
PADGAN: A GENERATIVE ADVERSARIAL NETWORK FOR PERFORMANCE AUGMENTED DIVERSE DESIGNS

ACCEPTED BY IDETC/CIE 2020

Wei Chen

Siemens Corporate Technology
Princeton, NJ 08540
chen.wei@siemens.com

Faez Ahmed

Northwestern University
Evanston, IL 10601
faez@northwestern.edu

ABSTRACT

Deep generative models are proven to be a useful tool for automatic design synthesis and design space exploration. When applied in engineering design, existing generative models face three challenges: 1) generated designs lack diversity and do not cover all areas of the design space, 2) it is difficult to explicitly improve the overall performance or quality of generated designs, and 3) existing models generally do not generate novel designs, outside the domain of the training data. In this paper, we simultaneously address these challenges by proposing a new Determinantal Point Processes based loss function for probabilistic modeling of diversity and quality. With this new loss function, we develop a variant of the Generative Adversarial Network, named “Performance Augmented Diverse Generative Adversarial Network” or PaDGAN, which can generate novel high-quality designs with good coverage of the design space. Using three synthetic examples and one real-world airfoil design example, we demonstrate that PaDGAN can generate diverse and high-quality designs. In comparison to a vanilla Generative Adversarial Network, on average, it generates samples with a 28% higher mean quality score with larger diversity and without the mode collapse issue. Unlike typical generative models that usually generate new designs by interpolating within the boundary of training data, we show that PaDGAN expands the design space boundary outside the training data towards high-quality regions. The proposed method is broadly applicable to many tasks including design space exploration, design optimization, and creative solution recommendation.

1 Introduction

A designer wants good design solutions which are creative and meet the performance requirements. By manually and iteratively exploring design ideas using experience and design heuristics, the designers take the risks of 1) wasting time on unfavorable or even invalid design candidates and 2) not exploring as deeply as they might want to. An ideal design space exploration tool should ensure that, with low cost, one can dive deep in the design space and explore all feasible design alternatives.

While recent advances in machine learning assisted automatic design synthesis and design space exploration are promising, the current methods are still far from this ideal picture. To model a design space, researchers have used deep generative models like variational autoencoders (VAEs) [1] and generative adversarial networks (GANs) [2], as they can learn the distribution of existing designs. The hope is that by learning an underlying *latent space*, which can represent most designs, one can automatically synthesize many new designs from the low-dimensional *latent vectors* and design exploration becomes more efficient due to the reduced dimensionality [3, 4, 5]. However, unlike image generation tasks where these generative models are commonly applied, engineering design problems have one or more performance (or quality) measures. The quality measures how well a design achieves its intended goals and is defined based on the specific problem. For example, beam design problems often define quality based on the compliance value (single-objective) [6] or both compliance and natural-frequency (multi-objective) [7]. For aerodynamic design, quality can be defined as the lift-to-drag ratio [5] or the inverse of the drag coefficient [8]. Current state-of-the-art generative models have no mechanism of explicitly promoting high-quality design generation. One may spend huge effort to train a generative model, only to find many generated designs are infeasible or do not meet design requirements. One way of working around this problem is to exclude low-quality data while training [8]. However, such an approach

may affect model performance due to reduced training sample size. This creates a need to explicitly embed the quality measurement into a generative model, so that it can learn to generate high-quality designs by making use of full data and their quality measurements.

In this work, we focus on addressing the problem of simultaneously maximizing diversity and quality of generated designs. Specifically, we develop a new loss function, based on Determinantal Point Processes (DPPs) [9], for generative models to encourage both high-quality and diverse design synthesis. Using this loss function, we develop a new variant of GAN, named PaDGAN. We show that it can generate high-quality new samples with a good coverage of the design space. More importantly, we found that PaDGAN can expand the existing boundary of the design space towards high-quality regions, which indicates its ability of generating novel high-quality designs.

With the ability of generating high-quality and diverse designs from a (reduced) latent representation, the proposed PaDGAN can then be used for improving the efficiency in design space exploration. While it is interesting to see how exploring the low-dimensional latent space of the PaDGAN can accelerate exploration or improve the performance of the optimal solution, we leave that to future work. In this paper, we focus on the architecture of PaDGAN and its performance in design synthesis.

2 Background and Related Work

Our work produces generative models that synthesize diverse designs from latent representations. There are primarily two streams of related research: 1) design synthesis and 2) diversity measurement. Within these two fields, we provide a brief background on two techniques we use in this paper — GANs and DPPs — and their applications in design. Readers interested in a more comprehensive understanding of their background are advised to read Kulesza *et al.*'s work [9] for DPPs and the chapter on “Deep Generative Models” in Ref. [10].

2.1 Deep Generative Model-Based Design Synthesis

To achieve automatic design synthesis, past researchers have used approaches based on shape grammar [11, 12, 13], graph enumeration [14, 15], functional models [16], analogy [17], and constraint programming [18, 19]. These methods often need to encode expert knowledge as either grammar rules, functional basis, or constraints. In recent years, data-driven design synthesis has become increasingly popular. Different from traditional design synthesis methods, data-driven methods do not necessarily require expert knowledge and can learn to generate plausible new designs from a database [20, 3, 21, 22, 5].

In the last few years, deep generative models have gained traction, due to their ability to learn complex feature representations. The family of deep generative models contains various methods like the Boltzmann machines, deep belief networks, and differentiable generator networks like VAEs and GANs. VAEs and GANs are the two most commonly used deep generative models for solving engineering design problems. For example, they have been used in applications like design exploration [23, 3, 4], surrogate modeling [24], and material microstructure design [25, 26].

Applications of Deep Generative Models in Design Synthesis. Many design applications have huge collections of unstructured design data (CAD models, images, microstructures, *etc.*) with hundreds of features and multiple functionalities. To learn from these complex datasets, deep generative models have increasingly been employed. For instance, Chen *et al.* [5, 27] proposed a BézierGAN model for airfoil parameterization and synthesis, and demonstrated significantly faster convergence to the optimum when optimizing over the latent space. Yang *et al.* [26] used a GAN to generate microstructures and performed design optimization over the latent space. Chen *et al.* [4] proposed a hierarchical GAN architecture to synthesize designs with inter-part dependencies. Oh *et al.* [28] integrated topology optimization and generative models to generate designs which are optimized for engineering performance. These methods either do not explicitly consider the quality of generated designs or use a separate optimization process to search for high-quality designs. Burnap *et al.* [29] used a VAE to generate new highly-rated automotive images, which are aesthetically pleasing. Shu *et al.* [8] proposed a GAN-based model to generate high-quality 3D designs, where they improve the quality of generated samples by retraining the model on an updated dataset with low performing designs removed. In contrast, our method improves the quality of generated designs while training the deep generative model, without retraining or discarding any samples in the training data. Also, to the best of our knowledge, there is no generative model that simultaneously encourages diversity and quality. While the methods we develop in this work are applicable to most deep generative models, we use GANs to demonstrate our results and will describe them next.

Generative Adversarial Networks. GANs [2] model a game between a generative model (*generator*) and a discriminative model (*discriminator*). The generative model maps an arbitrary noise distribution to the data distribution (*i.e.*, the distribution of designs in our scenario), thus can generate new data; while the discriminative model tries to perform

classification, *i.e.*, to distinguish between real and generated data. The generator G and the discriminator D are usually built with deep neural networks. As D improves its classification ability, G also improves its ability to generate data that fools D . Thus, a vanilla GAN (standard GAN with no bells and whistles) has the following objective function, which comprises of a discriminator loss term and a generator loss term:

$$\min_G \max_D V(D, G) = \mathbb{E}_{\mathbf{x} \sim P_{data}} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim P_z} [\log(1 - D(G(\mathbf{z})))] \quad (1)$$

where \mathbf{x} is sampled from the data distribution P_{data} , \mathbf{z} is sampled from the noise distribution P_z , and $G(\mathbf{z})$ is the generator distribution. A trained generator thus can map from a predefined noise distribution to the distribution of designs. The noise input \mathbf{z} is considered as the latent representation of the data, which can be used for design synthesis and exploration.

