_i = E_i^s(x_i)$ . The encoder  $E_i = (E_i^c, E_i^s)$  and decoder  $G_i$  are inverse operations such that  $E_i = G_i^{-1}$ . To generate converted sample  $x'_{j \leftarrow i}$ , we just extract and recombine the content code of  $x_i$  with the style code of domain  $j$ .

$$\begin{aligned} x'_{i \leftarrow j} &= G_i(c_j, s_i) = G_i(E_j^c(x_j), s_i) \\ x'_{j \leftarrow i} &= G_j(c_i, s_j) = G_j(E_i^c(x_i), s_j) \end{aligned} \quad (31)$$

It should be noted that the style code  $s_i$  is not inferred from one example, but learnt from the entire target domain  $j$ , which is a major difference from [29]. This is because the style extracted from a single example is ambiguous and may not capture the general characteristics of the target domain. To restrict the mapping  $p(x_j | x_i)$ , we use an constraint that is slightly different from the cycle consistency [44]. It assumes that an example converted to another domain and converted back should be unchanged, i.e.,  $x''_{i \leftarrow j \leftarrow i} = x_i$ . Instead, we apply a semi-cycle consistency in the latent space by assuming that only the latent codes remain unchanged  $E_i^c(x'_{i \leftarrow j}) = c_i$  and  $E_i^s(x'_{i \leftarrow j}) = s_i$ . The relaxed constraint provides diversity to the generated samples.

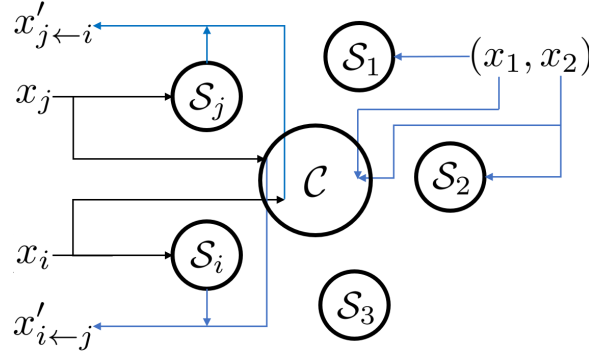


Figure 16: Autoencoder model with partially shared latent space. Sample  $x_i$  is encoded into a domain-specific space  $\mathcal{S}_i$  and a shared content space  $\mathcal{C}$ . Corresponding datapoints  $(x_1, x_2)$  are encoded in the same content code. Converted sample  $x'_{j \leftarrow i}$  is generated by recombining the original content code of  $x_i$  and the style code randomly sampled from  $\mathcal{S}_j$ .

We formulate the learning problem as a distributionally robust game. Each domain  $i$  has four agents  $E_i^c, E_i^s, G_i, D_i$ , in which  $D_i$  is a discriminator with two objectives. One is to distinguish between real samples and machine-generated samples, the other is to classify domain labels. On the other side, the generators  $G_i$  have two purposes: synthesize realistic samples and convert them into domain  $i$ . For example in voice conversion, the synthesized speech is evaluated on both the naturalness of its quality and the correctness of speaker's identity.

#### 5.4.1 Two-domain transformation

There are  $4n$  agents for  $n$  domains. For simplicity, we first investigate the conversion model between two domains. When  $n = 2$ , there are 8 agents:  $E_1^c, E_1^s, G_1, D_1$  for domain  $\mathbb{X}_1$  and  $E_2^c, E_2^s, G_2, D_2$  for domain  $\mathbb{X}_2$ . Each agent has a different objective, and the utility function of one agent depends on the action of other agents. Collaboration and competition exist among them. So this is a distributionally robust game with cooperative agents and uncertain utility functions. Based on the analysis above, there are five modules need to learn:

- ① domain-invariant content encoder
- ② domain-specific style encoders
- ③ real/fake discriminator

- 704 ④ domain classifier  $D_i$   
 705 ⑤ fake samples generator

706 The encoders and decoders form a group to synthesize converted samples in the target domain.  
 707 Another group is the discriminator and classifier. They work on the opposite side to distinguish  
 708 between real/fake samples and predict the domain label. The intergroup competition and intragroup  
 709 collaboration are listed in table 3.

710 How to define the game?

711 When Nash equilibrium reaches, the autoencoder  $(E_i^{c*}, E_i^{s*}, G_i^*)$  minimizes the reconstruction error  
 712  $L_{rec}^x \rightarrow 0, L_{rec}^c \rightarrow 0, L_{rec}^s \rightarrow 0$ . The GAN network  $(D_i^*, G_i^*)$  and adversarial loss  $L_{GAN}^{x_i}$  converge  
 713 at saddle points that minimize the distance between  $p(x_i)$  and  $p(x'_{j \leftarrow i})$ . The classifier  $D_i^{cls*}$  correctly  
 714 predicts the domain category of both real and fake samples  $D_i^{cls*}(x_i) = i, D_i^{cls*}(x'_{i \leftarrow j}) = i$ .

Table 3: Cooperative game

Intergroup competition	Learning module	Objective
$E_1^s, E_2^s$	②	$\min L_{cyc}^s$
$D_1, D_2$	④	$\min L_{cls}^x$
Intragroup collaboration	Learning module	Objective
$E_1^c, E_1^s, G_1$	<i>Autoencoder</i> <sub>1</sub>	$\min L_{rec}^{x_1}$
$E_2^c, E_2^s, G_2$	<i>Autoencoder</i> <sub>2</sub>	$\min L_{rec}^{x_2}$
$G_1, G_2$	⑤	$\min L_{GAN}^x$
$D_1, D_2$	③	$\max L_{GAN}^x$
$E_1^c, E_2^c$	①	$\min L_{cyc}^c$

715 The collaboration means, agents in different domains can help each other as they may have common  
 716 interests. For instance, a sample can be used to update the real/fake discriminator even if its class  
 717 label is missing. Except for the autoencoders  $E_i^c, E_i^s, G_i$ , other coalitions are cross-domain.  $(G_1, G_2)$   
 718 works on data synthesis,  $D_1, D_2$  works on real/fake discrimination,  $E_1^c, E_2^c$  extracts high-level content  
 719 information that we want to preserve.

