# Comprehensive Exploration of Synthetic Data Generation: A Survey

**André Bauer**
University of Chicago
Chicago, United States of America
`andrebauer@uchicago.edu`

**Simon Trapp, Michael Stenger, Robert Leppich, Samuel Kounev**
University of Würzburg
Würzburg, Germany
`{firstname}.{lastname}@uni-wuerzburg.de`

**Mark Leznik**
University of Ulm
Ulm, Germany
`mark.leznik@uni-ulm.de`

**Kyle Chard, Ian Foster**
Argonne National Laboratory
Lemont, United States of America
`{lastname}@anl.gov`

## Abstract

Recent years have witnessed a surge in the popularity of Machine Learning (ML), applied across diverse domains. However, progress is impeded by the scarcity of training data due to expensive acquisition and privacy legislation. Synthetic data emerges as a solution, but the abundance of released models and limited overview literature pose challenges for decision-making. This work surveys 417 Synthetic Data Generation (SDG) models over the last decade, providing a comprehensive overview of model types, functionality, and improvements. Common attributes are identified, leading to a classification and trend analysis. The findings reveal increased model performance and complexity, with neural network-based approaches prevailing, except for privacy-preserving data generation. Computer vision dominates, with GANs as primary generative models, while diffusion models, transformers, and RNNs compete. Implications from our performance evaluation highlight the scarcity of common metrics and datasets, making comparisons challenging. Additionally, the neglect of training and computational costs in literature necessitates attention in future research. This work serves as a guide for SDG model selection and identifies crucial areas for future exploration.

***Keywords*** Survey · Synthesis · Synthetic Data Generation

## 1 Introduction

In recent years, Artificial Intelligence (AI), particularly in subfields like Machine Learning (ML) and Deep Learning (DL), has experienced significant growth and popularity [1, 2]. As ML and DL models have evolved in complexity and efficiency, a persistent challenge has been the limited size of training datasets. This limitation stems from high labeling costs and privacy concerns, hindering the model's ability to generalize effectively. To address this issue, Synthetic Data Generation (SDG) emerges as a viable solution, providing substantial amounts of artificial data for model training. This artificial data includes novel, diverse, and realistic samples, alleviating the constraints imposed by traditional datasets [2].

Broadly speaking, SDG involves generating artificial data and labels, aiming to emulate authentic samples closely. This process is automated by utilizing generative models that estimate the probability distribution of their training data. This sets it apart from data augmentation, which manipulates existing data, as SDG generates new data by sampling from learned distributions. The advantages of synthetic data go beyond mere cost reduction, with its on-the-fly generation contributing to reduced computational time and addressing bias in data distribution. Synthetic data proves highly valuable when real data is insufficient, costly to label, or exhibits biased distributions.

The concept of synthetic data dates back to the 1960s and has evolved alongside the broader AI landscape. However, the dynamic nature of SDG, with new approaches emerging annually, poses a challenge for researchers, especially those

new to the field. Existing literature on SDG often lacks a comprehensive overview and classification of approaches, making it difficult to stay current with generative models, their applications, and the relationships between them. Recent studies have summarized the literature on specific model types such as Generative Adversarial Nets (GANs) [3, 4] and computer-rendered virtual 3D environments [5, 6]. Other works concentrate on distinct domains, including graph generation [7], computer vision [8, 9], text generation [10], music [11], privacy [12], and molecular science [13]. Additionally, some related literature has undertaken the task of compiling, comparing, or classifying various SDG approaches comprehensively [14, 15, 16, 17, 2]. However, these works have limitations in scope, often overlooking recent advancements such as self-attention [18]. They may focus exclusively on a single domain, classify models based on only a few aspects, or provide coarse coverage of literature and model architecture.

To address these gaps in the existing literature, this work endeavors to provide a comprehensive overview of SDG. More precisely, the contributions of this work include:

1. **Surveying Literature:** We conduct an extensive survey of the literature from the last decade, aiming to comprehensively cover all model types suitable for SDG. Our investigation involves scrutinizing 417 models, revealing 20 distinct model types, further categorized into 42 subtypes.

2. **Exploring Applications and Enhancements:** Within this survey, we delve into the applications and enhancements of the identified SDG model types, providing insights into their respective practical implementations.

3. **Introducing Classification Categories:** In addition to model identification, we introduce various categories for classifying the collected generative models. These categories include generated data types, performance, privacy considerations, and training processes.

4. **Knowledge Foundation:** The acquired knowledge from this exploration serves as a solid foundation, providing a comprehensive understanding of the diverse landscape of SDG model types.

5. **Guideline Development:** Building upon this foundation, we develop a practical guideline. This guideline is tailored to facilitate the selection of an appropriate SDG model type, offering valuable insights for researchers and practitioners in the field.

Our survey of 417 SDG models reveals notable trends. We classified these models based on over ten criteria, demonstrating an evident increase in complexity and performance over the years. The evolution is marked by the shift from simpler probabilistic models like Markov chains and Bayesian Networks (BNs) to more sophisticated neural network-based approaches. Notably, in the realm of SDG, computer vision stands out as the most popular application field. GANs and diffusion models have emerged as top performers in this domain. For handling sequential data, such as music or text, Recurrent Neural Networks (RNNs) dominate. Additionally, privacy-preserving data generation commonly employs models like Markov chains, BNs, and more advanced GANs. The work also highlights challenges such as non-standardized evaluation metrics and datasets, offering solutions like building graphs of predecessors based on the models' performance evaluations [19].

We are convinced that our work can serve as a valuable resource for newcomers and experienced researchers, providing an updated overview of the SDG landscape and aiding in identifying suitable models and datasets for specific tasks. Furthermore, we present the first comprehensive classification of numerous model implementations spanning multiple domains, facilitating the identification of strengths and weaknesses in these models.

The remainder of this work is structured as follows: In Section 2, we introduce 42 different SDG model types and highlight their usage in literature. In Section 3, we classify and compare the found generative models according to different categories before discussing a guideline for choosing which SDG model type is suitable for a given scenario. In Section 4, review existing surveys and distinguish the scope of our study from existing surveys. In Section 5, we conclude our paper. In Appendix A, we list all used acronyms.

## 2 Overview of Generative Models

Various approaches exist to generate synthetic data, ranging from models based on graphs or simple probabilistic assumptions to deep neural networks. The following sections describe the different available model architectures in general and provide a chronological overview of the recent literature regarding the usage of these models for SDG. The structure is inspired by the work done by Harshvardhan et al. [16] and extended further to include missing models and some novel and important implementations we found. We focus on approaches released in the last ten years to keep this work within viable and reasonable bounds. We mainly restrict ourselves to models referenced by the works mentioned in Section 4, extended by some of their often-cited literature and additional material from our research on Google Scholar.

## 2.1 Gaussian Mixture Models

Gaussian Mixture Models (GMMs) are density estimation algorithms mostly used for data clustering but can also be used as a generative probabilistic model. A GMM consists of $N$ Gaussian distributions (components), which are normal distributions that are continuous for a real-valued random variable and symmetric about their mean. Each Gaussian can be characterized by a mean $\mu_i$ and a variance $\rho_i$ and has a probability/weight $\pi_i$ in the mixture model so that $\sum_{i=1}^{N} \pi_i = 1$. For one-dimensional data, $\pi_i$ and $\mu_i$ are numbers, but in two-dimensional space, $\pi_i$ is a vector, and $\mu_i$ is a covariance matrix. The probability of a data point $d$ belonging to the cluster represented by Gaussian $i$ is $P(d = i) = \pi_i$ and the observation likelihood of $d$ in Gaussian $i$ is $P(d|d = i, \mu_i, \rho_i) = N(d|\mu_i, \rho_i)$ where $N(.)$ is the normal distribution function. [16]



(a) One-dimensional GMM with three components. (Source: [16])

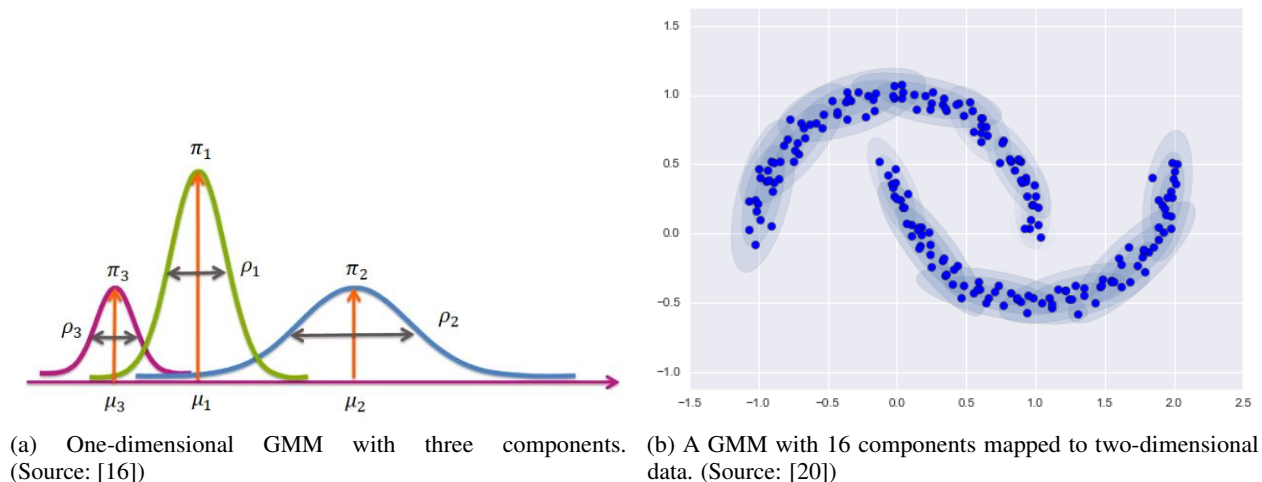(b) A GMM with 16 components mapped to two-dimensional data. (Source: [20])

Figure 1: Illustrations of GMMs used for one and two-dimensional data.