Problems in Using GANs for Design Synthesis. Learning in GANs can be difficult in practice, which may be one of the reasons that they are less widely used in design compared to VAEs. Despite an enormous amount of recent work in the machine learning community, GANs are notoriously unstable to train, and it has been observed that they often suffer from *mode collapse* [30], in which the generator network learns how to generate samples from a few modes of the data distribution but misses many other modes. For instance, when training on multiple categories of designs, a GAN model would sometimes generate designs only for a single category [31]. Recent approaches [32, 33, 34] tackled mode collapse in one of two different ways: 1) modifying the learning of the system to reach a better convergence point; or 2) explicitly enforcing the models to capture diverse modes or map back to the true-data distribution. Solutions to the mode collapse problem range from designing a reconstructor network in VEEGAN [33] to matching the similarity matrix of generated samples with data [35]. In contrast, PaDGAN addresses the mode collapse problem implicitly by virtue of promoting generation of diverse solutions, which encourages samples to cover different modes, hence alleviating mode collapse.

2.2 Measuring Design Coverage

Massive highly redundant sources of audio, video, speech, text documents, and sensor data have become commonplace and are expected to become larger and more preponderant in the future [36]. This brings a need to measure diversity of a set of items, such that redundancy in data can be reduced and machine learning models can be trained using data with a smaller sample size and which are not biased in favor of a few classes. *Diversity* (also called coverage or variety) is a measure of how different a set of items are from each other. Quantitatively, it is measured using two predominant ways — submodular functions or DPPs. Submodular functions are set functions with diminishing marginal gain property, which naturally model notions of coverage and diversity. They achieved among the top results on common automatic document summarization benchmarks (*e.g.*, at the Document Understanding Conference [37]). In design, too, researchers have used submodular functions-based diversity measures to understand design space exploration using terms like *variety* [38, 39, 40]. These functions have helped designers sift through large sets of ideas by ranking them [41] or selecting a diverse subset [42]. Ahmed *et al.* [41] compared DPPs [9] with certain commonly used submodular functions. They concluded that unlike submodular functions, DPPs are more flexible, since they only need a valid similarity kernel as an input rather than an underlying Euclidean space or clusters. In this paper, we will use DPPs as a measure of diversity, which we will describe next.

Determinantal Point Processes. DPPs, which arise in quantum physics, are probabilistic models that model the likelihood of selecting a subset of diverse items as the determinant of a kernel matrix. Viewed as joint distributions over the binary variables corresponding to item selection, DPPs essentially capture negative correlations and provide a way to elegantly model the trade-off between often competing notions of quality and diversity. The intuition behind DPPs is that the determinant of a kernel matrix roughly corresponds to the volume spanned by the vectors representing the items. Points that “cover” the space well should capture a larger volume of the overall space, and thus have a higher probability. As shown by Kulesza *et al.* [43], one of DPPs’ advantages is that computing marginals, computing certain conditional probabilities, and sampling can all be done in polynomial time. In this paper, we focus on another advantage of DPPs, which is the decomposition of DPP kernels into quality and similarity terms.

For the purposes of modeling real data, the most relevant construction of DPPs is through L-ensembles [44]. An L-ensemble defines a DPP via a positive semi-definite matrix L indexed by the elements of a subset S . The kernel matrix L defines a global measure of similarity between pairs of items, so that more similar items are less likely to co-occur. The probability of a set S occurring under a DPP is calculated as:

$$\mathbb{P}_L(S) = \frac{\det(L_S)}{\det(L + I)}, \quad (2)$$

where $L_S \equiv [L_{ij}]_{ij \in S}$ denotes the restriction of L to the entries indexed by elements of S , I is an $N \times N$ identity matrix, and N is the total number of items. For any set size, the most probable subset under a DPP will have the maximum likelihood over $\mathbb{P}_L(S)$ or (equivalently) the highest determinant (the denominator can be ignored for maximizing determinant of a fixed set size). Similar to sub-modular functions, one of the main applications of DPP is extractive document summarization, where it provided state-of-the-art results. In Section 3, we show how the decomposition of DPP kernels can be used to design a DPP-based loss function, which promotes quality and diversity of generated samples in a generative model.

2.3 Comparison with State-of-the-Art and Our Contributions

The work closest to ours is the GDPP method [35] by Elfeki *et al.*. The authors devised an objective term that encourages the GAN to synthesize data with diversity similar to the training data. PaDGAN differs from their method in three aspects. First, PaDGAN is stable against scaling of data while on validating GDPP for multiple test problems, we found that their method does not work for problems with training data at different scales. Second, while PaDGAN aims to maximize the diversity of generated samples, GDPP aims to achieve a similar diversity value as the training data. By avoiding the goal of mimicking the diversity of the training data, PaDGAN will generate diverse samples even when the original training dataset is biased in favor of a few modes, while GDPP is designed to mimic the bias in generated samples. Finally, whereas we maximize the quality of generated samples, whereas GDPP does not have such consideration. This feature of PaDGAN is helpful for design exploration as it can help discover novel high-quality designs (demonstrated in Section 5.2).

The scientific contributions and novelty of this work are as follows:

1. We propose a novel design synthesis method that simultaneously encourage synthesis of diverse and high-performance designs.
2. We find that PaDGAN can expand the design space boundary towards high-quality regions that it had not seen from existing data.
3. We propose a way to control the trade-off between quality and diversity in DPPs. Our method extends past work on decomposing a DPP kernel by providing a way to tune the relative importance of quality over diversity.
4. We provide easy-to-verify test cases and metrics to validate any generative models, whose goal is to maximize sample quality and/or coverage over a dataset with multiple modes.

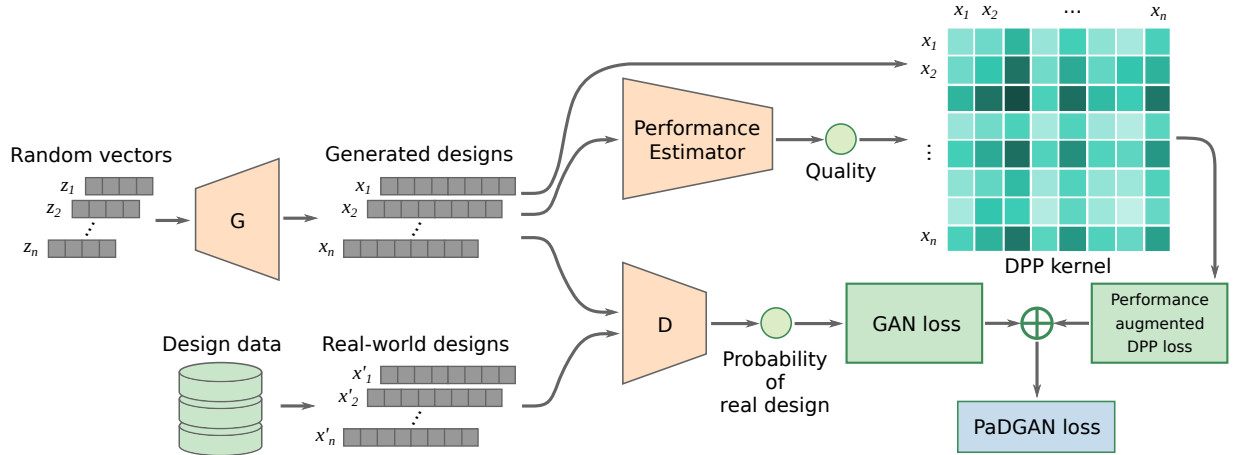


Figure 1: Architecture of PaDGAN.