720 We jointly train the encoders, decoders and GAN’s discriminators with multiple objectives. To  
 721 keep encoder and decoder as inverse operations, a reconstruction loss is applied in the direction  
 722  $x_i \rightarrow (c_i, s_i) \rightarrow x'_i, (i, j \in 1, 2)$ . Sample  $x_i$  should not be changed after encoding and decoding.

$$L_{rec}^{x_i} = \mathbb{E}_{x_i}(\|x_i - x'_i\|_1), \quad x'_i = G_i(E_i^c(x_i), E_i^s(x_i)) \quad (32)$$

723 In our model, the latent space is partially shared. Thus the cycle consistency constraint [44] is not  
 724 preserved, i.e.,  $x''_{1 \leftarrow 2 \leftarrow 1} \neq x_1$ . We apply a semi-cycle loss in the coding direction  $c_1 \rightarrow x'_{2 \leftarrow 1} \rightarrow$   
 725  $c'_{2 \leftarrow 1}$  and  $s_2 \rightarrow x'_{2 \leftarrow 1} \rightarrow s'_{2 \leftarrow 1}$ .

$$\begin{aligned} L_{cyc}^{c_1} &= \mathbb{E}_{c_1, s_2}(\|c_1 - c'_{2 \leftarrow 1}\|_1), \quad c'_{2 \leftarrow 1} = E_2^c(x'_{2 \leftarrow 1}) \\ L_{cyc}^{s_2} &= \mathbb{E}_{c_1, s_2}(\|s_2 - s'_{2 \leftarrow 1}\|_1), \quad s'_{2 \leftarrow 1} = E_2^s(x'_{2 \leftarrow 1}) \end{aligned} \quad (33)$$

726 Moreover, we add a GAN module to ensure the quality of generated samples. They should be  
 727 indistinguishable from the real samples in the target domain  $\mathbb{X}_i$ . GAN loss is computed between  $x_j$   
 728 and  $x'_{i \leftarrow j}$  to represent the distance between two distributions  $p(x_j), p(x'_{i \leftarrow j})$ .

$$L_{GAN}^{x_i} = \mathbb{E}_{c_j, s_i}[\log(1 - D_i(x'_{i \leftarrow j}))] + \mathbb{E}_{x_i}[\log D_i(x_i)] \quad (34)$$

729 The full loss is the weighted sum of  $L_{recon}, L_{cycle}, L_{GAN}$ .

$$\begin{aligned} &\min_{E_1^c, E_1^s, E_2^c, E_2^s, G_1, G_2} \max_{D_1, D_2} L(E_1^c, E_1^s, E_2^c, E_2^s, G_1, G_2, D_1, D_2) \\ &= \lambda_s(L_{cyc}^{s_1} + L_{cyc}^{s_2}) + \lambda_c(L_{cyc}^{c_1} + L_{cyc}^{c_2}) + \lambda_x(L_{rec}^{x_1} + L_{rec}^{x_2}) + \lambda_g(L_{GAN}^{x_1} + L_{GAN}^{x_2}) \end{aligned} \quad (35)$$

730 where  $\lambda_s, \lambda_c, \lambda_x, \lambda_g$  control the weights of the components.



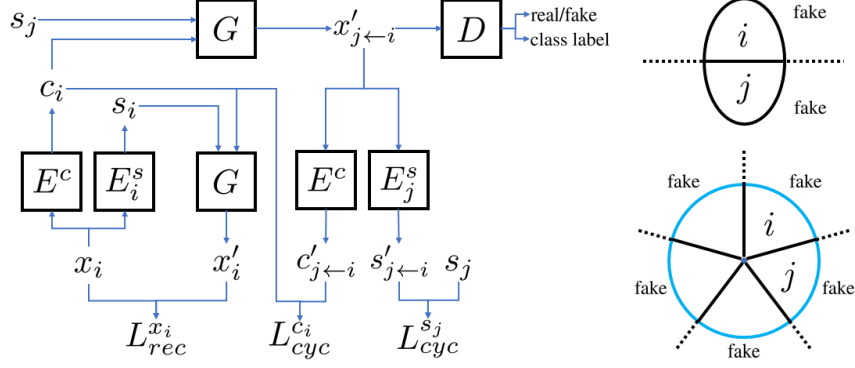


Figure 17: Learn multi-domain transformation. Left: conversion model with shared content encoder, decoder and classifier. Top-right: two domains. Bottom-right: multiple domains ( $n=5$ ).

#### 5.4.2 Multi-domain transformation

In multi-domain case, there are  $4n$  agents in the game. To reduce complexity, we replace the domain-specific models  $E_i^c, G_i, D_i$  with a shared content encoder  $E^c$ , a shared decoder  $G$ , and a single multiclass classifier  $D^{cls}$ . Thus, only  $n + 3$  agents left. (form coalition?) Figure 17 shows the multi-domain transformation model.

When  $n = 2$ , there are 3 kinds of data: real in  $\mathbb{X}_1$ , real in  $\mathbb{X}_2$  and fake in  $\mathbb{X}_{fake}$ . Two real/fake discriminators  $D_1, D_2$  are enough for classification. If this idea is extended to multiple domains, there will be  $n$  binary discriminators and  $2^n$  outputs. Instead, we replace them with one binary real/fake discriminator  $D$  and one multiclass domain classifier  $D^{cls}$ . The two-step classification not only reduces complexity ( $2n$  outputs) but also makes the most of the training data. An example in domain  $i$  is also useful to train the generative model for other domains, as it is the common interest of all agents  $G_1, G_2, \dots G_n$  to synthesize realistic data. Therefore, the multi-domain transformation model can be trained on different datasets with partially labeled data.

In the following sections, we will apply the proposed method on two real-world applications: image style transfer and emotinal speech conversion.

#### 5.5 Case Study: Image Style Transfer

Image style transfer, or image-to-image translation is a hot topic in computer vision. The objective is to transfer the visual style of an image while keep its semantic content. Similar tasks includes texture synthesis, artistic style transfer, photorealistic image style transfer, etc.