To map data points to $N$ clusters/components of a GMM during training, an expectation-maximization (EM) algorithm is used:

1. Choose (random) locations $\mu_i$ and shapes $\rho_i$ for each component

2. Repeat until convergence:

    (a) *E-step*: For each data point, find $\pi_i$ encoding the membership probability

    (b) *M-step*: For each cluster, update $\mu_i$, $\rho_i$ and $\pi_i$ based on all data points

Multiple GMMs can be initialized with different random values to reduce the chance of missing the globally optimal solution. [20].

Van der Oord [21] propose a deep GMM, which contains multiple layers of linear transformations $x = Az + b$ applied to the normal variable $z \sim \mathcal{N}(0, I_n)$. For each sample, a random path through the graph is taken (see Figure 2), and the transformations are concatenated. In theory, deep GMMs can be represented by normal GMMs, but training would be more complex, and the deep models generalize better. The model is trained with an EM algorithm and parallelizable by design. Its capabilities are demonstrated by generating low-resolution greyscale images.

Zen and Senior [22] applied mixture density networks for speech synthesis to enable multimodal regression and the prediction of variances. To this end, the authors model the conditional probability distribution using a GMM, with parameters predicted by a fully connected multi-layer neural network.

VanderPlas [20] demonstrates how GMMs can be fitted well on two-dimensional data (see Figure 1b) by trying different component sizes $N$ and searching for the model with the optimal Akaike information criterion (AIC) or Bayesian information criterion (BIC). Synthetic data can be obtained from this trained model by randomly selecting a Gaussian according to probability $\pi_i$ and sampling from its normal distribution. The author uses this approach on binary black-and-white images of handwritten digits and generates authentic artificial samples.
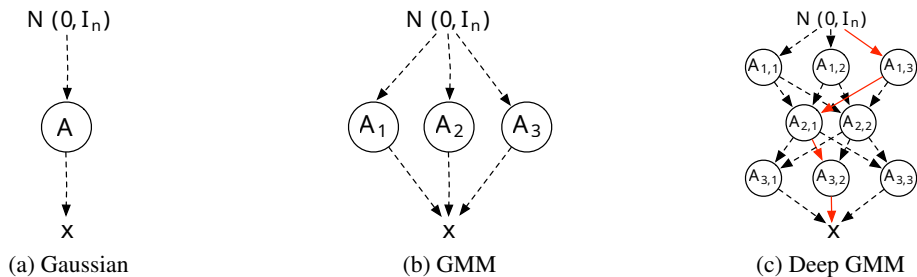
(a) Gaussian          (b) GMM          (c) Deep GMM

Figure 2: Comparison between a Gaussian, GMM and a deep GMM with transformation biases $b_{i,j}$ not shown. (Source: [21])

## 2.2 Kernel Density Estimators

Kernel Density Estimatorss (KDEs) are models that approximate the probability density function $p(X = x)$ of a random variable $X$ in between a set of observations. For observations $x_i$ and an optimizable smoothing parameter $h > 0$, the density estimation is defined as

$$\hat{p}(x) = \frac{1}{n} \sum_{i=1}^{n} K(x - x_i, h) \tag{1}$$

with a positive kernel function $K(y, h)$, for which various implementations exist. One kernel often used is the Gaussian kernel $K(y, h) = \exp(-\frac{y^2}{2h^2})$, but many others are available to model different kinds of data. [23]

Bhattacharya et al. [23] synthesize Photoplethysmograms (PPGs), which are time series over volumetric blood flow in the human body, by decomposing it into components (pulse length, peak position, amplitude, etc.), training a KDE for each component and sampling from the resulting probability distributions, from which new PPGs can be constructed.

## 2.3 Markov Chain Models

Order $n$ Markov chains are probabilistic models for infinite sequences of symbols where the probability for each symbol only depends on the previous $n$ symbols [11]. For $n = 1$, they can be drawn as a graph with the states (symbols) as nodes and the transition probabilities $a_{st} = P(x_i = t | x_{i-1} = s)$ as edges (see Figure 3). The probability of a whole sequence can be computed by applying $P(X, Y) = P(X|Y)P(Y)$ many times; that is, the probability of a sequence $x$ with length $L$ is $P(x) = P(x_1) \prod_{i=2}^{L} a_{x_{i-1}x_i}$. [24]



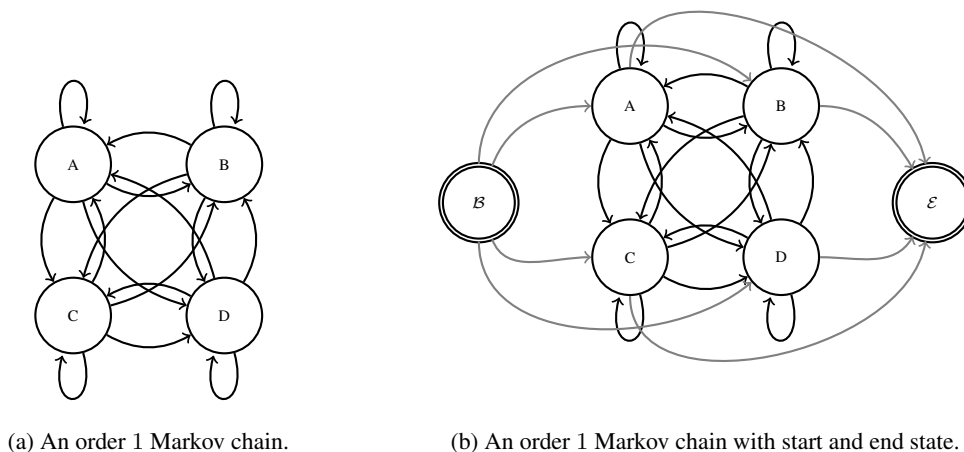(a) An order 1 Markov chain.          (b) An order 1 Markov chain with start and end state.

Figure 3: Illustrations of graphs of Markov chains with symbols as nodes and transition probabilities as edges. (Adapted from: [24])

4

The advantage of Markov chains is their simplicity. They can be easily understood because the conditional probabilities can be computed by counting relative symbol appearances. Further, they can be interpreted as automata on which additional control mechanisms like Markov constraints and factor graphs can be imposed. On the downside, simple order 1 models can not capture long-term temporal structures, and order $n$ models tend to overfit, require significantly bigger amounts of training data, and need to compute exponentially more conditional probabilities for large $n$. [11, 25]

Pachet et al. [26] introduce the Continuator, a system for interactive music generation in the user's style without a priori musical knowledge, allowing a musician to "jam" in real-time with the computer. The Continuator is powered by an *analysis* module and a *generator*. The analysis module first detects the ends of musical phrases, then builds a Markovian model of these phrases and detects global properties like tempo, meter, and note density. The generator uses the Markov model and the properties to generate music in the input style and continue it note-by-note.