3 Methodology

Built on a standard GAN architecture, PaDGAN introduces a *performance augmented DPP loss* which measures the diversity and quality of a batch of generated designs during training. The overall model architecture of PaDGAN is shown in Fig. 1. In this section, we begin by describing how to decompose a DPP kernel, then proceed on how to create

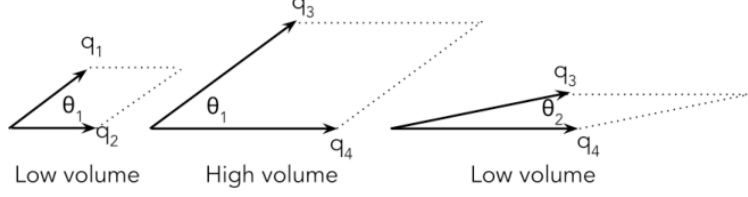


Figure 2: Visualization of volume in 2-D space, where q_i represent the quality of an item and θ_i shows how similar they are. Comparing the leftmost figure to the central figure, we observe that the similarity is same but the quality magnitude increases, leading to a higher volume. From the central figure to the rightmost figure, we observe that the quality magnitude is the same, but the similarity increases, leading to a lower volume (or diversity) encompassed by the two vectors.

a DPP loss which augments high performing designs, and finally provide a method to balance diversity and quality using a quality dial. We also add a note on improving training stability at the end.

3.1 Decomposition of a DPP kernel

DPP kernels can be decomposed into quality and diversity parts [9]. The $(i, j)^{th}$ entry of a positive semi-definite DPP kernel L can be expressed as:

$$L_{ij} = q_i \phi(i)^T \phi(j) q_j. \quad (3)$$

We can think of $q_i \in R^+$ as a scalar value measuring the quality (or performance) of an item i , and $\phi(i)^T \phi(j)$ as a signed measure of similarity between items i and j . The decomposition enforces L to be positive semidefinite. Suppose we select a subset S of samples, then this decomposition allows us to write the probability of this subset S as the square of the volume spanned by $q_i \phi_i$ for $i \in S$ using the equation below:

$$\mathbb{P}_L(S) \propto \prod_{i \in S} (q_i^2) \det(K_S), \quad (4)$$

where K_S is the similarity matrix of S .

The first term increases with the quality of the selected items, and the second term increases with the diversity of the selected items. As item i 's quality q_i increases, so do the probabilities of sets containing item i . As two items i and j become more similar, $\phi_i^T \phi_j$ increases and the probabilities of sets containing both i and j decrease. From a geometric intuition, the determinant of L_Y is equal to the squared volume of the parallelepiped spanned by the vectors $q_i \phi_i$ for $i \in Y$. We show an illustration of this intuition in Fig. 2. The magnitude of the vector representing item i is q_i , and its direction is ϕ_i . It shows how DPPs decomposed into quality and diversity naturally balance the two objectives of high-quality and high diversity.

When selecting a subset S of items, without the diversity term, we would choose high-quality items, but we would tend to choose similar high-quality items over and over. Without the quality term, we would get a very diverse set, but we might fail to include the most important items in S , focusing instead on low-quality outliers. By combining the two models, we can achieve a more balanced result. The key intuition of PaDGAN is that if we can find a way to add the term from Eq. (4) to the objective function of any generative model, then while training it will be encouraged to generate high probability subsets, which will be both diverse and high-quality. In the next section, we define such a loss function.

While, the authors used this decomposition to find quality and similarity terms from a known kernel, we reverse this procedure to create the kernel L for a sample of points generated by PaDGAN from known inter-sample similarity values and quality. Note that in a DPP model, the quality or performance of an item is a scalar value, like compliance, displacement, drag-coefficient, *etc.* The quality can be estimated using an external model (like a physics-based simulator) or by finding the distance of current performance of a design from a target performance. For multi-dimensional cases, quality can be derived by taking the norm of multiple dimensions. The similarity terms $\phi(i)^T \phi(j)$ can be derived using any similarity kernel, which we represent using $k(\mathbf{x}_i, \mathbf{x}_j) = \phi(i)^T \phi(j)$ and $\|\phi(i)\| = \|\phi(j)\| = 1$. Here \mathbf{x}_i is a vector representation of a design.

3.2 Performance Augmented DPP Loss

Our performance augmented DPP loss models diversity and quality simultaneously and gives a lower loss to sets of designs which are both high-quality and diverse. Specifically, we construct a kernel matrix L_B for a generated batch B based on Eq. (3). For each entry of L_B , we have

$$L_B(i, j) = k(\mathbf{x}_i, \mathbf{x}_j) (q(\mathbf{x}_i)q(\mathbf{x}_j))^{\gamma_0}, \quad (5)$$

where $\mathbf{x}_i, \mathbf{x}_j \in B$, $q(\mathbf{x})$ is the quality function at \mathbf{x} , and $k(\mathbf{x}_i, \mathbf{x}_j)$ is the similarity kernel between \mathbf{x}_i and \mathbf{x}_j . We add γ_0 term as a dial to control the weight of quality, which is further explained in Section 3.3.

The performance augmented DPP loss is expressed as

$$\mathcal{L}_{\text{PaD}}(G) = -\frac{1}{|B|} \log \det(L_B) = -\frac{1}{|B|} \sum_{i=1}^{|B|} \log \lambda_i \quad (6)$$

where λ_i is the i -th eigenvalue of L_B . By adding this loss to the vanilla GAN's objective from Eq. (1), the problem becomes:

$$\min_G \max_D V(D, G) + \gamma_1 \mathcal{L}_{\text{PaD}}(G), \quad (7)$$

where γ_1 controls the weight of $\mathcal{L}_{\text{PaD}}(G)$. To update any weight θ_G^i in the generator in terms of $\mathcal{L}_{\text{PaD}}(G)$, we descend its gradient based on the chain rule:

$$\frac{\partial \mathcal{L}_{\text{PaD}}(G)}{\partial \theta_G^i} = \sum_{j=1}^{|B|} \left(\frac{\partial \mathcal{L}_{\text{PaD}}(G)}{\partial q(\mathbf{x}_j)} \frac{dq(\mathbf{x}_j)}{d\mathbf{x}_j} + \frac{\partial \mathcal{L}_{\text{PaD}}(G)}{\partial \mathbf{x}_j} \right) \frac{\partial \mathbf{x}_j}{\partial \theta_G^i}, \quad (8)$$

where $\mathbf{x}_j = G(\mathbf{z}_j)$.

Equation (8) indicates a need for $dq(\mathbf{x})/d\mathbf{x}$, which is the gradient of the quality function. In practice, this gradient is accessible when the quality is evaluated through any performance estimator that is differentiable, like adjoint-based solver methods. If the gradient of a performance estimator is not available, one can either use numerical differentiation or approximate the quality function using a differentiable surrogate model (*e.g.*, a neural network-based surrogate model). In our experiments in Section 5.2, we use a neural network-based surrogate model. We will explore the possibility of using an automatic differentiation enabled simulator (*e.g.*, an adjoint solver) as the performance estimator in future studies.

3.3 Introducing a quality dial for DPP kernels

Note that we modified the original objective to introduce γ_0 as a parameter. We found that traditional DPP decomposition does not allow us to change the importance of quality versus diversity within a given kernel. This means that if we fix the quality scores and similarity scores, the trade-off between the two cannot be controlled. A naïve way to increase the importance of quality would be to multiply the quality scores by a large constant and expect it to increase its importance relative to diversity. However, with careful observation one would realize that this approach would not work. Using the geometric interpretation of the DPPs, this would be equivalent to scaling all lengths by the same factor, which will not affect the volumes relative value. As quality and diversity objectives are multiplied together to get the probability of the set (Eq. (4)), to change the relative importance, we need to adjust the dynamic range of the quality scores. We do this by using an exponent to change the distribution of quality. When $\gamma_0 = 0$, all quality scores collapse to one and the resultant PaDGAN model only generates diverse designs. In contrast, for large values of γ_0 , the highest quality scores have the largest probability mass and PaDGAN only generates the highest quality designs, ignoring diversity. This method of balancing diversity and quality provides more flexibility to PaDGAN and in general, can be used for many applications of DPPs.

3.4 Improving PaDGAN stability

Stabilization of GAN learning remains an open problem and in this section, we provide a heuristic method to improve GAN stability, when using a surrogate model for evaluating quality. Note that in Eq. (8), the quality gradient is used in the back propagation step. If the quality gradients are not accurate, the generator learning can go astray. This is not a problem when the quality estimator is a simulator that can reasonably evaluate (even with low-fidelity) any design in the design space, irrespective of the designs being invalid or unrealistic. However, it creates problems when we use a surrogate model. A surrogate model is normally trained only on realistic designs and hence may perform unreliably on unrealistic ones. In the initial stages of training, a GAN model will not always generate realistic designs during training.