Gatys et al. [62] introduced an algorithm to separate and recombine the content and style of natural images. They claimed that a Convolutional Neural Network (CNN) is the ideal representation to factorize semantic content and artistic style. High-level features are extracted by higher layers of the network, while low-level features are captured by the correlations between filter responses in various layers. However, the style is learnt from a single image, which limits the ability to capture the general theme of the target domain.

Instead of example-based style transfer, we will focus on learning the translation model between two image distributions. There are a broad range of researches on image synthesis and representation learning. We will expore the state-of-the-art neural network architecture and implement our model to learn the disentangled representations of image style and content. The proposed method will be tested on a typical image style transfer task.

#### 5.6 Case Study: Emotional Voice Conversion

Voice transformation (VT) is a technique to modify some properties of human speech while preserving its linguistic information. VT can be applied to change the speaker identity, i.e., voice conversion (VC) [85], or to transform the speaking style of a speaker, such as emotion and accent conversion [92]. In this section, we will focus on emotion voice transformation. The goal is to change emotion-related

766 characteristics of a speech signal while preserving its linguistic content and speaker identity. Emotion  
 767 conversion techniques can be applied to various tasks, such as hiding negative emotions for customer  
 768 service agents, helping film dubbing, and creating more expressive voice messages on social media.

769 Existing VC approaches cannot be applied directly because they change speaker identity by assuming  
 770 pronunciation and intonation to be a part of the speaker-independent information. Since the speaker’s  
 771 emotion is mainly conveyed by prosodic aspects, some studies have focused on modelling prosodic  
 772 features such as pitch, tempo, and volume [93, 94]. In [95], a rule-based emotional voice conversion  
 773 system was proposed. It modifies prosody-related acoustic features of neutral speech to generate  
 774 different types of emotions. A speech analysis-synthesis tool STRAIGHT [90] was used to extract  
 775 fundamental frequency ( $F_0$ ) and power envelope from raw audio. These features were parameterized  
 776 and modified based on Fujisaki model [96] and target prediction model [97]. The converted features  
 777 were then fed back into STRAIGHT to re-synthesize speech waveforms with desired emotions.  
 778 However, this method requires temporal aligned parallel data that is difficult to obtain in real  
 779 applications; and the accurate time alignment needs manual segmentation of the speech signal at  
 780 phoneme level, which is very time consuming.

781 To address these issues, we propose a nonparallel training method. Instead of learning example based  
 782 one-to-one mapping between paired emotional utterances  $(x_1, x_2)$ , we learn the conversion model  
 783 between two emotion domains  $(\mathbb{X}_1, \mathbb{X}_2)$ .

784 Inspired by the disentangled representation learning in image style transfer [47, 62], we assume that  
 785 each speech signal  $x_i \in \mathbb{X}_i$  can be decomposed into a content code  $c \in \mathcal{C}$  that represents emotion-  
 786 invariant information and a style code  $s_i \in \mathcal{S}_i$  that represents emotion-dependent information.  $\mathcal{C}$   
 787 is shared across domains and contains the information we want to preserve.  $\mathcal{S}_i$  is domain-specific  
 788 and contains the information we want to change. In conversion stage, we extract content code of  
 789 the source speech and recombine it with style code of the target emotion. A generative adversarial  
 790 network (GAN) [31] is added to improve the quality of converted speech. Our approach is nonparallel,  
 791 text-independent, and does not rely on any manual operation.

792 For implementation, we use deep neural networks to learn the latent representations of speech.  
 793 Specifically, the encoders and decoders are implemented with one-dimensional CNNs to capture the  
 794 temporal dependences. The GAN discriminators are implemented with two-dimensional CNNs to  
 795 capture the spectra-temporal patterns. All networks are equipped with gated linear units (GLU) [98]  
 796 as activation functions.

797 We plan to test our approach on IEMOCAP [99] to learn the conversion models for four emotions:  
 798 angry, happy, neutral, sad. IEMOCAP is a nonparallel dataset widely used in emotional speech  
 799 recognition and analysis. It contains scripted and improvised dialogs in five sessions; each has labeled  
 800 emotional sentences pronounced by two professional English speakers. The emotions in scripted  
 801 dialogs have strong correlation with the lingual content. Since our task is to change emotion but keep  
 802 the speaker identity and linguistic content, we only use the improvised dialogs of the same speaker.

803 There are three metrics for performance evaluation: emotion correctness, voice quality and the ability  
 804 to retain speaker identity. For subjective evaluation, we will conduct listening tests on Amazon  
 805 MTurk to evaluate the converted speech. Each example is evaluated by a group of random listeners.  
 806 They will be asked to manually classify the emotion, and give 1-to-5 opinion scores on voice quality  
 807 and the similarity with the original speaker. The classification result and mean opinion score (MOS)  
 808 are two major measurements. For objective evaluation, we plan to use the state-of-the-art speech  
 809 emotion classifier [100] to check the emotion category of generated speech. It indicates a success if  
 810 our model can increase the proportion of the target emotions and reduce the original emotions. To our  
 811 knowledge, this is the first work for nonparallel emotion conversion. If there’s time, we will develop  
 812 multidomain emotion conversion models for unseen speakers.