### 2.3.1 Hidden Markov Models

Hidden Markov Models (HMMs) are extensions of Markov chains where the Markov chain state sequence $\pi$ is *hidden* and each hidden state $k$ has emission probabilities $e_k(b) = P(x_i = b | \pi_i = k)$ for the observable symbols $x$ (see Figure 4a). This model can depict many issues with reduced complexity compared to simple Markov chains, but the inference of hidden states $\pi$ for observation $x$ with joint probability $P(x, \pi) = a_{0\pi_1} \prod_{i=1}^{L} e_{\pi_i}(x_i) a_{\pi_i \pi_{i+1}}$ is more complex (see for example Figure 4b). The most probable state sequence $\pi^* = \arg \max_\pi P(x, \pi)$ for an observation $x$ can be recursively computed with the Viterbi algorithm. [16, 24]



(a) Architecture of an HMM. (Source: [16])

(b) HMM of a dishonest casino switching between fair and unfair dices. (Source: [24])
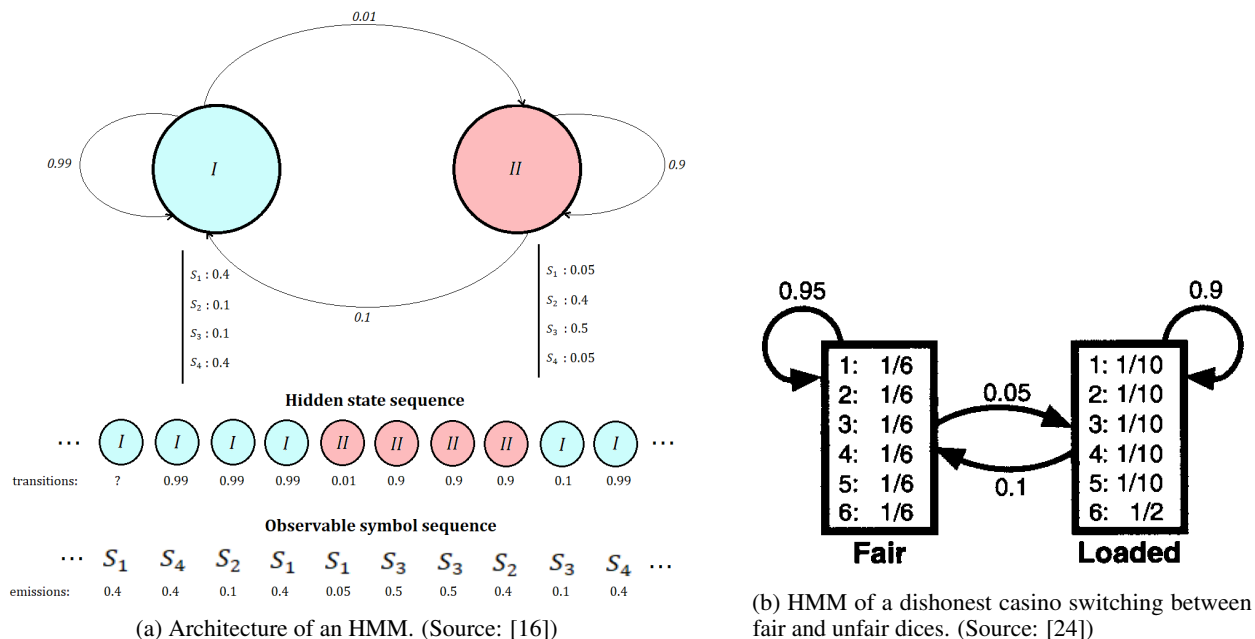
Figure 4: Illustrations of HMMs.

Durbin et al. [24] use HMMs to model biological (genetic) sequences and demonstrate the labeling of unannotated and generation of new data for this topic. They show that these simple models can learn truthful models even from observations where the hidden paths are unknown using the Baum-Welch and Viterbi algorithms. They also discuss different model topologies for different sequence lengths and find the careful topology construction of a HMM validated by human experts to be essential for good model performance.

Racyzński et al. [27] interpolate results using multiple learned sub-models (namely, a bigram Markov chain model $P(C_t | C_{t-1})$, a tonality relation $P(C_t | T_t)$, and a melody relation model $P(C_t | M_t)$) to sequentially generate chords $C_i$ as accompaniment for music.

Kaliakatsos-Papakostas et al. [28] propose the constrained HMM (CHMM), which allows intermediate states of the sequence to be fixed to a specific value (anchor chords) with probability 1. Then, the Viterbi algorithm is used to find the most likely path between these checkpoints. The model is used to generate harmonic chord sequences for a melody by mapping the hidden note states between the anchors defined by the melody to chords.

Bindschaedler and Shokri [29] use HMMs to generate plausible and private location traces by clustering the available locations and synthesizing a sequence of locations from the same cluster label sequence as a real seed trace. The generated trace is further evaluated using a plausibility and a privacy (geographic similarity to the seed trace) test.

### 2.3.2 N-Grams

The *n-gram* is a tuple of $n$ values, for instance, a sequence of $n$ words of a sentence $(w_1, ..., w_n)$. They are usually used to efficiently model the probabilities of words based on the already written text:

$$p(X_t|X_{t-1}, ..., X_{t-n+1}) = \frac{count(X_{t-n+1}, ..., X_t) + 1}{count(X_{t-n+1}, ..., X_{t-1}) + V} \tag{2}$$

under the Markovian assumption $p(X_t|X_{t-1}, ..., X_1) = p(X_t|X_{t-1}, ..., X_{t-n+1})$, that is, the probability of a word at position $t$ only depends on the previous $n-1$ words. Adding Laplace smoothing with the $+1$ in the numerator and the word vocabulary size $V$ in the denominator also allows previously unseen word combinations to be created with low probability. N-grams are normally used for language modeling and synthetic text generation. [30]

Bengio et al. [31] extend the idea of n-grams with a Multilayer Perceptron (MLP): First, each word $X_k$ is encoded into a 1-hot vector $\mathbf{1}(X_k)$. Then, the vectors are linearly embedded using matrix $W_x$, concatenated, and fed into the MLP, which is trained to predict the next word probability. Finally, a softmax function $SM$ is applied, resulting in the neural network language model:

$$p(X_t|X_{t-1}, ..., X_{t-n+1}) = SM(MLP[W_x\mathbf{1}(X_{t-1}), ..., W_x\mathbf{1}(W_{t-n+1})]). \tag{3}$$

Barbieri et al. [32] implement Markov constraints to generate lyrics in a given style and rhythm for music. They train a Markov process using relative frequencies of n-grams of lyrics, for example, of a specific author and then restructure it as a finite-length sequence of constrained variables with assigned probability for each value. The constraints are related to rhyme, rhythm, syntax (part-of-speech templates), and semantics (relations between words).

Short overview of other usages of n-grams:

| Approach | Description | Year |
|---|---|---|
| [33] | Extension of the constrained Markov model from [32] to allow a global meter constraint to be efficiently implemented. | 2013 |
| [34] | Construction of a probabilistic finite state automat from n-grams of preprocessed MIDI files to harmonically accompany a melody. | 2013 |
| [35] | Avoidance of plagiarism in generated sequences in high-order Markov chains with a *MaxOrder* global Markov constraint that prevents chunks longer than *MaxOrder* from being replicated from the training data by building Markov automatons with restricted maximum path lengths. | 2014 |
| [36] | FlowComposer: A web tool consisting of two collaborating constrained Markov models (melody and chords) for generation, re-harmonization, and interactive composition of music lead sheets. | 2016 |
| [37] | A multiple viewpoint system consisting of Markov chains obtained from n-grams generates music. The agents are responsible for different aspects of the music and are ordered sequentially and hierarchically. | 2016 |

### 2.4 Bayesian Networks

[38]) A Bayesian Network (BN) is a Directed Acyclic Graph (DAG), a special type of graphical model in which random variables are the nodes and dependencies between these variables are the edges. The directed edges run from the "cause" or parent node to the "effect" or child node and define the conditional dependency of the nodes. Each random variable has a continuous or discrete probability distribution function. An example for a Bayesian Network can be seen in Figure 5. [39]

If we define $\mathbf{V} = \{V_j : j \in \{1, ..., N\}\}$ as the set of random variables of a Bayesian network, the probability distribution of $V_j$ as $p(V_j|\mathbf{Pa}_j)$ and $\mathbf{Pa}_j$ as the set of parents of $V_j$, the joint distribution over $\mathbf{V}$ is defined as [41]:

$$p(\mathbf{V}) = p(V_1, ..., V_N) = \prod_{j=1}^{N} p(V_j|\mathbf{Pa}_j). \tag{4}$$

| C | P(R) |
|---|---|
| T | .80 |
| F | .20 |

P(C) = .50

Rain

Cloudy

WetGrass

Sprinkler

| C | P(S) |
|---|---|
| T | .10 |
| F | .50 |

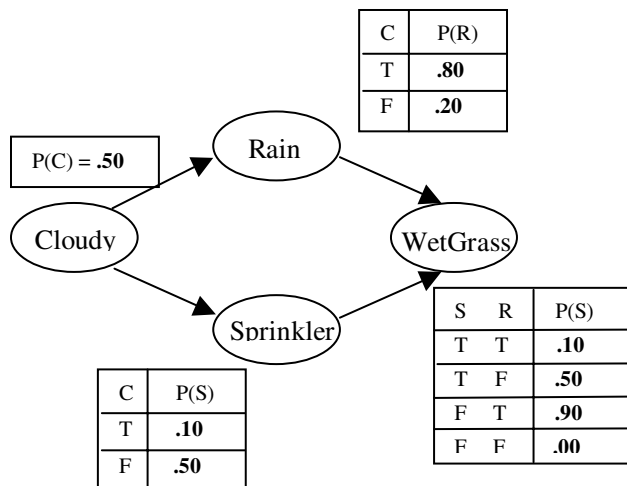| S | R | P(S) |
|---|---|---|
| T | T | .10 |
| T | F | .50 |
| F | T | .90 |
| F | F | .00 |