This makes it difficult for the surrogate model to correctly guide the generator’s update and may cause stability issues. To avoid this problem, we propose two small modifications to PaDGAN:

1. *Realistic quality weighted quality*. Specifically, we weight the predicted quality at \mathbf{x} by the probability of \mathbf{x} being the real design (predicted by the discriminator):

$$q(\mathbf{x}) = D(\mathbf{x})q'(\mathbf{x})$$

where $q'(\mathbf{x})$ is the predicted quality (by a surrogate model for example), and $D(x)$ is the discriminator’s output at \mathbf{x} .

2. An *escalating schedule* for setting γ_1 (the weight of DPP loss). A GAN is more likely to generate unrealistic designs in its early stage of training. Thus, we initialize γ_1 at 0 and increase it during training, so that PaDGAN focuses on learning to generate realistic designs at the early stage, and takes quality into consideration later when the generator can produce more realistic designs. The schedule is set as:

$$\gamma_1 = \gamma'_1 \left(\frac{t}{T} \right)^p$$

where γ'_1 is the value of γ_1 at the end of training, t is the current training step, T is the total number of training steps, and p is a factor controlling the steepness of the escalation.

We can also consider the uncertainty of the quality estimation and put a lower weight on the quality score when the uncertainty is high. However, we only consider the above two modifications in this paper and leave others to future work. Note that these modifications are only needed if one is using a performance estimator (*e.g.*, a surrogate model) which gives unreliable quality predictions for unrealistic designs.

4 Experiment

So far, we have shown how the mathematical components of PaDGAN will encourage it to generate high-quality and diverse samples. In this section, we will describe experiments, which can help us validate our claims. These experiments are carefully designed such that the outcome of any generative models can be verified easily. This section introduces the experimental settings for each example. To show the merit of modeling quality and diversity simultaneously, we compare the PaDGAN with alternative models where those two attributes are modeled separately. In the following sections, we show that for three multi-modal synthetic problems, PaDGAN outperforms all other methods by achieving both high-quality and high diversity. Finally, after showing that the claims hold on three test cases, we apply PaDGAN on a real-world airfoil synthesis problem. We find that PaDGAN can discover new regions of high-quality designs, which are outside the design domain over which it was trained.

4.1 Data and Quality Measure

Synthetic example I. The purpose of creating 2D synthetic examples is to test the performance of PaDGAN given known ground truth and visualize the results in terms of diversity and quality. These examples are analogical to any 2D design problem, where designs are represented by two variables. In this synthetic example I, we generate a ring-shaped dataset, with data uniformly distributed between two origin-centered circles of 0.25 and 0.5 in radius, respectively (Fig. 3). We use a density function of an unnormalized Gaussian mixture as the quality function:

$$q(\mathbf{x}) = \sum_{k=1}^K \exp \left(-\frac{(\mathbf{x} - \mu_k)^T (\mathbf{x} - \mu_k)}{2\sigma^2} \right), \quad (9)$$

where μ_k is the mode of the k -th mixture component and σ is the standard deviation. The centers μ_1, \dots, μ_K are evenly spaced around a circle centered at the origin and with a radius of 0.4. We set $K = 6$ and $\sigma \approx 0.1$. Hence, there are six peaks of quality and points are evenly spread between two concentric circles in the training data. Ideally, by simultaneously maximizing diversity and quality, we expect generating more samples near the six local optima (*i.e.*, modes) of the quality function, and those samples should be spread out and evenly distributed among all six mixture components.

Synthetic example II. The data in this example are nine clusters placed on a 3×3 grid (Fig. 3). Similar to synthetic example I, we use Eq. (9) as the quality function. Here we set $K = 4$ and $\sigma \approx 0.16$. Four out of nine clusters (modes) of the data overlap with local optima of the quality function. We expect that if both diversity and quality are considered, the generator should produce most samples in all the four high-quality clusters and few samples in other clusters (instead of generating most samples from a single high-quality cluster).

Synthetic example III. This example is the same as example I, except that data is bounded within two origin-centered circles of 0.325 and 0.375 in radius, respectively (Fig. 3). The purpose of decreasing the coverage of data is to demonstrate PaDGAN’s capability of extrapolating in the high-quality regions (*i.e.*, expanding the boundary of existing design space towards high-quality regions).

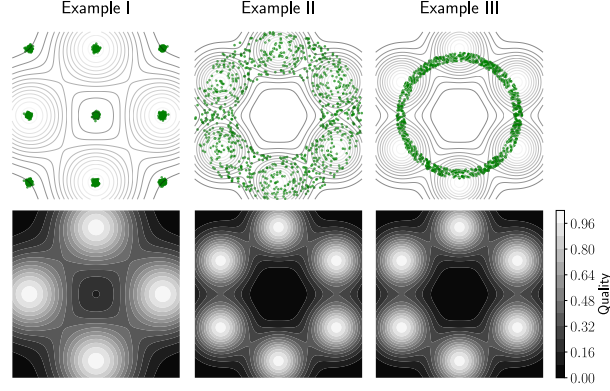


Figure 3: Data and quality functions in synthetic examples. The green dots in the top plots represent data and the contours represent quality functions.

Airfoil example. An airfoil is the cross-sectional shape of a wing or a propeller/rotor/turbine blade. In this example, we use the UIUC airfoil database¹ as our data source. It provides the geometries of nearly 1,600 real-world airfoil designs. Each design is represented by discrete 2D coordinates along their upper and lower surfaces. We preprocessed and augmented the dataset based on Ref. [5] to generate a dataset of 38,802 airfoils. The lift to drag ratio C_L/C_D is a common objective in aerodynamic design optimization problems. Thus we used C_L/C_D as the performance measure, which can be computed using XFOIL software [45]. To provide the gradient of the quality function for Eq. (8), we trained a neural network-based surrogate model on all 38,802 airfoils to approximate the quality. Note that for all the examples, we scaled the quality scores between 0 and 1. We show a subset of 100 randomly chosen example airfoils from the training data in the left plot of Fig. 10.

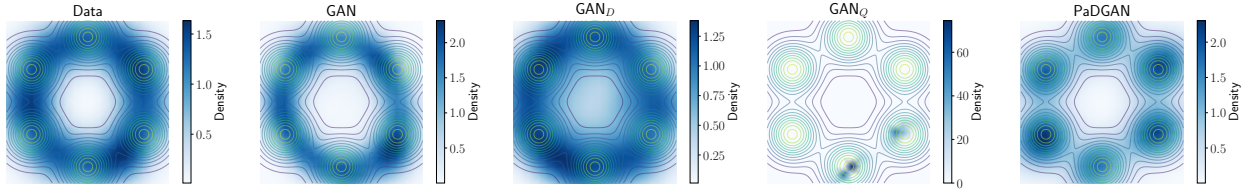


Figure 4: Results on synthetic example I. The leftmost plot show training data (green dots) and quality functions (contour plots). The rest of the plots show the density of samples generated by different models.

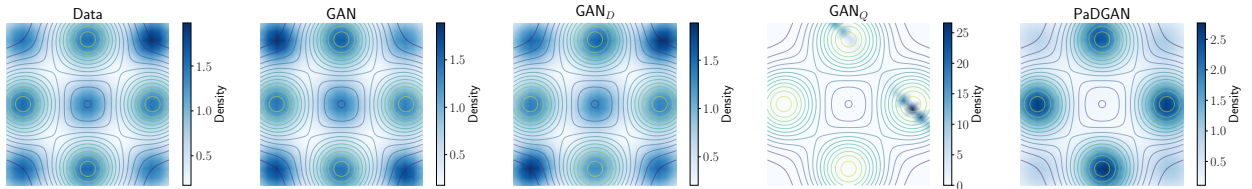


Figure 5: Results on synthetic example II. The leftmost plot show training data (green dots) and quality functions (contour plots). The rest of the plots show the density of samples generated by different models.