## 813 6 Dissertation Outline and Research Timeline

### 814 6.1 Dissertation Outline

815 Below is the proposed outline of the planned dissertation.

- 816 1. Introduction
  - 817 1.1. Generative Learning Background
  - 818 1.2. Game Theory Background
  - 819 1.3. Motivation
  - 820 1.4. Summary of Contributions
- 821 2. Distributionally Robust Games
  - 822 2.1. Introduction
  - 823 2.2. Related Work
  - 824 2.3. Problem Formulation
  - 825 2.4. Minimax Robust Game
    - 826 2.4.1. Infinite Dimension
    - 827 2.4.2. Triality Theory
    - 828 2.4.3. Dimension Reduction
  - 829 2.5. Case Study: Unsupervised Generative Learning
- 830 3. Wassertein Metric
  - 831 3.1. Introduction
  - 832 3.2. Optimal Transportation Problem
  - 833 3.3. From KL divergence to Wasserstein Metric
  - 834 3.4. Displacement Interpolation
  - 835 3.5. Case Study: A Toy Example Simulation
- 836 4. Learning Generative Models
  - 837 4.1. Distributionally Robust Optimization
  - 838 4.2. Numerical Investigation
  - 839 4.3. Train a Deep Generative Model
  - 840 4.4. Case Study: Unsupervised Learning for Clustering
  - 841 4.5. Case Study: Generative Modeling for Image Synthesis
- 842 5. Image Style Transfer
  - 843 5.1. Introduction
  - 844 5.2. Motivation
  - 845 5.3. Related Work
  - 846 5.4. Method
    - 847 5.4.1. Two-domain Transformation
    - 848 5.4.2. Multi-domain Transformation
  - 849 5.5. Experiments
- 850 6. Emotional Voice Conversion
  - 851 6.1. Introduction
  - 852 6.2. Related Work
    - 853 6.2.1. Emotion-related Features
    - 854 6.2.2. Nonparallel Training Approaches
    - 855 6.2.3. Disentangled Representation Learning
  - 856 6.3. Method
    - 857 6.3.1. Assumptions
    - 858 6.3.2. Model
    - 859 6.3.3. Loss functions

860	6.4. Experiments
861	6.4.1. Corpus
862	6.4.2. Network Structure
863	6.4.3. Performance Evaluation
864	6.5. Discussions
865	7. Conclusions and Future Work
866	Bibliography

## 867 6.2 Research Timeline

Fall 2018	Study on Generative Models for Voice Conversion
Spring 2019	Write Thesis
Fall 2019	Thesis Defense

## 868 7 Conclusion

869 This proposal contains my plan for graduating with a PhD in the early Fall of 2019, and described how  
870 game theory models can be used in generative learning. While this thesis focus on image synthesis  
871 and speech conversion, it is a general model and can apply to other kinds of data by replacing the  
872 data representation network. Its contributions can help shed a light to a much unexplored field, and  
873 help create a path towards fully automatic data generation and manipulation.

## 8 List of Publications

### 8.1 Thesis Related Publications

The following publications are directly related to my dissertation.

1. Jian Gao, Deep Chakraborty, Hamidou Tembine, and Olaitan Olaleye, "Nonparallel emotional speech conversion," CoRR, vol. abs/1811.01174. (submitted to ICASSP 2019)
2. Jian Gao and Hamidou Tembine, Distributionally Robust Games for Deep Generative Learning, July 2018. DOI: 10.13140/RG.2.2.15305.44644
3. Jian Gao, Yida Xu, Julian Barreiro-Gomez, Massa Ndong, Michalis Smyrnakis and Hamidou Tembine (September 5th 2018) Distributionally Robust Optimization. In Jan Valdman, Optimization Algorithms, IntechOpen. DOI: 10.5772/intechopen.76686. ISBN: 978-1-78923-677-4
4. Jian Gao and Hamidou Tembine, Distributionally Robust Games: Wasserstein Metric, International Joint Conference on Neural Networks (IJCNN), Rio de Janeiro, Brazil, July 2018
5. Jian Gao and Hamidou Tembine, Bregman Learning for Generative Adversarial Networks, Chinese Control and Decision Conference (CCDC), Shenyang, China, June 2018 (*Best Paper Finalist Award*)
6. Dario Bauso, Jian Gao and Hamidou Tembine, Distributionally Robust Games: f-Divergence and Learning, 11th EAI International Conference on Performance Evaluation Methodologies and Tools (VALUETOOLS), Venice, Italy, Dec 2017

### 8.2 Relevant Publications

These papers are relevant publications on other generative models.

1. J. Gao and H. Tembine, "Distributed Mean-Field-Type Filters for Traffic Networks," in IEEE Transactions on Intelligent Transportation Systems. doi: 10.1109/TITS.2018.2816811
2. J. Gao and H. Tembine, "Correlative mean-field filter for sequential and spatial data processing," IEEE EUROCON 2017 -17th International Conference on Smart Technologies, Ohrid, 2017, pp. 243-248. doi: 10.1109/EUROCON.2017.8011113
3. Jian Gao and Hamidou Tembine, Distributed Mean-Field-Type Filter for Vehicle Tracking, in American Control Conference (ACC), Seattle, USA, May 2017 (*Student Travel Award*)
4. J. Gao and H. Tembine, "Distributed Mean-Field-Type Filters for Big Data Assimilation," 2016 IEEE 18th International Conference on High Performance Computing and Communications; IEEE 14th International Conference on Smart City; IEEE 2nd International Conference on Data Science and Systems (HPCC/SmartCity/DSS), Sydney, NSW, 2016, pp. 1446-1453. doi: 10.1109/HPCC-SmartCity-DSS.2016.0206

### 8.3 Other Publications

The papers below were published during my PhD, but will not be included in the thesis, as they diverge too much from the main topic of the research.

1. J. Gao and H. Tembine, "Empathy and berge equilibria in the forwarding dilemma in relay-enabled networks," 2017 International Conference on Wireless Networks and Mobile Communications (WINCOM), Rabat, 2017, pp. 1-8. doi: 10.1109/WINCOM.2017.8238199 (*Best paper Award*)
2. Fanhuai Shi, Jian Gao, Xixia Huang, An affine invariant approach for dense wide baseline image matching. International Journal of Distributed Sensor Networks (IJDSN) 12(12) (2016)