Figure 5: Example of a Bayesian Network with discrete random variables. (Source: [40])

There are many algorithms to train BNs under different conditions where the net structure is known or unknown in advance, and the reference data is fully or partially observable. In the context of synthetic data, we usually either have the trivial case of available expert knowledge, from which the structure of a BN can be constructed, or available real-world data, but no information about the structure of a BN that corresponds well to the data. For the latter case, two problems need to be solved:

1. Finding a metric to compare potential structures of BNs.
2. Searching for potential BN structures algorithmically.

A solution for the first problem is provided by the joint probability $p(D, S^h)$ for the data $D$ and a hypothetical structure $S^h$ of a BN:

$$\log p(D, S^h) = \log p(D|S^h) + \log p(S^h). \tag{5}$$

The Bayesian information criterion (BIC) [42] can be used to calculate $\log p(D|S^h)$ while the prior probability $p(S^h)$ of a structure can be determined for example by assigning probabilities to a predefined set of possible structures or defining a prior network and measuring the deviation of $S^h$ from it. [39]

The second problem, searching for structures, is NP-hard if done for all possible combinations, so different greedy algorithms are employed. In general, these algorithms increase $p(S^h)$ step-wise by performing actions on the graph (adding, removing, or reversing an edge) until a maximum is reached. [39]

In recent years, more sophisticated learning algorithms have been developed. They use dynamic programming [43, 44, 45] to split the global learning problem into small subproblems. Others define the learning task as a shortest path problem and solve it with an A* algorithm [46].

Young et al. [47] use BNs to anonymize a data set so it can be disclosed to the public. They limit their networks to discrete variables and map continuous values (e.g., age) to discrete intervals to facilitate the learning process, which consists of multiple steps:

1. The user defines a prior network and conditional probability distributions for each node. Also, an imaginary database is supplied to generate confidence in the prior structure.
2. The probability distributions of the nodes are updated using the training data.
3. A greedy search algorithm [48] starts the search from the prior network, creates all possible networks with one change (edge addition, removal, or reversion), and selects the one with the highest Bayes factor (likelihood of the model according to the data) as the new baseline until no further improvement occurs.

The trained BN is now used as an imputation model [49] to generate synthetic data by consecutively drawing random samples from the nodes in the hierarchy.

Suzuki et al. [50] generate fixed-size four-part (alto, tenor, bass, soprano) symbolic harmonies based on the melody of a soprano voice and experiment with the conditioning on chord nodes. The pitches for each voice are classified jointly based on the previous value (Markov property) and the current soprano or chord value. The soprano network without chords produces smoother results.

Zhang et al. [51] propose *PrivBayes*, a differentially private[1] method to release high-dimensional data sets. First, a BN with succinct attribute correlations is created. Then noise is injected into the low-dimensional marginal distributions of the BN, and the now noisy approximated data distribution is used to sample a private synthetic data set.

Draghi et al. [38] approach the bias problem in the training data of BNs by identifying under-represented cases and over-sampling them with synthetic data. They start by training a BN on a subset of the original data, modifying it, and generating a new biased data set. Then, a BN is trained on this biased data and generates the data set $D_{bias}$ with the same size as the original. Next, a classifier is trained on $D_{bias}$ and tries to predict values from a validation set, which is a subset of the original data and adds samples with an uncertain outcome (low probability) to a data set $D_{unc}$. Finally, the *BayesBoost* is performed by generating $m$ similar samples for each sample in $D_{unc}$ with a BN trained on $D_{bias}$ and adding the results to $D_{bias}$, resulting in a new less biased dataset.

Short overview of other usages of BNs:

| Approach | Description | Year |
|---|---|---|
| [53] | Training a BN on motion capture data to produce realistic human 3D poses that are rendered with different textures (e.g., clothes) to produce training data for human 3D pose estimation. | 2016 |
| [54] | DataSynthesizer: Creating a Bayesian Network from data with a *DataDescriber*, injecting noise into the distributions and sampling from it with the *DataGenerator*. The *ModelInspector* compares the properties of the synthetic data to the real one. | 2017 |
| [55] | Modeling heterogeneous (continuous and discrete variables) medical patient data as a BN to be able to incorporate expert knowledge and generate private data sets. They experiment with three ways to deal with missing values: Deleting the entire entry, adding "miss states/nodes" to the BN, and using the FCI algorithm [56] to infer the missing variables. Using probabilistic graphical modeling, the model produces high-fidelity results with a low risk of patient re-identification (cloning of training data). | 2020 |
| [38] | Exploring the data bias problem and under-representation in underlying ground truth samples. Specifically, it is an important problem in medical data, where synthetic data generation is used to mask sensitive patient data. The authors propose an approach to identifying under-sampled data and improving data synthesis to correct this problem. | 2021 |

## 2.5 Genetic Algorithms

Genetic Algorithms (GAs) are ML algorithms that mimic the natural selection process over time. The population consists of a crowd of individuals at each time step or *generation*. To create the population for the next generation, three actions are performed:

**Selection** Selection of suitable candidates from the population, using a *fitness function* to eliminate worse candidates and increase the chance of survival of better ones so they can pass on their good properties. Individuals can be selected multiple times to maintain population size.

**Crossover** Information exchange between candidates.

**Mutation** Perturb the candidates' information by randomly changing some properties, usually according to some distribution.

This generation process (illustrated in Figure 6) continues until the system converges (i.e., all candidates are identical) or a user-defined criterion is met. The speed of the GA model is measured by the number of generations needed to meet the requirements. [57]

Liu et Ting [59] use a GA for polyphonic accompaniment generation given a dominant music melody. They remove the impractical need for a human evaluation criterion of previous approaches by building a fitness function based on music theory.

You et Liu [58] propose a GA that finds similar and suitable chord variations for a given target melody and some example chord progressions from other songs provided as MIDI files. The initial population consists of the exemplar

---

[1]Differential privacy introduces randomness to the data provision process, resulting in "plausible deniability of any outcome" [52]
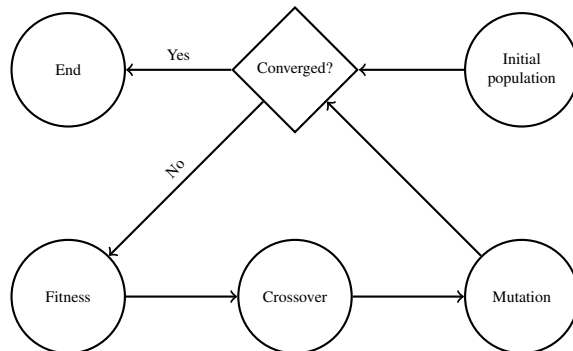
Figure 6: Process of a GA. (Adapted from [58] and modified)

chord patterns with keys shifted to match the key of the target melody. The fitness function and crossover process further incorporate music theory to ensure that the evolving chord patterns harmonize with the melody and conform to basic rules.

Chen et al. [57] implement a GA for categorical tabular data with non-hierarchical variables. They start by independently computing the univariate distributions of all columns of the original data; then, they sample a user-defined amount of synthetic tables from these distributions, which are the population of the first generation. The GA process then optimizes the statistics of the data sets iteratively until the desired similarity to the original data is reached. The computational workload of table GAs is significantly higher, and they are more error-prone than GAs for string data due to the increase in variables and their relationships.