¹http://m-selig.ae.illinois.edu/ads/coord_database.html

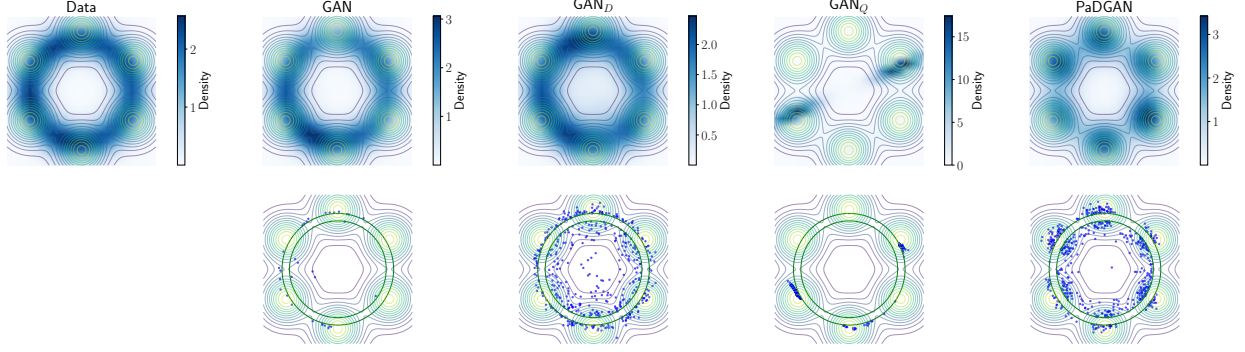


Figure 6: Results on synthetic example III. The leftmost plot show training data (green dots) and quality functions (contour plots). The remaining plots in the top row show the density of samples generated by different models. We observe that GAN and GAN_D generate samples similar to the data, while GAN_Q suffers with mode collapse and only generates samples for two clusters. PaDGAN generates more samples in the high quality region, ignoring the low-quality areas of the training data. Plots in the bottom row show only the generated samples (blue dots) which are “outside” the region of training data (indicated by two green circles). We observe that PaDGAN generates unseen data in the high-quality areas while other methods either do not explore outside the domain or generate low-quality samples.

4.2 Model Configuration and Training

To demonstrate the effectiveness of the PaDGAN, we compare it with the following three models:

1. GAN: a vanilla GAN with the objective of Eq. (1).
2. GAN_D : PaDGAN with $\gamma_0 = 0$ in Eq. (5), *i.e.*, which only optimizes for diversity and ignores the quality.
3. GAN_Q : a vanilla GAN which ignores diversity and only optimizes for quality using the following additional term $\mathcal{L}_Q(G) = -\frac{1}{|B|} \sum_{i=1}^{|B|} q(\mathbf{x}_i)$. The training objective is then set to:

$$\min_G \max_D V(D, G) + \gamma_2 \mathcal{L}_Q(G)$$

where γ_2 controls the weight of the quality objective.

To find similarity between designs, we use a RBF kernel with a bandwidth of 1.0 when constructing L_B in Eq. (5), *i.e.*, $k(\mathbf{x}_i, \mathbf{x}_j) = \exp(-0.5\|\mathbf{x}_i - \mathbf{x}_j\|^2)$. This gives a value between 0 to 1, with a higher value for more similar designs. In synthetic examples, we set $\gamma_0 = 2$ and $\gamma_1 = 0.5$ for PaDGAN and $\gamma_2 = 10$ for GAN_Q (these settings were chosen based on a few initial trials and domain knowledge). The generators and discriminators are fully connected neural networks. In the airfoil example, we set $\gamma_0 = 2$ and $\gamma_1 = 0.2$ for PaDGAN. We used a residual neural network (ResNet) [46] as the surrogate model and a BézierGAN [5, 27] to generate airfoils. For simplicity, we refer to the BézierGAN as a vanilla GAN and the BézierGAN with loss \mathcal{L}_{PaD} as a PaDGAN in the airfoil example in the rest of the paper. Detailed network architecture and hyperparameter settings can be found in our open-source code ².

4.3 Evaluation

We use the *diversity score* and the *quality score* of generated samples to measure the performance of generative models. The diversity score is expressed as the mean log determinant of the similarity matrix:

$$s_{\text{div}} = \frac{1}{n} \sum_{i=0}^n \log \det(L_{S_i}), \quad (10)$$

where n is the number of times diversity is evaluated, $S_i \subseteq Y$ is a random subset of Y (the set of generated samples), and L_{S_i} is the similarity matrix of S_i with entries $L_{S_i}(j, k) = k(\mathbf{x}_j, \mathbf{x}_k)$ for each $\mathbf{x}_j, \mathbf{x}_k \in S_i$. The quality score is

²<https://github.com/wchen459/PaDGAN>

computed by taking the average quality of generated samples:

$$s_{\text{qa}} = \frac{1}{|Y|} \sum_{i=0}^{|Y|} q(\mathbf{x}_i), \quad (11)$$

where $\mathbf{x}_i \in Y$ is a randomly generated design.

For synthetic examples, we define the *overall score*, to measure the overall performance by combining measures for diversity and quality of generated samples:

$$s_{\text{overall}} = - \sum_k \frac{m_k}{|Y|} \log \left(\frac{m_k}{|Y|} \right), \quad (12)$$

where m_k is the number of generated samples within the one-sigma interval of the k -th mixture component of the quality function. The overall score is affected by both the amount of high-quality samples and the spread of those samples. The highest score occurs when there are the same number of generated samples within the one-sigma interval of each mixture component and no samples are outside those intervals.

In the experiments, we set $|Y| = 1000$, $|S_i| = 10$, and $n = 1000$. To take into consideration the stochasticity of the model training, for each type of model (PaDGAN, GAN, GAN_D , and GAN_Q), we train them ten times for each experimental setting, and report the performance statistics for all those ten models (Figs. 7, 8, and 12). We report and discuss the results in the next section.

5 Results and Discussion

In this section, we compare the performance of PaDGAN to its alternatives (*i.e.*, GAN, GAN_D , and GAN_Q) and discuss the implication of these results.

5.1 Synthetic Examples

Figures 4, 5, and 6 show the density plots of generated samples for each model, which represents their *generative distribution*. Ideally, when we sample designs from the generator, we want these designs to have a good coverage over real-world designs (*i.e.*, the training data) and most of them should have high-quality. In Fig. 4, the generative distribution learned by a vanilla GAN fails to cover the entire training data (non-uniform contours). However, in both examples, the generative distribution of GAN_D has a good coverage of the training data due to its diversity objective. This shows that the diversity objective by itself is capable of avoiding mode collapse. By replacing the diversity objective with a quality objective, GAN_Q only generates samples near one of the optima of the quality functions, ignoring the others. In practice, this will give many high-quality samples but they all look very similar to each other. In contrast, the generative distribution of PaDGAN exhibits has a higher density near high-quality regions and also good coverage of the design space.

Figure 6 shows that both GAN_D and PaDGAN expands the boundary of training data. Particularly, PaDGAN expands the boundary towards high-quality regions. If these samples represent designs, it basically indicates that PaDGAN can expand the boundary of existing designs. We will further demonstrate this with a real design problem later. This promising result shows that by diversifying generated samples, PaDGAN is capable of expanding the design space towards the direction of high-quality regions. Note that this is not only filling the “holes” of the design space by interpolation, but also *extrapolation* on the right direction. It is not surprising that the generator knows which direction to expand since it receives from the performance estimator the information of quality gradients.

Figures 7, 8, and 9 show the statistics of ten trained models for each method. Both figures tell that GAN_D has the best performance in the diversity score and the worst performance in the quality score. GAN_Q generates the highest quality samples, but has the lowest diversity scores, showing that all the samples very similar to each other. PaDGAN has the highest overall score in both examples, which shows that it generates high-quality samples that spread over different optima. The lowest variance indicates a consistent performance over multiple runs of PaDGAN training.

5.2 Airfoil Example

We synthesized 100 airfoil designs from a vanilla GAN and 100 from a PaDGAN, computed their quality (C_L/C_D values) using XFOIL³, and used the t-Distributed Stochastic Neighbor Embedding (t-SNE) to map these designs onto

³We set $C_L/C_D = 0$ when the simulation fails.