## References

- [1] J. Deng, W. Dong, R. Socher, L. Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pages 248–255, June 2009.
- [2] Tero Karras, Timo Aila, Samuli Laine, and Jaakko Lehtinen. Progressive growing of gans for improved quality, stability, and variation. *CoRR*, abs/1710.10196, 2017.
- [3] Hao-Wen Dong, Wen-Yi Hsiao, Li-Chia Yang, and Yi-Hsuan Yang. Musegan: Multi-track sequential generative adversarial networks for symbolic music generation and accompaniment. In *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence, (AAAI-18), the 30th innovative Applications of Artificial Intelligence (IAAI-18), and the 8th AAAI Symposium on Educational Advances in Artificial Intelligence (EAAI-18), New Orleans, Louisiana, USA, February 2-7, 2018*, pages 34–41, 2018.
- [4] Yunjey Choi, Min-Je Choi, Munyoung Kim, Jung-Woo Ha, Sunghun Kim, and Jaegul Choo. Stargan: Unified generative adversarial networks for multi-domain image-to-image translation. *CoRR*, abs/1711.09020, 2017.
- [5] Afroz Ahamad. Generating text through adversarial training using skip-thought vectors. *CoRR*, abs/1808.08703, 2018.
- [6] Noam Mor, Lior Wolf, Adam Polyak, and Yaniv Taigman. A universal music translation network. *CoRR*, abs/1805.07848, 2018.
- [7] Cédric Villani. *Optimal Transport: Old and New*. Grundlehren der mathematischen Wissenschaften. Springer, 2009 edition, September 2008.
- [8] Solomon Kullback. *Information theory and statistics*. Peter Smith, 1978.
- [9] B. Fuglede and F. Topsøe. Jensen-shannon divergence and hilbert space embedding. In *International Symposium on Information Theory, 2004. ISIT 2004. Proceedings.*, pages 31–, June 2004.
- [10] Michele Aghassi and Dimitris Bertsimas. Robust game theory. *Mathematical Programming*, 107(1):231–273, Jun 2006.
- [11] Dimitris Bertsimas, David B. Brown, and Constantine Caramanis. Theory and applications of robust optimization. *SIAM Review*, 53(3):464–501, 2011.
- [12] H.E. Scarf. *A Min-max Solution of an Inventory Problem*. P (Rand Corporation). Rand Corporation, 1957.
- [13] Maurice Sion. On general minimax theorems. *Pacific J. Math.*, 8(1):171–176, 1958.
- [14] Allen L. Soyster. Technical note - convex programming with set-inclusive constraints and applications to inexact linear programming. *Operations Research*, 21(5):1154–1157, 1973.
- [15] Imre Csiszár. Eine informationstheoretische ungleichung und ihre anwendung auf den beweis der ergodizität von markoffschen ketten. *Magyar. Tud. Akad. Mat. Kutató Int. Közl.*, 8:85–108, 1963.
- [16] Tetsuzo Morimoto. Markov processes and the h-theorem. *Journal of the Physical Society of Japan*, 18(3):328–331, 1963.
- [17] S. M. Ali and S. D. Silvey. A general class of coefficients of divergence of one distribution from another. *Journal of the Royal Statistical Society. Series B*, 28(1):131–142, 1966.
- [18] Hamidou Tembine. *Distributed Strategic Learning for Wireless Engineers*. CRC Press, Inc., Boca Raton, FL, USA, 2012.
- [19] Evolutionary dynamics over continuous action spaces for population games that arise from symmetric two-player games. *Journal of Economic Theory*, 148(2):743 – 777, 2013.

- 963 [20] Lăcrăşău Pavel. *Game Theory for Control of Optical Networks*. Static & Dynamic Game Theory:  
964 Foundations & Applications. 1 edition.
- 965 [21] Gradient dynamics in population games: Some basic results. *Journal of Mathematical*  
966 *Economics*, 46(5):691 – 707, 2010. Mathematical Economics: Special Issue in honour of  
967 Andreu Mas-Colell, Part 1.
- 968 [22] B. D. Nichols. A comparison of action selection methods for implicit policy method reinforce-  
969 ment learning in continuous action-space. In *2016 International Joint Conference on Neural*  
970 *Networks (IJCNN)*, pages 3785–3792, July 2016.
- 971 [23] Stochastic fictitious play with continuous action sets. *Journal of Economic Theory*, 152:179 –  
972 213, 2014.
- 973 [24] Pairwise comparison dynamics for games with continuous strategy space. *Journal of Economic*  
974 *Theory*, 153:344 – 375, 2014.
- 975 [25] Dario Bauso, Hamidou Tembine, and Tamer Başar. Robust mean field games. *Dynamic Games*  
976 *and Applications*, 6(3):277–303, Sep 2016.
- 977 [26] I. Csiszár. Information-type measures of difference of probability distributions and indirect  
978 observations. *Studia Scientiarum Mathematicarum Hungarica*, 2:299–318, 1967.
- 979 [27] Ian Goodfellow and et al. Generative adversarial nets. In *NIPS*. 2014.
- 980 [28] Yuhuai Wu, Yuri Burda, Ruslan Salakhutdinov, and Roger B. Grosse. On the quantitative  
981 analysis of decoder-based generative models. *CoRR*, abs/1611.04273, 2016.
- 982 [29] Leon A. Gatys, Alexander S. Ecker, and Matthias Bethge. A neural algorithm of artistic style.  
983 *CoRR*, abs/1508.06576, 2015.
- 984 [30] Diederik P. Kingma and Max Welling. Auto-encoding variational bayes. *CoRR*, abs/1312.6114,  
985 2013.
- 986 [31] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil  
987 Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In Z. Ghahramani,  
988 M. Welling, C. Cortes, N. D. Lawrence, and K. Q. Weinberger, editors, *Advances in Neural*  
989 *Information Processing Systems 27*, pages 2672–2680. Curran Associates, Inc., 2014.
- 990 [32] Martín Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein GAN. *CoRR*,  
991 abs/1701.07875, 2017.
- 992 [33] Martín Arjovsky and Léon Bottou. Towards principled methods for training generative  
993 adversarial networks. *CoRR*, abs/1701.04862, 2017.
- 994 [34] Tim Salimans, Ian J. Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, and  
995 Xi Chen. Improved techniques for training gans. In *Advances in Neural Information Processing*  
996 *Systems 29: Annual Conference on Neural Information Processing Systems 2016, December*  
997 *5-10, 2016, Barcelona, Spain*, pages 2226–2234, 2016.
- 998 [35] M. Lichman. UCI machine learning repository, 2013.
- 999 [36] Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the  
1000 wild. In *Proceedings of International Conference on Computer Vision (ICCV)*, 2015.
- 1001 [37] Alec Radford, Luke Metz, and Soumith Chintala. Unsupervised representation learning with  
1002 deep convolutional generative adversarial networks. *CoRR*, abs/1511.06434, 2015.
- 1003 [38] Richard Sinkhorn and Paul Knopp. Concerning nonnegative matrices and doubly stochastic  
1004 matrices. *Pacific J. Math.*, 21(2):343–348, 1967.
- 1005 [39] Sebastian Ruder. An overview of gradient descent optimization algorithms. *CoRR*,  
1006 abs/1609.04747, 2016.