## 2.6 Boltzmann Machines

A Boltzmann machine is an undirected network consisting of binary visible nodes $\mathbf{v} \in \{0, 1\}^D$ and hidden nodes $\mathbf{h} \in \{0, 1\}^P$. The model parameters $\theta = \{\mathbf{W}, \mathbf{L}, \mathbf{J}\}$ are the visible-to-hidden, visible-to-visible and hidden-to-hidden symmetric interaction terms (matrices) of the graph (see Figure 7). The model parameters $\theta$ of a Boltzmann machine are trained using gradient ascent. They assign probabilities to the visible units, where the training data is put in, based on the states of the hidden units. [60]
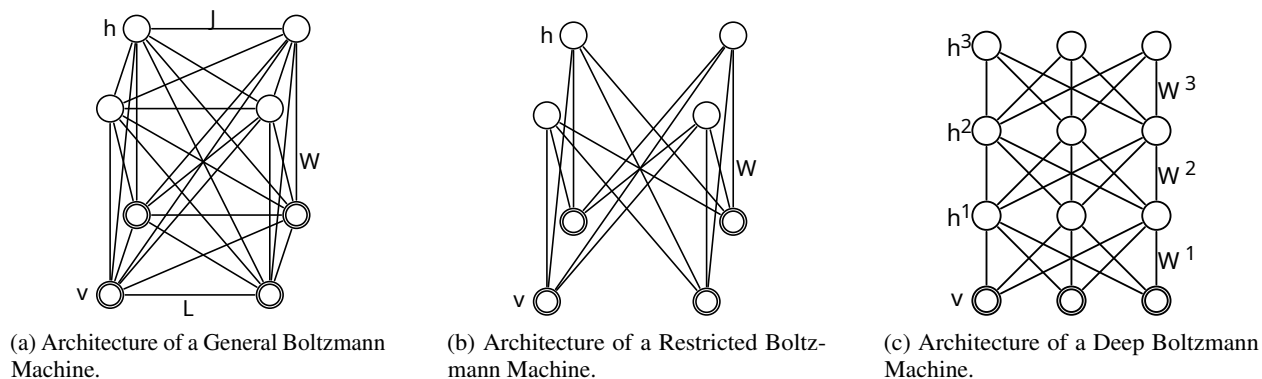


(a) Architecture of a General Boltzmann Machine.

(b) Architecture of a Restricted Boltzmann Machine.

(c) Architecture of a Deep Boltzmann Machine.

Figure 7: Architectures of different kinds of Boltzmann machines. (Source: [60])

### 2.6.1 Restricted Boltzmann Machines

Restricted Boltzmann Machines (RBMs) [61] are a subset of general Boltzmann machines, where only visible-to-hidden ($\mathbf{W}$) connections are allowed, so both $\mathbf{J}$ and $\mathbf{L}$ are set to zero. These models have the advantage that inference is exact, and learning is significantly more efficient [60]. The method used to train a RBM is unsupervised learning, so the data is unlabeled. Each training sample is provided as input $\mathbf{v}$ and $\theta$ is modified to increase the likelihood function $p(\mathbf{v}, \theta)$ using, for example, a gradient method or Contrastive Divergence [62].

Lee et al. [63] extend the RBM with convolution filters to process two-dimensional high-resolution images. The visible input units of size $N_V \times N_V$ are processed by $K$ filters with size $N_W \times N_W$ to produce $K$ hidden layers with size

$N_H \times N_H$. Each hidden layer is then partitioned into $C \times C$ blocks that are each connected to exactly one binary unit in the *max-pooling layer* $P$ to shrink the representation. The architecture is depicted in Figure 8.
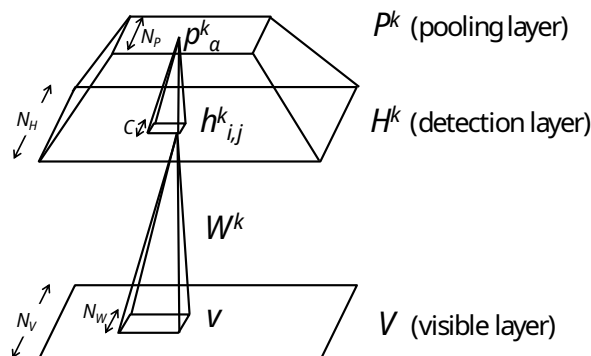


Figure 8: The convolutional RBM architecture with only one of the $K$ hidden and max-pooling layers shown for readability. (Source: [63])

Lattner et al. [64] propose a multi-component model to generate music with consistent local and global structural properties. First, a convolutional RBM [63] learns the local structure of musical pieces based on training data. Then, *constrained sampling* is applied to a randomly initialized "piano roll" music representation matrix $\mathbf{v} \in [0,1]^{T \times P}$ consisting of probabilities of active pitches $1 < p < P$ over time steps $1 < t < T$: The sampling process first applies 20 gradient descent steps to $\mathbf{v}$ with a loss function involving self-similarity, tonality and meter constraints regarding a template piece $\mathbf{x} \in [0,1]^{T \times P}$ before performing alternating steps of Gibbs sampling[2] with the convolutional RBM and one step of gradient descent with lower learning rate to $\mathbf{v}$. This process is repeated until the solution no longer improves on the RBM and constraints jointly.

### 2.6.2 Deep Belief Networks

A Deep Belief Network (DBN), similar to its successor Deep Boltzmann Machine (DBM) (see Section 2.6.4), is a combination of multiple RBMs, where the hidden layer of one RBM becomes the input (visible layer) of the higher-level RBM, but only the top two hidden layers are undirected and form an *associative memory* while the lower layers form a DAG (see Figure 9). Like the RBM, the DBN can be used to learn high-level representations of unlabeled data and convert representations back to visible data (e.g., images), but it can also be extended to supervised learning tasks by appending the label as an input to the visible layer. [65]



(a) Architecture of a Deep Belief Network.

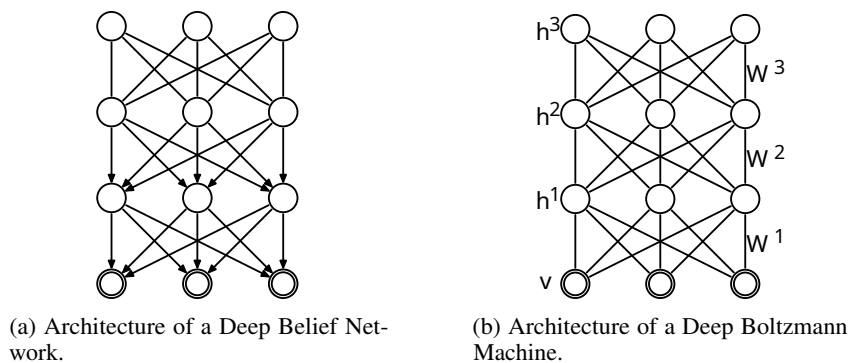(b) Architecture of a Deep Boltzmann Machine.

Figure 9: Comparison of the architectures of a DBN and a DBM. (Source: [60])

The layer-wise greedy learning process starts by training the lowest-level RBM with the visible and first hidden layer normally. Then, the output of the first hidden layer becomes the input (visible layer) of the next RBM, and now the second model is trained. This continues until a sufficient number of layers are reached. The second step introduces a

---

[2]Gibbs sampling is a Markov chain Monte Carlo (MCMC) algorithm used for approximate statistical inference that iteratively samples from the conditional probability distribution of each variable in a multivariate distribution while holding the others fixed.

fine-tuning algorithm with a bottom-up and a top-down pass. Now, the "recognition" weights $\mathbf{W}_x^\top$ and "generative" weights $\mathbf{W}_x$ are decoupled and modified independently. During the bottom-up pass, fixed recognition weights are used to stochastically determine all hidden values and update the generative weights on the directed connections according to a likelihood metric. The top-down pass starts with a state of the top-level associative memory and uses fixed generative weights to determine the visible layers on which the recognition weights are updated similarly. [65]

Hinton et al. [65] use alternating Gibbs sampling [66] in the DBN's associative memory "until the Markov chain converges to the equilibrium distribution". Then, images of digits are generated in the DBN that was trained on the MNIST dataset by drawing a sample from this distribution and passing it down the generative weights. By fixing the label units, certain digits can be synthetically generated, and with repeated iterations of Gibbs sampling between down-passes, the digits become more realistic.

Lee et al. [63], who introduced convolutional RBMs, stack these on top of one another to create convolutional DBNs. Different from [65], they use undirected connections between all layers, like DBMs. The experimental results with a two-layer convolutional DBN on the Kyoto natural image data set show that the first layer learns edge filters and the second filters for contours, corners, angles, and surface boundaries. The hierarchical representations obtained from these layers improve classification results on the Caltech-101 object and MNIST digit classification tasks.

Bickerman et al. [67] apply DBNs to create jazz melodies in an unsupervised manner. They divide each beat into 12 slots, each consisting of 30 bits (12 chords, 18 melodies) that encode the note pitch, octave, and length. A two-layer DBN can produce short stylistically plausible jazz samples based on a random sequence input but cannot capture regularities in music. Other problems are the large training time and complicated sampling procedure of DBNs.

Sun [68] employs two DBNs (one flipped) with pair-wise pre-trained layers as an autoencoder to generate random piano roll bars of music from binary piano roll matrices (rows represent note pitches, columns a 16th note of playing time) of complete or incomplete music. On average, the network copies 56,9% of the notes during the reconstruction, which still results in noticeably different works being created.