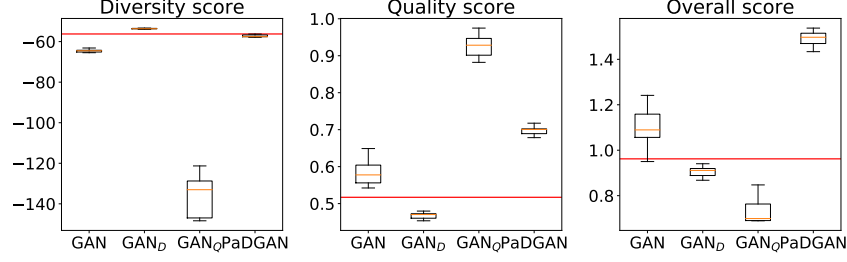


Figure 7: Scores for synthetic example I. The red horizontal line denotes the diversity/quality score of the training data. The box plots show the statistics of ten models for each method.

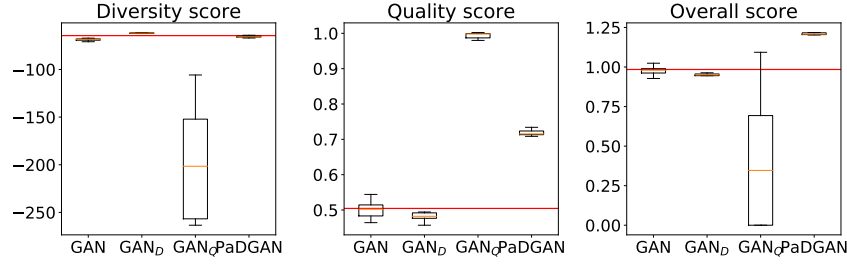


Figure 8: Scores for synthetic example II. The red horizontal line denotes the diversity/quality score of the training data. The box plots show the statistics of ten models for each method.

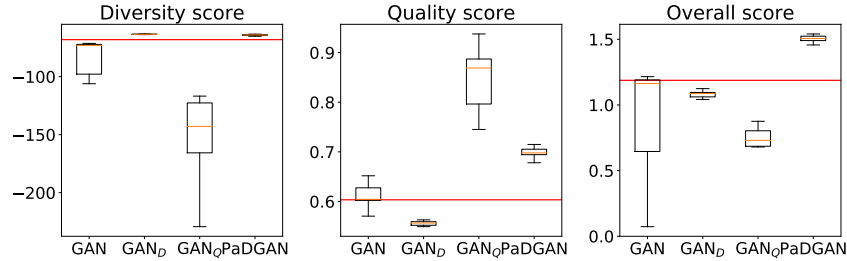


Figure 9: Scores for synthetic example III. The red horizontal line denotes the diversity/quality score of the training data. The box plots show the statistics of ten models for each method.

the same two-dimensional space, as shown in Fig. 10. The quality is indicated by the shades of plotted designs, where dark shaded airfoils are of higher quality. We also show 100 designs from the training data in the left most figure to represent the original design space. Both the GAN and the PaDGAN generate realistic airfoil designs. We observe that the vanilla GAN (middle figure) generates a few airfoils that fill in the gaps of the training data (*i.e.*, interpolation). However, PaDGAN discovers new high-quality designs, which are outside the boundary of the training data. We mark these regions in by ellipses in the leftmost part of Fig. 10. This shows that the diversity promoting part of PaDGAN encourages it to discover new unseen design areas while the quality promoting part helps it find areas where high-quality designs are found, as is also demonstrated by synthetic example III. In future work, we will explore if PaDGAN can be used as a tool to assist in design discovery by generating novel high-quality designs for more complex design domains.

We show the quality (*i.e.*, C_L/C_D) distributions of training data and generated designs by vanilla GAN and PaDGAN in Fig. 11. We observe that the quality distribution of data has two modes (large number of samples) — one near 0 and one near 70. The vanilla GAN’s quality distribution mimics these two modes but has a larger probability mass near 0. Comparing with both the training data and the vanilla GAN, PaDGAN’s quality distribution has a larger mass over the higher-quality region. This shows that PaDGAN generates most samples which are of significantly higher quality than the training data.

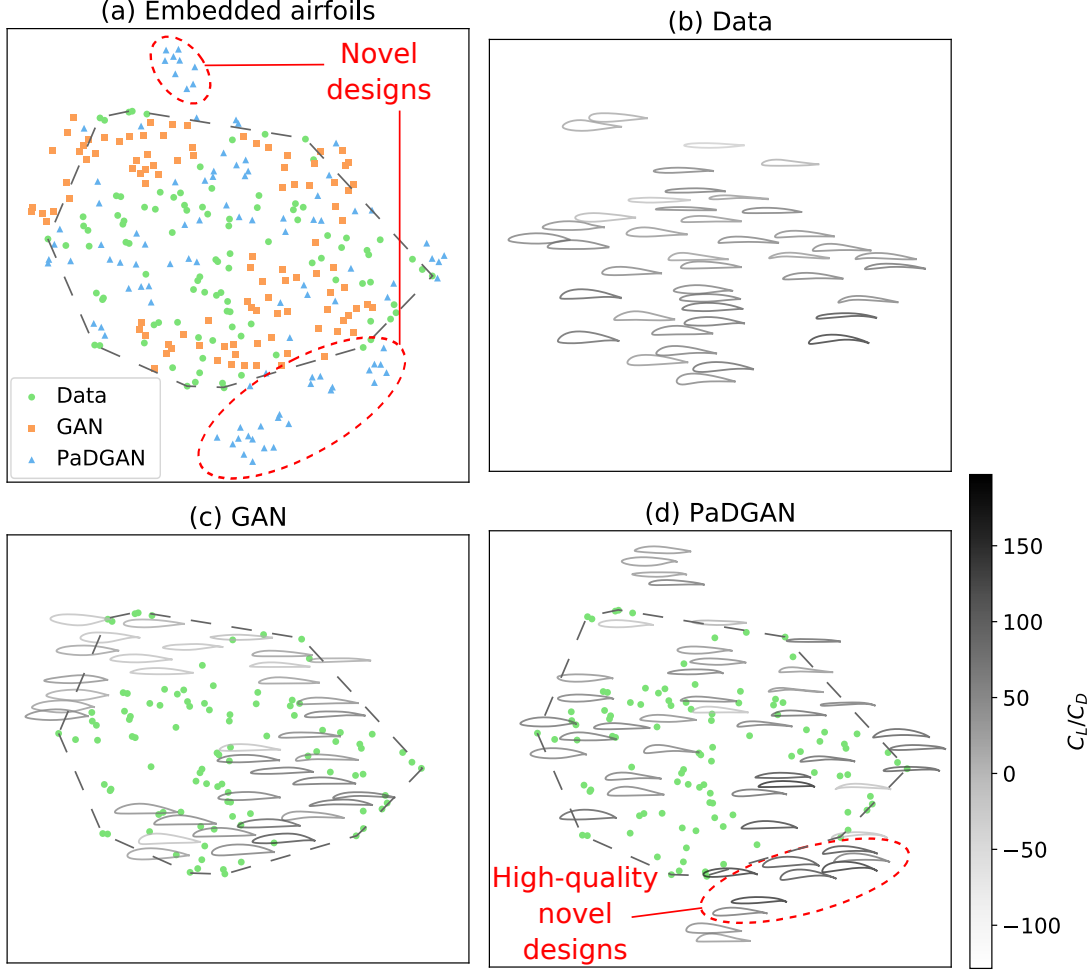


Figure 10: To compare the distribution of real and synthetic airfoils, we map airfoil designs sampled randomly from training data, vanilla GAN, and PaDGAN through t-SNE into the same 2D space (shown in (a)). Plots (b)-(d) visualizes the airfoil geometries, where the shades represent quality (*i.e.*, C_L/C_D). The dots in (c) and (d) represent training data. We label the convex hull of the sampled training data in Plots (a), (c), and (d), which roughly indicates the boundary of the original design space.