- 1007 [40] Lucas Theis, Aäron van den Oord, and Matthias Bethge. A note on the evaluation of generative  
1008 models. *CoRR*, abs/1511.01844, 2015.
- 1009 [41] K. Bousmalis, N. Silberman, D. Dohan, D. Erhan, and D. Krishnan. Unsupervised pixel-  
1010 level domain adaptation with generative adversarial networks. In *2017 IEEE Conference on*  
1011 *Computer Vision and Pattern Recognition (CVPR)*, pages 95–104, July 2017.
- 1012 [42] Ashish Shrivastava, Tomas Pfister, Oncel Tuzel, Joshua Susskind, Wenda Wang, and Russell  
1013 Webb. Learning from simulated and unsupervised images through adversarial training. In *The*  
1014 *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017.
- 1015 [43] Yaniv Taigman, Adam Polyak, and Lior Wolf. Unsupervised cross-domain image generation.  
1016 *CoRR*, abs/1611.02200, 2016.
- 1017 [44] JY Zhu, T. Park, P. Isola, and A.A. Efros. Unpaired image-to-image translation using cycle-  
1018 consistent adversarial networks. In *IEEE ICCV*, Oct 2017.
- 1019 [45] Mingming Gong, Kun Zhang, Tongliang Liu, Dacheng Tao, Clark Glymour, and Bernhard  
1020 Schölkopf. Domain adaptation with conditional transferable components. In *Proceedings of*  
1021 *the 33rd International Conference on International Conference on Machine Learning - Volume*  
1022 *48, ICML’16*, pages 2839–2848. JMLR.org, 2016.
- 1023 [46] M.Y. Liu, T. Breuel, and J. Kautz. Unsupervised image-to-image translation networks. In  
1024 *Advances in Neural Information Processing Systems (NIPS)*, pages 700–708. Curran Associates,  
1025 Inc., 2017.
- 1026 [47] X. Huang, MY. Liu, S. Belongie, and J. Kautz. Multimodal unsupervised image-to-image  
1027 translation. In *The European Conference on Computer Vision (ECCV)*, September 2018.
- 1028 [48] S. Si, D. Tao, and B. Geng. Bregman divergence-based regularization for transfer subspace  
1029 learning. *IEEE Transactions on Knowledge and Data Engineering*, 22(7):929–942, July 2010.
- 1030 [49] Yaroslav Ganin, Evgeniya Ustinova, Hana Ajakan, Pascal Germain, Hugo Larochelle, François  
1031 Laviolette, Mario Marchand, and Victor Lempitsky. Domain-adversarial training of neural  
1032 networks. *J. Mach. Learn. Res.*, 17(1):2096–2030, January 2016.
- 1033 [50] Bharath Bhushan Damodaran, Benjamin Kellenberger, Rémi Flamary, Devis Tuia, and Nicolas  
1034 Courty. Deepjdot: Deep joint distribution optimal transport for unsupervised domain adaptation.  
1035 In Vittorio Ferrari, Martial Hebert, Cristian Sminchisescu, and Yair Weiss, editors, *Computer*  
1036 *Vision – ECCV 2018*, pages 467–483, Cham, 2018. Springer International Publishing.
- 1037 [51] Aäron Van Den Oord, Nal Kalchbrenner, and Koray Kavukcuoglu. Pixel recurrent neural  
1038 networks. In *Proceedings of the 33rd International Conference on International Conference*  
1039 *on Machine Learning - Volume 48, ICML’16*, pages 1747–1756. JMLR.org, 2016.
- 1040 [52] Mehdi Mirza and Simon Osindero. Conditional generative adversarial nets. *CoRR*,  
1041 abs/1411.1784, 2014.
- 1042 [53] Emily L Denton, Soumith Chintala, arthur szlam, and Rob Fergus. Deep generative image  
1043 models using a laplacian pyramid of adversarial networks. In C. Cortes, N. D. Lawrence,  
1044 D. D. Lee, M. Sugiyama, and R. Garnett, editors, *Advances in Neural Information Processing*  
1045 *Systems 28*, pages 1486–1494. Curran Associates, Inc., 2015.
- 1046 [54] Ishaan Gulrajani, Faruk Ahmed, Martin Arjovsky, Vincent Dumoulin, and Aaron C Courville.  
1047 Improved training of wasserstein gans. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wal-  
1048 lach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information*  
1049 *Processing Systems 30*, pages 5767–5777. Curran Associates, Inc., 2017.
- 1050 [55] Anders Boesen Lindbo Larsen, Søren Kaae Sønderby, Hugo Larochelle, and Ole Winther.  
1051 Autoencoding beyond pixels using a learned similarity metric. In *Proceedings of the 33rd*  
1052 *International Conference on International Conference on Machine Learning - Volume 48,*  
1053 *ICML’16*, pages 1558–1566. JMLR.org, 2016.