### 2.6.3 Temporal Restricted Boltzmann Machines and Related Models

The Temporal Restricted Boltzmann Machine (TRBM) [69] is a sequence of RBMs where the biases of one RBM depend on the hidden state of the previous RBM (see Figure 10). Training the model works similar to a normal RBM, but on sequences instead of fixed-size samples. A significant problem of TRBMs is that for computing probabilities during inference, the evaluation of all possible states (partition functions) of two RBMs is required, making a heuristic approach necessary [70].
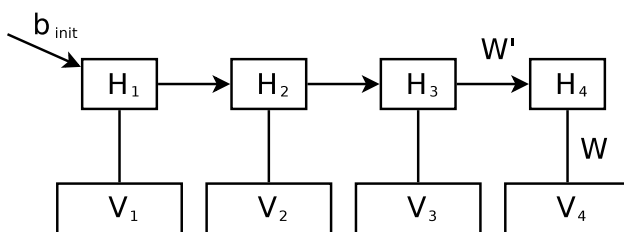


Figure 10: Architecture of a TRBM. (Source: [70])

Recurrent Temporal Boltzmann Machines (RTRBMs) [70] are slightly modified TRBMs which offer exact inference and feasible gradient computation. RTRBMs can be expressed as a RNN with the same parameters as the RTRBM and its log-likelihood as the cost function, so gradients can be computed using the Backpropagation Through Time (BPTT) algorithm.

The Structured Recurrent Temporal Boltzmann Machine (SRTRBM) [71] learns a dependency structure between pairs of visible and hidden units instead of using full connectivity like the RTRBM. The model encourages sparse graphs and can reveal the structure of the underlying time-series data.

Sutskever et al. [69], the original authors of the TRBM, train their model on 10,000 video sequences with 100 frames and achieve good results on future frame prediction and online denoising (removing artifacts). In their later work on RTRBMs [70], they improve on the video generation task.

Mittelmann et al. [71] compare their SRTRBM to the previously proposed models by Sutskever et al. [69, 70] by predicting frames on synthetic bouncing ball videos. Further, they use the SRTRBM to make predictions on motion capture and weather data. They improve on the performance of the predecessors in all temporal modeling tasks.

### 2.6.4 Deep Boltzmann Machines

A DBM [60] is a combination of multiple RBMs where the hidden layer of one RBM becomes the visible layer of the next one, resulting in a model with one visible and a stack of multiple hidden layers (see Figure 9). Using a greedy layer-wise learning approach, such a multi-layer model can efficiently learn high-level representations from unlabeled data.

Salakhutdinov et al. [60] use a DBM trained on the MNIST dataset to generate synthetic handwritten digits by initializing the model with random binary states and running a Gibbs sampler for 100,000 steps. Further, they demonstrate the creation of greyscale toy images on the NORB [72] dataset and improve the results by increasing the amount of training data with simple pixel translations.

### 2.6.5 Gated Boltzmann Machines

A gated Boltzmann machine encodes transformations between two observations using its hidden layer. It can be described as a conditional RBM with a visible input $\mathbf{x}$, a hidden transformation representation layer $\mathbf{h}$ that acts as a *gate* (see Figure 11) and a visible output layer $\mathbf{y}$. Conditional on the input, the inference and learning procedures are tractable [73]. The trainable parameters are stored in a three-dimensional parameter "tensor" $\mathbf{W}$ and the compatibility between $\mathbf{x}$, $\mathbf{y}$ and $\mathbf{h}$ is computed by an energy function $E(\mathbf{y}, \mathbf{h}; \mathbf{x})$ [74]:

$$E(\mathbf{y}, \mathbf{h}; \mathbf{x}) = -\sum_{ijk} W_{ijk} x_i y_j h_k \tag{6}$$

The values for $\mathbf{y}$ and $\mathbf{h}$ can be computed in the same way [74]:

$$p(h_k|\mathbf{x}, \mathbf{y}) = \frac{1}{1 + \exp(-\sum_{ij} W_{ijk} x_i y_j)} \tag{7}$$

$$p(y_j|\mathbf{x}, \mathbf{h}) = \frac{1}{1 + \exp(-\sum_{ik} W_{ijk} x_i h_k)}. \tag{8}$$



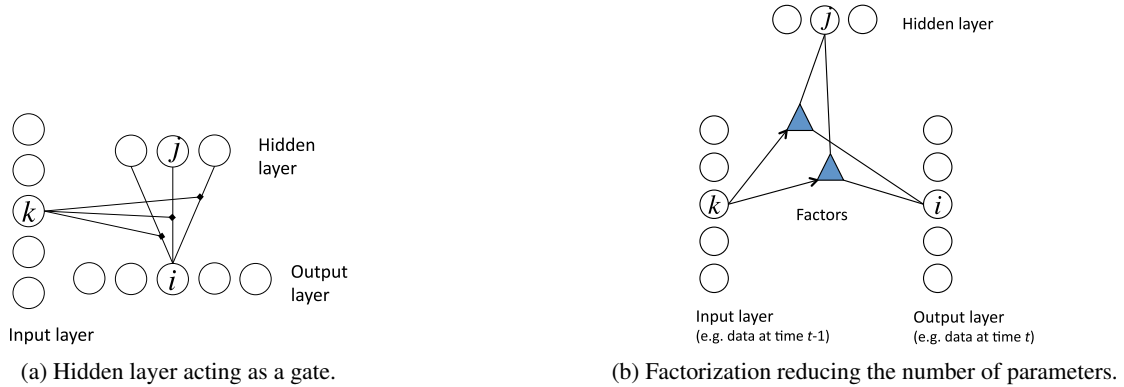(a) Hidden layer acting as a gate.　　　　(b) Factorization reducing the number of parameters.

Figure 11: Gate approaches in a gated Boltzmann machine. (Source: [73])

The model parameters $\mathbf{W}$ are trained with gradient-based optimization maximizing the average conditional log-likelihood $L = \frac{1}{N} \sum_{\alpha} \log p(\mathbf{y}^{\alpha}|\mathbf{x}^{\alpha})$ for training pairs $(\mathbf{x}^{\alpha}, \mathbf{y}^{\alpha})$. Because parts of the gradient are intractable, Gibbs sampling is used to approximate partial results with the help of the conditional distributions described in Equation 7 and Equation 8 in a scheme called contrastive divergence. [74]

Memisevic et al. [74] use the gated Boltzmann machine to learn transformation representations on randomly transformed $8 \times 8$ pixel images and predict the next images on video patches of size $22 \times 22$ pixels. They notice that the model becomes intractable for larger images due to the cubic parameter space $\mathbf{W}$ and modify it to make it iteratively applicable to smaller image patches. Finally, the model's abilities are demonstrated by analogy making, which means obtaining a transformation from a source image and applying the transformation to a target image using the patch-wise approach.

Memisevic et al. [75] tackle the problem of the rapidly expanding cubic parameter tensor $\mathbf{W}$ for large inputs by approximating the results using three matrices $w^x$, $w^y$ and $w^h$, so $w_{ijk} = \sum_{f=1}^{F} w_{if}^x w_{jf}^y w_{kf}^h$. Suppose the number of

factors $F$ is similar to the number of hidden and visible variables. In that case, the model now only requires $O(N^2)$ instead of $O(N^3)$ parameters (see Figure 11 for a comparison). Like [74], the model learns filters as transformation representations on larger $40 \times 40$ pixel images. Also, the analogy experiments are repeated, and motion extraction is successfully performed.

Taylor et al. [73] propose three key properties for time series models:

1. Distributed (i.e., componential) hidden state instead of sampling from a single category like HMMs to retain high flexibility and capacity.

2. Undirected, bipartite graph as the model structure to make inference simple and efficient.

3. Ability to form deep networks by stacking models and learning one layer at a time to capture more abstract data features.

Guided by these constraints, they introduce the Conditional Restricted Boltzmann Machine (CRBM), which takes one or more visible layers from previous time steps as additional inputs to a normal RBM (see Figure 12a), and the Conditional Deep Belief Network (CDBN), which stacks multiple CRBMs similar to a DBN, to build a generative model for time series initialized by real data (see Figure 12b). Further, the gated CRBM is introduced to enable multiplicative interactions between time steps, allowing the learned transformations to be highly nonlinear. Also, the factorization method from [75] is reintroduced to reduce parameters, and predefined style labels are added to the gate process. The models are successfully evaluated using the CMU motion data set where motions are continued by the CRBMs or new motions with mixed or changing styles are created.
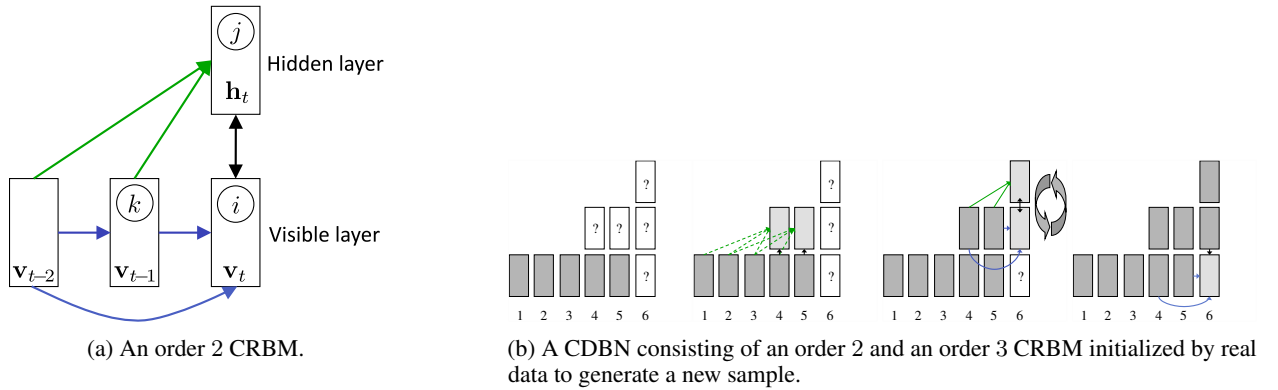


(a) An order 2 CRBM.