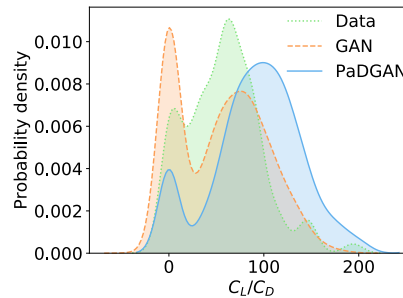


Figure 11: The distribution of quality (C_L/C_D) for training data, vanilla GAN, and PaDGAN.

Figure 12 shows the statistics of quality and diversity scores over ten runs of model training. The PaDGAN's diversity score is always higher than the training data's (shown by a red horizontal line), whereas the vanilla GAN almost always

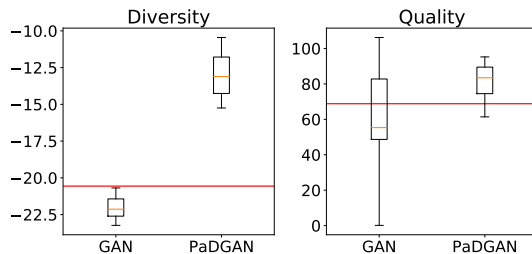


Figure 12: Scores for the airfoil example. The red horizontal line denotes the diversity/quality score of the training data. The box plots show the statistics of ten models for each method.

has a lower diversity score than the data. The quality score of the PaDGAN has a higher mean and lower variance than the vanilla GAN. These results demonstrate the effectiveness of PaDGAN as a design exploration tool.

To show the evaluation scores in Figs. 7-9 and Fig. 12 more clearly, we list the means and 95% confidential intervals of all scores in Appendix A.

6 Conclusion and Future Work

In this paper, we proposed a new loss function for generative models based on Determinantal Point Processes. With this loss function, we developed a new GAN model, named PaDGAN. To the best of authors’ knowledge, this is the first GAN model that can simultaneously encourage the generation of diverse and high-quality designs. We use both synthetic and real-world examples to demonstrate the effectiveness of PaDGAN and show that by diversifying generated samples, PaDGAN expands the existing boundary of the design space towards high-quality regions. This model is particularly useful when we want to thoroughly explore different high-quality design alternatives or discover novel solutions. For example, when performing design optimization, one may accelerate the search for global optimal solutions by sampling start points from the proposed model. Also, this method can be a tool in the early conceptual design stage to aid the creative process. It can generate new designs which are learnt from previous generations of designs, while introducing novelty and taking into account the desired quality metrics. The resultant designs can be used as inspirations to steer designers in exploring novel designs. Although we demonstrated the effectiveness of our method via a GAN-based model, the proposed framework also generalizes to other generative models like variational autoencoders and can be used for various design synthesis problems.

Note that by trying to mimic the training data, PaDGAN captures design constraints implicitly. For instance, in Fig.6 (Example III), it captures the inner and outer ring of the training data and generates the majority of the points inside the two circular rings. However, we still observe a few points outside the rings, as we do not explicitly define this as a constraint boundary. To explicitly capture design constraints, one can train a differentiable classifier (*e.g.*, a neural network-based classifier) which predicts constraint satisfaction and use it as a second discriminator. However, this approach of explicitly capturing the constraints is outside the scope of this work.

While we developed this method for design applications, it can generalize to many other domains, where quality and coverage over a domain are needed. For example, in molecule discovery, our model can be integrated with the generative model developed by Gómez-Bombarelli *et al.* [47], who combined a generative model with the search over latent space to generate new molecules. In 3D shape synthesis, our model can be trained on large datasets like ShapeNet and used as a recommender system within CAD software. The loss function we develop can also be integrated with human face synthesis methods, to generate new human faces, which are high quality (depending on any criteria like beauty) and from different groups (regions, race, gender, age *etc.*). Overall, the method provides a new direction of research, where generative models focus on the unbiased generation of high-quality items.

References

- [1] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- [2] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in neural information processing systems*, pages 2672–2680, 2014.

- [3] Wei Chen, Mark Fuge, and Jonah Chazan. Design manifolds capture the intrinsic complexity and dimension of design spaces. *Journal of Mechanical Design*, 139(5), 2017.
- [4] Wei Chen and Mark Fuge. Synthesizing designs with interpart dependencies using hierarchical generative adversarial networks. *Journal of Mechanical Design*, 141(11), 2019.
- [5] Wei Chen, Kevin Chiu, and Mark Fuge. Aerodynamic design optimization and shape exploration using generative adversarial networks. In *AIAA SciTech Forum*, San Diego, USA, Jan 2019. AIAA.
- [6] Martin Philip Bendsoe and Ole Sigmund. *Topology Optimization: Theory, Methods and Applications*. Springer, February 2004.
- [7] Faez Ahmed, Kalyanmoy Deb, and Bishakh Bhattacharya. Structural topology optimization using multi-objective genetic algorithm with constructive solid geometry representation. *Applied Soft Computing*, 39:240–250, 2016.
- [8] Dule Shu, James Cunningham, Gary Stump, Simon W Miller, Michael A Yukish, Timothy W Simpson, and Conrad S Tucker. 3d design using generative adversarial networks and physics-based validation. *Journal of Mechanical Design*, 142(7), 2020.
- [9] Alex Kulesza and Ben Taskar. Determinantal point processes for machine learning. *arXiv preprint arXiv:1207.6083*, 2012.
- [10] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep learning*. MIT press, 2016.
- [11] Thomas Gmeiner and Kristina Shea. A spatial grammar for the computational design synthesis of vise jaws. In *ASME 2013 International Design Engineering Technical Conferences and Computers and Information in Engineering Conference*. American Society of Mechanical Engineers Digital Collection, 2013.
- [12] Corinna Königseder, Tino Stanković, and Kristina Shea. Improving design grammar development and application through network-based analysis of transition graphs. *Design Science*, 2, 2016.
- [13] Kristina Shea, Robert Aish, and Marina Gourtovaia. Towards integrated performance-driven generative design tools. *Automation in Construction*, 14(2):253–264, 2005.
- [14] Daniel R Herber, Tinghao Guo, and James T Allison. Enumeration of architectures with perfect matchings. *Journal of Mechanical Design*, 139(5), 2017.
- [15] Vinjamuri Venkata Kamesh, Kuchibhotla Mallikarjuna Rao, Srinivasa Rao, and Annambhotla Balaji. Topological synthesis of epicyclic gear trains using vertex incidence polynomial. *Journal of Mechanical Design*, 139(6), 2017.
- [16] Cari R Bryant, Robert B Stone, Daniel A McAdams, Tolga Kurtoglu, Matthew I Campbell, et al. Concept generation from the functional basis of design. In *ICED 05: 15th International Conference on Engineering Design: Engineering Design and the Global Economy*, page 1702. Engineers Australia, 2005.
- [17] Amaresh Chakrabarti, Kristina Shea, Robert Stone, Jonathan Cagan, Matthew Campbell, Noe Vargas Hernandez, and Kristin L Wood. Computer-based design synthesis research: an overview. *Journal of Computing and Information Science in Engineering*, 11(2), 2011.
- [18] David F Wyatt, David C Wynn, Jerome P Jarrett, and P John Clarkson. Supporting product architecture design using computational design synthesis with network structure constraints. *Research in Engineering Design*, 23(1):17–52, 2012.
- [19] Jan Wijnkiet and Theo Hofman. Modified computational design synthesis using simulation-based evaluation and constraint consistency for vehicle powertrain systems. *IEEE Transactions on Vehicular Technology*, 67(9):8065–8076, 2018.
- [20] Xi Chen, Matteo Diez, Manivannan Kandasamy, Zhiguo Zhang, Emilio F Campana, and Frederick Stern. High-fidelity global optimization of shape design by dimensionality reduction, metamodels and deterministic particle swarm. *Engineering Optimization*, 47(4):473–494, 2015.
- [21] Danny D’Agostino, Andrea Serani, Emilio F Campana, and Matteo Diez. Nonlinear methods for design-space dimensionality reduction in shape optimization. In *International Workshop on Machine Learning, Optimization, and Big Data*, pages 121–132. Springer, 2017.
- [22] Danny D’Agostino, Andrea Serani, Emilio F Campana, and Matteo Diez. Deep autoencoder for off-line design-space dimensionality reduction in shape optimization. In *2018 AIAA/ASCE/AHS/ASC Structures, Structural Dynamics, and Materials Conference*, page 1648, 2018.
- [23] Alexander Burnap, Ye Liu, Yanxin Pan, Honglak Lee, Richard Gonzalez, and Panos Y Papalambros. Estimating and exploring the product form design space using deep generative models. In *ASME 2016 International Design Engineering Technical Conferences and Computers and Information in Engineering Conference*. American Society of Mechanical Engineers Digital Collection, 2016.