- 1054 [56] A. Van Den Oord, S. Dieleman, H. Zen, K. Simonyan, O. Vinyals, A. Graves, N. Kalchbrenner,  
1055 A.W. Senior, and K. Kavukcuoglu. Wavenet: A generative model for raw audio. In *SSW*, page  
1056 125, 2016.
- 1057 [57] Aaron van den Oord, Nal Kalchbrenner, Lasse Espeholt, koray kavukcuoglu, Oriol Vinyals, and  
1058 Alex Graves. Conditional image generation with pixelcnn decoders. In D. D. Lee, M. Sugiyama,  
1059 U. V. Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information Processing*  
1060 *Systems* 29, pages 4790–4798. Curran Associates, Inc., 2016.
- 1061 [58] J. E. Kyprianidis, J. Collomosse, T. Wang, and T. Isenberg. State of the "art": A taxonomy of  
1062 artistic stylization techniques for images and video. *IEEE Transactions on Visualization and*  
1063 *Computer Graphics*, 19(5):866–885, May 2013.
- 1064 [59] M. Elad and P. Milanfar. Style transfer via texture synthesis. *IEEE Transactions on Image*  
1065 *Processing*, 26(5):2338–2351, May 2017.
- 1066 [60] Jun-Yan Zhu, Philipp Krähenbühl, Eli Shechtman, and Alexei A. Efros. Generative visual  
1067 manipulation on the natural image manifold. In Bastian Leibe, Jiri Matas, Nicu Sebe, and  
1068 Max Welling, editors, *Computer Vision – ECCV 2016*, pages 597–613, Cham, 2016. Springer  
1069 International Publishing.
- 1070 [61] Yijun Li, Ming-Yu Liu, Xueting Li, Ming-Hsuan Yang, and Jan Kautz. A closed-form solution  
1071 to photorealistic image stylization. In *The European Conference on Computer Vision (ECCV)*,  
1072 September 2018.
- 1073 [62] L.A. Gatys, A.S. Ecker, and M. Bethge. Image style transfer using convolutional neural  
1074 networks. In *CVPR*, pages 2414–2423. IEEE, 2016.
- 1075 [63] Christian Ledig, Lucas Theis, Ferenc Huszar, Jose Caballero, Andrew Cunningham, Alejandro  
1076 Acosta, Andrew Aitken, Alykhan Tejani, Johannes Totz, Zehan Wang, and Wenzhe Shi. Photo-  
1077 realistic single image super-resolution using a generative adversarial network. In *The IEEE*  
1078 *Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017.
- 1079 [64] Richard Zhang, Phillip Isola, and Alexei A. Efros. Colorful image colorization. In Bastian  
1080 Leibe, Jiri Matas, Nicu Sebe, and Max Welling, editors, *Computer Vision – ECCV 2016*, pages  
1081 649–666, Cham, 2016. Springer International Publishing.
- 1082 [65] Takuhiro Kaneko, Kaoru Hiramatsu, and Kunio Kashino. Generative attribute controller with  
1083 conditional filtered generative adversarial networks. In *The IEEE Conference on Computer*  
1084 *Vision and Pattern Recognition (CVPR)*, July 2017.
- 1085 [66] Yijun Li, Chen Fang, Jimei Yang, Zhaowen Wang, Xin Lu, and Ming-Hsuan Yang. Universal  
1086 style transfer via feature transforms. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fer-  
1087 gus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing*  
1088 *Systems* 30, pages 386–396. Curran Associates, Inc., 2017.
- 1089 [67] Ting-Chun Wang, Ming-Yu Liu, Jun-Yan Zhu, Andrew Tao, Jan Kautz, and Bryan Catanzaro.  
1090 High-resolution image synthesis and semantic manipulation with conditional gans. *CoRR*,  
1091 abs/1711.11585, 2017.
- 1092 [68] Levent Karacan, Zeynep Akata, Aykut Erdem, and Erkut Erdem. Learning to generate images  
1093 of outdoor scenes from attributes and semantic layouts. *CoRR*, abs/1612.00215, 2016.
- 1094 [69] Zhifei Zhang, Yang Song, and Hairong Qi. Age progression/regression by conditional adver-  
1095 sarial autoencoder. In *The IEEE Conference on Computer Vision and Pattern Recognition*  
1096 *(CVPR)*, July 2017.
- 1097 [70] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A. Efros. Image-to-image translation  
1098 with conditional adversarial networks. In *The IEEE Conference on Computer Vision and*  
1099 *Pattern Recognition (CVPR)*, July 2017.
- 1100 [71] Justin Johnson, Alexandre Alahi, and Li Fei-Fei. Perceptual losses for real-time style transfer  
1101 and super-resolution. In Bastian Leibe, Jiri Matas, Nicu Sebe, and Max Welling, editors, *Com-*  
1102 *puter Vision – ECCV 2016*, pages 694–711, Cham, 2016. Springer International Publishing.