(b) A CDBN consisting of an order 2 and an order 3 CRBM initialized by real data to generate a new sample.

Figure 12: Architecture of CRBM and the deep version CDBN. (Source: [73])

## 2.7 Autoencoders

The basic autoencoder is a network that has as input a vector $\mathbf{x} \in [0, 1]^d$ and maps it to a hidden representation $\mathbf{y} \in [0, 1]^{d'}$ with a function $\mathbf{y} = f_\theta(\mathbf{x}) = s(\mathbf{Wx} + \mathbf{b})$. $\mathbf{W}$ is a weight matrix of size $d' \times d$, $\mathbf{b}$ a bias vector and together they are the parameters $\theta = \{\mathbf{W}, \mathbf{b}\}$ of $f$. $s(x) = \frac{1}{1+e^{-x}}$ is the sigmoid function. The hidden representation is then mapped back to a vector $\mathbf{z} \in [0, 1]^d$ where $\mathbf{z} = g_{\theta'}(\mathbf{y}) = s(\mathbf{W'y} + \mathbf{b'})$ with $\theta' = \{\mathbf{W'}, \mathbf{b'}\}$. Optionally, the constraint $\mathbf{W'} = \mathbf{W}^\top$ can be applied. [76]

During the unsupervised representation learning, the autoencoder adapts its parameters to minimize the *average reconstruction error* for training samples $\mathbf{x}^{(i)}$ and corresponding $\mathbf{y}^{(i)}$ and $\mathbf{z}^{(i)}$:

$$\theta^*, \theta'^* = \arg\min_{\theta, \theta'} \frac{1}{n} \sum_{i=1}^{n} L(\mathbf{x}^{(i)}, \mathbf{z}^{(i)}) = \arg\min_{\theta, \theta'} \frac{1}{n} \sum_{i=1}^{n} L(\mathbf{x}^{(i)}, g_{\theta'}(f_\theta(\mathbf{x}^{(i)}))). \tag{9}$$

The loss function $L$ can be anything compatible with the data, for example, the *squared error* $L(\mathbf{x}, \mathbf{z}) = \|\mathbf{x} - \mathbf{z}\|^2$. Similar to DBNs and DBMs, the hidden layer of one autoencoder can become the input layer of another one to learn higher-level representations or pre-train the weights to generate a neural network from them later. [76]

RBMs and autoencoders are very similar in structure (RBMs are undirected, while autoencoders are usually directed graphs), but they differ in training procedure and hidden representation: The autoencoder considers the real-valued

mapping from the input as its representation, while the Boltzmann machine samples a binary representation from that real-valued mapping. Autoencoders can also be seen as deterministic feedforward neural networks, while RBMs can be defined as the generative stochastic variant. [77, 78]

Short overview of other usages of autoencoders:

| Approach | Description | Year |
|---|---|---|
| [79] | Deep Autoregressive Network (DARN): A deep autoencoder with a mixture of stochastic and deterministic hidden units that incorporate autoregressive connections in the same layer, so $p(h) = \prod_{j=1}^{n_h} p(h_j \| h_{1:j-1})$. Outperforms RBM and NADE on image log-likelihood. | 2014 |
| [78] | The first evaluation of one and two-layer autoencoders on music audio spectrograms. | 2014 |
| [80] | An hierarchical Long Short-Term Memory (LSTM) autoencoder with an attention mechanism that builds and reconstructs embeddings of words, sentences, and paragraphs. The model is capable of coherent multi-sentence generation. | 2015 |
| [81] | Using so-called Ladder Networks, a semi-supervised learning model combining supervised and unsupervised learning in deep neural networks. It enhances learning by adding denoising tasks at each level of the model, fostering more robust feature learning. The model's effectiveness is demonstrated on several benchmark datasets, where it achieves state-of-the-art performance in semi-supervised learning tasks. | 2015 |
| [82] | Extending several state-of-the-art network architecture approaches by introducing auxiliary variables to deep generative models, which improve variational distribution approximation. | 2016 |
| [83] | DEFactor: Conditional molecule generation with optimal properties using a graph convolutional network [84] as an encoder to produce a latent graph representation, a LSTM creates a sequence of node embeddings autoregressively, and a decoder determines the edge and node types using a similarity measure of the node embeddings. Additionally, an existence module MLP stops the LSTM generator when a non-informative embedding is encountered to generate graphs of arbitrary size. | 2018 |
| [85] | Adversarial Generator-Encoder Network (AGE): A generator creates data from a specified latent distribution $z \sim p(z)$ and an adversarial encoder converts real and generated data to latent vectors. The AGE learns by comparing the distributions of the real and fake latent vectors with $p(z)$. The model is suitable for conditional and unconditional generation, converges quickly, and does not require a discriminator. | 2018 |
| [86] | Augmenting time series data through time-warped autoencoders. The authors introduce two techniques - independent and dependent - and showcase their effectiveness in producing synthetic data samples. This method utilizes the characteristics of auto-encoders to create high-quality, realistic time series data. | 2021 |

### 2.7.1 Helmholtz Machines

The Helmholtz machine [87] is a directed deep generative model [88] with binary units, biases, and generative and recognition weights for the respective direction (see Figure 13). The weights are trained with the "wake-sleep" algorithm [89]: During the "wake" (recognition) phase, the input values are propagated bottom-up through the layers to create a representation. The state $s_v$ of unit $v$ with bias $b_v$ and incoming weights $w_{uv}$ is defined as

$$p(s_v = 1) = \frac{1}{1 + \exp(-b_v - \sum_u s_u w_{uv})}. \tag{10}$$

Then the generative weights of the hidden states $s_i^\alpha$ between layers $k$ and $j$ are updated with the delta rule $\Delta w_{kj} = \epsilon s_k^\alpha (s_j^\alpha - p_j^\alpha)$ and learning rate $\epsilon$.

The "sleep" phase generates the states of the lower layers from the highest layer in a top-down approach. It updates the recognition weights of all states with another delta rule $\Delta w_{jk} = \epsilon s_j (s_k - q_k)$ where $q_k$ is the probability that unit $k$ is activated by the recognition weights of the states $s_j$ of the lower layer.

### 2.7.2 Denoising Autoencoder

Denoising Autoencoders (DAEs) modify the training process of basic autoencoders by randomly corrupting the input vector $\mathbf{x}$ (i.e., setting a fixed proportion of values to zero), resulting in $\tilde{\mathbf{x}}$, and then trying to restore it with the transformations $f$ and $g$ (see Figure 14). The data corruption improves the model's generalization abilities and, therefore, the representations' robustness. Also, the constraint $d' < d$ of the basic autoencoder to prevent overfitting can now be omitted. [76]
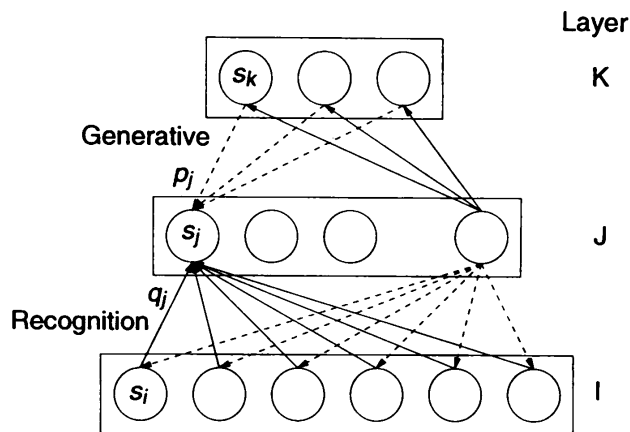
Figure 13: Architecture of a Helmholtz machine with three layers K, L, and I and probabilities $p$ for generation and $q$ for recognition. (Source: [89])
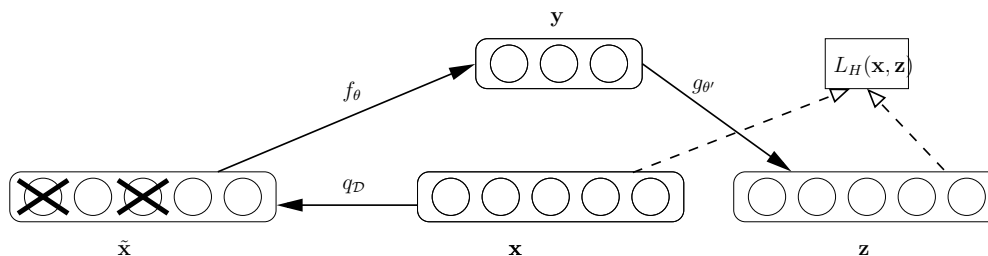