- [24] James D Cunningham, Timothy W Simpson, and Conrad S Tucker. An investigation of surrogate models for efficient performance-based decoding of 3d point clouds. *Journal of Mechanical Design*, 141(12), 2019.
- [25] Ruijin Cang, Aditya Vipradas, and Yi Ren. Scalable microstructure reconstruction with multi-scale pattern preservation. In *ASME 2017 International Design Engineering Technical Conferences and Computers and Information in Engineering Conference*. American Society of Mechanical Engineers Digital Collection, 2017.
- [26] Zijiang Yang, Xiaolin Li, L Catherine Brinson, Alok N Choudhary, Wei Chen, and Ankit Agrawal. Microstructural materials design via deep adversarial learning methodology. *Journal of Mechanical Design*, 140(11), 2018.
- [27] Wei Chen and Mark Fuge. Béziergan: Automatic generation of smooth curves from interpretable low-dimensional parameters. *arXiv preprint arXiv:1808.08871*, 2018.
- [28] Sangeun Oh, Yongsu Jung, Seongsin Kim, Ikjin Lee, and Namwoo Kang. Deep generative design: Integration of topology optimization and generative models. *Journal of Mechanical Design*, 141(11), 2019.
- [29] Alex Burnap, John R Hauser, and Artem Timoshenko. Design and evaluation of product aesthetics: A human-machine hybrid approach. *Available at SSRN 3421771*, 2019.
- [30] Tim Salimans, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, and Xi Chen. Improved techniques for training gans. In *Advances in neural information processing systems*, pages 2234–2242, 2016.
- [31] Xudong Mao, Qing Li, Haoran Xie, Raymond YK Lau, Zhen Wang, and Stephen Paul Smolley. Least squares generative adversarial networks. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2794–2802, 2017.
- [32] Duhyeon Bang and Hyunjung Shim. Mggan: Solving mode collapse using manifold guided training. *arXiv preprint arXiv:1804.04391*, 2018.
- [33] Akash Srivastava, Lazar Valkov, Chris Russell, Michael U Gutmann, and Charles Sutton. Veegan: Reducing mode collapse in gans using implicit variational learning. In *Advances in Neural Information Processing Systems*, pages 3308–3318, 2017.
- [34] Xi Chen, Yan Duan, Rein Houthoofd, John Schulman, Ilya Sutskever, and Pieter Abbeel. Infogan: Interpretable representation learning by information maximizing generative adversarial nets. In *Advances in neural information processing systems*, pages 2172–2180, 2016.
- [35] Mohamed Elfeki, Camille Couprie, Morgane Riviere, and Mohamed Elhoseiny. Gdpp: Learning diverse generations using determinantal point processes. In *International Conference on Machine Learning*, pages 1774–1783, 2019.
- [36] Apramey Dube, Anu Helkkula, et al. Customer approach to the use of big data: Wearables for service. In *Proceedings of SERVSIG 2016 Conference*. Maastricht University, 2016.
- [37] Hui Lin and Jeff Bilmes. A class of submodular functions for document summarization. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies-Volume 1*, pages 510–520. Association for Computational Linguistics, 2011.
- [38] Jami J Shah, Santosh V Kulkarni, and Noe Vargas-Hernandez. Evaluation of idea generation methods for conceptual design: effectiveness metrics and design of experiments. *Journal of Mechanical Design*, 122(4):377–384, 2000.
- [39] Mark Fuge, Josh Stroud, and Alice Agogino. Automatically inferring metrics for design creativity. *ASME Paper No. DETC2013-12620*, 2013.
- [40] Faez Ahmed, Sharath Kumar Ramachandran, Mark Fuge, Sam Hunter, and Scarlett Miller. Measuring and optimizing design variety using herfindahl index. In *ASME 2019 International Design Engineering Technical Conferences and Computers and Information in Engineering Conference*. American Society of Mechanical Engineers Digital Collection, 2019.
- [41] Faez Ahmed and Mark Fuge. Ranking ideas for diversity and quality. *Journal of Mechanical Design*, 140(1), 2017.
- [42] Faez Ahmed, Mark Fuge, and Lev D Gorbunov. Discovering diverse, high quality design ideas from a large corpus. In *ASME 2016 International Design Engineering Technical Conferences and Computers and Information in Engineering Conference*. American Society of Mechanical Engineers Digital Collection, 2016.
- [43] Alex Kulesza and Ben Taskar. k-dpps: Fixed-size determinantal point processes. In *Proceedings of the 28th International Conference on Machine Learning (ICML-11)*, pages 1193–1200, 2011.
- [44] Alexei Borodin. Determinantal point processes. *arXiv preprint arXiv:0911.1153*, 2009.

- [45] Mark Drela. Xfoil: An analysis and design system for low reynolds number airfoils. In *Low Reynolds number aerodynamics*, pages 1–12. Springer, 1989.
- [46] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [47] Rafael Gómez-Bombarelli, Jennifer N Wei, David Duvenaud, José Miguel Hernández-Lobato, Benjamín Sánchez-Lengeling, Dennis Sheberla, Jorge Aguilera-Iparraguirre, Timothy D Hirzel, Ryan P Adams, and Alán Aspuru-Guzik. Automatic chemical design using a data-driven continuous representation of molecules. *ACS central science*, 4(2):268–276, 2018.

Appendix A: Table of Evaluation Metrics

We list the means and 95% confidence intervals of all evaluation scores (Figs. 7-9 and Fig. 12) in Table 1. It shows that PaDGAN received best overall score for all cases, and atleast the second best score for both diversity and quality in all four examples.

Table 1: Diversity and Quality Scores for all experiments. The best score for each example is marked by * symbol and the second best scores by ** symbol.

	Model	Diversity Score	Quality Score	Overall Score
Example I	GAN	-69.0885 ± 27.4357	0.5826 ± 0.0620	$1.0678 \pm 0.2938^{**}$
	GAN_D	$-53.6183 \pm 0.3601^*$	0.4672 ± 0.0162	0.9060 ± 0.0414
	GAN_Q	-152.1945 ± 77.5582	$0.9131 \pm 0.1065^*$	0.5913 ± 0.5875
	PaDGAN	$-57.2489 \pm 1.16202^{**}$	$0.6955 \pm 0.0269^{**}$	$1.4897 \pm 0.0624^*$
Example II	GAN	-68.2451 ± 3.9807	0.5004 ± 0.0447	$0.9752 \pm 0.0531^{**}$
	GAN_D	$-61.8685 \pm 0.7375^*$	0.4811 ± 0.0239	0.9473 ± 0.0269
	GAN_Q	-197.0243 ± 119.0978	$0.9914 \pm 0.0280^*$	0.4261 ± 0.8780
	PaDGAN	$-65.5501 \pm 1.53222^{**}$	$0.7176 \pm 0.0344^{**}$	$1.2117 \pm 0.0320^*$
Example III	GAN	-93.1237 ± 67.4991	0.6219 ± 0.0847	0.9019 ± 0.8333
	GAN_D	$-63.5478 \pm 0.4637^*$	0.5560 ± 0.0152	$1.0811 \pm 0.0468^{**}$
	GAN_Q	-150.3896 ± 65.4173	$0.8492 \pm 0.1219^*$	0.6356 ± 0.6366
	PaDGAN	$-63.9063 \pm 1.76832^{**}$	$0.6968 \pm 0.0248^{**}$	$1.4993 \pm 0.0624^*$
Airfoil	GAN	$-22.0326 \pm 1.4982^{**}$	$56.1703 \pm 66.4480^{**}$	N/A
	PaDGAN	$-12.5455 \pm 4.8970^*$	$76.0670 \pm 43.9802^*$	N/A