- 1103 [72] Xiaolong Wang and Abhinav Gupta. Generative image modeling using style and structure  
1104 adversarial networks. In Bastian Leibe, Jiri Matas, Nicu Sebe, and Max Welling, editors, *Com-  
1105 puter Vision – ECCV 2016*, pages 318–335, Cham, 2016. Springer International Publishing.
- 1106 [73] Mingsheng Long, Yue Cao, Jianmin Wang, and Michael I. Jordan. Learning transferable  
1107 features with deep adaptation networks. In *Proceedings of the 32Nd International Conference  
1108 on International Conference on Machine Learning - Volume 37*, ICML’15, pages 97–105.  
1109 JMLR.org, 2015.
- 1110 [74] Hemanth Venkateswara, Jose Eusebio, Shayok Chakraborty, and Sethuraman Panchanathan.  
1111 Deep hashing network for unsupervised domain adaptation. In *The IEEE Conference on  
1112 Computer Vision and Pattern Recognition (CVPR)*, July 2017.
- 1113 [75] N. Courty, R. Flamary, D. Tuia, and A. Rakotomamonjy. Optimal transport for domain  
1114 adaptation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 39(9):1853–  
1115 1865, Sept 2017.
- 1116 [76] Nicolas Courty, Rémi Flamary, Amaury Habrard, and Alain Rakotomamonjy. Joint distribution  
1117 optimal transportation for domain adaptation. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wal-  
1118 lach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information  
1119 Processing Systems 30*, pages 3730–3739. Curran Associates, Inc., 2017.
- 1120 [77] L. Sun, K. Li, H. Wang, S. Kang, and H. Meng. Phonetic posteriorgrams for many-to-one  
1121 voice conversion without parallel data training. In *2016 IEEE International Conference on  
1122 Multimedia and Expo (ICME)*, pages 1–6, July 2016.
- 1123 [78] M. Abe, S. Nakamura, K. Shikano, and H. Kuwabara. Voice conversion through vector  
1124 quantization. In *ICASSP-88., International Conference on Acoustics, Speech, and Signal  
1125 Processing*, pages 655–658 vol.1, April 1988.
- 1126 [79] Y. Stylianou, O. Cappe, and E. Moulines. Continuous probabilistic transform for voice  
1127 conversion. *IEEE Transactions on Speech and Audio Processing*, 6(2):131–142, March 1998.
- 1128 [80] Toru Nakashika, Tetsuya Takiguchi, and Yasuo Ariki. High-order sequence modeling using  
1129 speaker-dependent recurrent temporal restricted boltzmann machines for voice conversion. In  
1130 *INTERSPEECH 2014, 15th Annual Conference of the International Speech Communication  
1131 Association, Singapore, September 14-18, 2014*, pages 2278–2282, 2014.
- 1132 [81] Z. Wu, T. Virtanen, E. S. Chng, and H. Li. Exemplar-based sparse representation with residual  
1133 compensation for voice conversion. *IEEE/ACM Transactions on Audio, Speech, and Language  
1134 Processing*, 22(10):1506–1521, Oct 2014.
- 1135 [82] D. Griffin and Jae Lim. Signal estimation from modified short-time fourier transform. In  
1136 *ICASSP ’83. IEEE International Conference on Acoustics, Speech, and Signal Processing*,  
1137 volume 8, pages 804–807, April 1983.
- 1138 [83] Daniel Erro, Iñaki Sainz, Eva Navas, and Inma Hernández. Improved hnm-based vocoder for  
1139 statistical synthesizers. In *INTERSPEECH 2011, 12th Annual Conference of the International  
1140 Speech Communication Association, Florence, Italy, August 27-31, 2011*, pages 1809–1812,  
1141 2011.
- 1142 [84] A. Tamamori, T. Hayashi, K. Kobayashi, K. Takeda, and T. Toda. Speaker-dependent wavenet  
1143 vocoder. In *Proc. Interspeech*, volume 2017, pages 1118–1122, 2017.
- 1144 [85] S.H. Mohammadi and A. Kain. An overview of voice conversion systems. *Speech Communi-  
1145 cation*, 88:65–82, 2017.
- 1146 [86] D. Sundermann, H. Hoge, A. Bonafonte, H. Ney, A. Black, and S. Narayanan. Text-independent  
1147 voice conversion based on unit selection. In *2006 IEEE International Conference on Acoustics  
1148 Speech and Signal Processing Proceedings*, volume 1, pages I–I, May 2006.
- 1149 [87] Feng-Long Xie, Frank K. Soong, and Haifeng Li. A kl divergence and dnn-based approach to  
1150 voice conversion without parallel training sentences. In *INTERSPEECH*, 2016.

- 1151 [88] Chin-Cheng Hsu, Hsin-Te Hwang, Yi-Chiao Wu, Yu Tsao, and Hsin-Min Wang. Voice  
1152 conversion from unaligned corpora using variational autoencoding wasserstein generative  
1153 adversarial networks. In *INTERSPEECH*, 2017.
- 1154 [89] Jaime Lorenzo-Trueba, Junichi Yamagishi, Tomoki Toda, Daisuke Saito, Fernando Villavicen-  
1155 cio, Tomi Kinnunen, and Zhen-Hua Ling. The voice conversion challenge 2018: Promoting  
1156 development of parallel and nonparallel methods. *CoRR*, abs/1804.04262, 2018.
- 1157 [90] H. Kawahara, I. Masuda-Katsuse, and A. De Cheveigne. Restructuring speech representations  
1158 using a pitch-adaptive time–frequency smoothing and an instantaneous-frequency-based f0  
1159 extraction: Possible role of a repetitive structure in sounds1. *Speech communication*, 27(3-  
1160 4):187–207, 1999.
- 1161 [91] T. Kaneko and H. Kameoka. Parallel-data-free voice conversion using cycle-consistent adver-  
1162 sarial networks. *arXiv preprint arXiv:1711.11293*, 2017.
- 1163 [92] G. Zhao, S. Sonsaat, J. Levis, E. Chukharev-Hudilainen, and R. Gutierrez-Osuna. Accent  
1164 conversion using phonetic posteriorgrams. In *ICASSP*, pages 5314–5318. IEEE, 2018.
- 1165 [93] Miaomiao Wang, Miaomiao Wen, Keikichi Hirose, and Nobuaki Minematsu. Emotional voice  
1166 conversion for mandarin using tone nucleus model – small corpus and high efficiency. 2012.
- 1167 [94] Z. Wang and Y. Yu. Multi-level prosody and spectrum conversion for emotional speech  
1168 synthesis. In *Signal Processing (ICSP)*, pages 588–593. IEEE, 2014.
- 1169 [95] Y. Xue, Y. Hamada, and M. Akagi. Voice conversion for emotional speech: Rule-based  
1170 synthesis with degree of emotion controllable in dimensional space. *Speech Communication*,  
1171 102:54–67, 2018.
- 1172 [96] H. Fujisaki and K. Hirose. Analysis of voice fundamental frequency contours for declarative  
1173 sentences of japanese. *ASJs Japan (E)*, 5(4):233–242, 1984.
- 1174 [97] Y. Xue and M. Akagi. A study on applying target prediction model to parameterize power  
1175 envelope of emotional speech. In *RISP workshop NCSP’16*, 2016.
- 1176 [98] Y.N. Dauphin, A. Fan, M. Auli, and D. Grangier. Language modeling with gated convolutional  
1177 networks. In *ICML*, pages 933–941, 2017.
- 1178 [99] C. Busso, M. Bulut, CC Lee, A. Kazemzadeh, E. Mower, S. Kim, J.N. Chang, S. Lee, and  
1179 S.S. Narayanan. Iemocap: Interactive emotional dyadic motion capture database. *Language*  
1180 *resources and evaluation*, 42(4):335, 2008.
- 1181 [100] S. Mirsamadi, E. Barsoum, and C. Zhang. Automatic speech emotion recognition using  
1182 recurrent neural networks with local attention. In *ICASSP*, pages 2227–2231. IEEE, 2017.