Figure 14: Training process of a Denoising Autoencoder. (Source: [76])

Bengio et al. [90] propose a generalized probabilistic interpretation of DAEs, where an observed random variable $X$ is corrupted using a conditional distribution $\mathcal{C}(\overline{X}|X)$ and the DAE is trained to estimate the reverse conditional $P_\theta(X|\overline{X})$, where $\theta$ are the trainable parameters. The sampling process from this DAE is realized using a Markov chain where $X_t \sim P_\theta(X|\overline{X}_{t-1})$ and $\overline{X}_t \sim \mathcal{C}(\overline{X}|X_t)$, which the authors prove generates the data-generating distribution $P(X)$ with a properly trained model. They introduce *walkback* training, where the default corruption process $\mathcal{C}$ is replaced with a walkback process $\overline{\mathcal{C}}$, which generates one or more boosted "negative examples" $\overline{X}^*$ for a training sample $X$ by going a random-length walk through the aforementioned Markov chain with the current model parameters. The training with these greatly divergent corruptions, which may not even be represented by the training data or simple corruptions, prevents the DAE from deviating too far from the plausible prediction range. The experiments on the MNIST dataset (see Figure 15) show that the results of the Markov chain of the walkback-trained model look more natural than the ones obtained from the model trained with a simple corruption process.



(a) Training with simple corruption.

(b) Training with walkback.

Figure 15: Comparison of the MNIST results of a DAE trained normally or with walkback. (Source: [90])

### 2.7.3 Contractive Autoencoder

The Contractive Autoencoder (CAE) uses the constraint $\mathbf{W'} = \mathbf{W}^\top$ and improves the robustness of the hidden representation $\mathbf{y}$ by penalizing sensitivity to the input $\mathbf{x}$ with the Frobenius norm of the Jacobian matrix $J_f(\mathbf{x})$, which is the sum of squares of all partial derivatives of the extracted features $\mathbf{y} = f(\mathbf{x})$ concerning the input $\mathbf{x}$ [91, 92]:

$$J_f(\mathbf{x}) = \frac{\partial f(\mathbf{x})}{\partial \mathbf{x}} \tag{11}$$

$$\|J_f(\mathbf{x})\|_F^2 = \sum_{ij} (\frac{\partial f_j(\mathbf{x})}{\partial \mathbf{x}_i})^2. \tag{12}$$

In the case of the sigmoid activation function being used, this results in the following:

$$\|J_f(\mathbf{x})\|_F^2 = \sum_{i=1}^{d'} (\mathbf{y}_i(1 - \mathbf{y}_i))^2 \sum_{j=1}^{d} \mathbf{W}_{ij}^2. \tag{13}$$

This objective function is minimized on the data $D_n$, with hyperparameter $\lambda \geq 0$, parameters $\theta = \{\mathbf{W}, \mathbf{b}, \mathbf{b'}\}$ and cross-entropy loss $L(\mathbf{x}, \mathbf{z})$:

$$\mathcal{J}_{CAE}(\theta) = \sum_{\mathbf{x} \in D_n} \left(L(\mathbf{x}, g(f(\mathbf{x}))) + \lambda\|J_f(\mathbf{x})\|_F^2\right) \tag{14}$$

$$L(\mathbf{x}, \mathbf{z}) = -\sum_{i=1}^{d} \mathbf{x}_i \log(\mathbf{z}_i) + (1 - \mathbf{x}_i) \log(1 - \mathbf{z}_i) \tag{15}$$

Rifai et al. [92] propose an algorithm (see Algorithm 1) to generate samples from a pre-trained CAE that provides an ergodic Harris chain with a stationary distribution $\pi$ under the condition that $J_t J_t^\top$ is full rank.

---

**Algorithm 1** Sampling algorithm for a CAE. (Source: [92])

---

**Require:** $f$, $g$, step size $\sigma$ and chain length $T$
**Ensure:** Sequence $(x_1, y_1), (x_2, y_2),...,(x_T, y_T)$
    Initialize $x_0$ arbitrarily and $y_0 = f(x_0)$.
    **for** $t = 0$ **to** $T$ **do**
        Let Jacobian $J_t = \frac{\partial f(x_t)}{\partial x_t}$.
        Let $\epsilon \sim N(0, \sigma I_k)$ isotropic Gaussian noise.
        Let perturbation $\Delta y = J_t J_t^\top \epsilon$.
        Let $x_t = g(y_{t-1} + \Delta y)$ and $y_t = f(x_t)$.
    **end for**

---

They test their technique on a CAE with two-layer stacks and compare it against a 2-layer DBN, resulting in slightly better/worse performance on the Toronto Face Database (TFD)/MNIST dataset regarding the log-likelihood of the generated samples but significantly reduced sensitivity to deformations of MNIST digits.

Bengio et al. [93] hypothesize that sampling from higher-level representations improves the quality (log-likelihood) and class variation of obtained samples. They prove their claims by sampling from high-level representations generated by a Markov chain process from various depths of multi-layer CAEs and DBNs and measuring the log-likelihood. Further, they test the mixing of representations of digits at various depths with linear interpolation, which also gives more plausible samples at higher levels.

### 2.7.4 Generative Stochastic Network

A Generative Stochastic Network (GSN) [94] is a generalized framework of a deep DAE that has the structure of a Markov chain and can be trained with back-propagated gradients and without layer-wise pre-training. The transition operator $P(x_t, h_t|x_{t-1}, h_{t-1})$ is responsible for generating the next visible state $X_t$ and hidden state $H_t$ of the Markov chain (see Figure 16).

To enable the back-propagation of the reconstruction log-likelihood $\log P(X_1 = x_0|H_1)$ into all the parameters of the encoding function $f_{\theta_1}$ and reconstruction function $g_{\theta_2}$, a deterministic function is used to define $H_{t+1} = f_{\theta_1}(X_t, Z_t, H_t)$ with $Z_t$ being an independent noise source so $X_t$ cannot be exactly recovered from $H_{t+1}$. This resembles the masking of values in the input layer of a DAE.
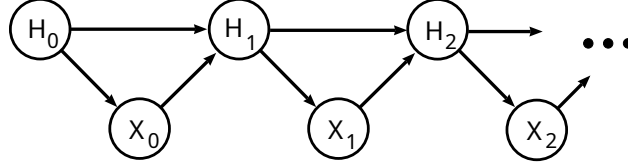
Figure 16: The Markov chain structure of a GSN. (Source: [94])

Bengio et al. [94] use the GSN framework to adapt the Gibbs sampling process of a DBM (see Figure 17), but with the ability to use the GSN's backpropagation at each layer. The chain starts with a training sample $X = x_0$ and encodes and reconstructs intermediate samples $x_i$ several times. The training of the model is realized using the sum of all log-likelihoods to the target $X = x_0$, inspired by the *walkback* objective (see [90]).
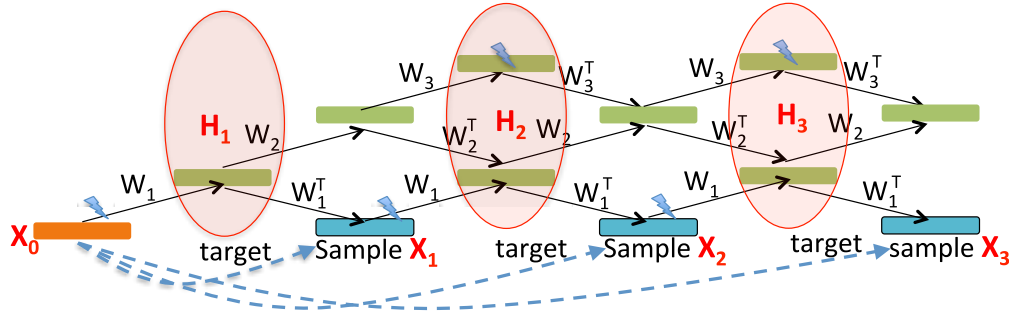


Figure 17: A GSN inspired by the Gibbs sampling process of a DBM. The lightning symbols indicate the corruption of the samples with salt-and-pepper noise. (Source: [94])

### 2.7.5 Variational Autoencoder

Variational Autoencoders (VAEs) is an autoencoder with encoder $\mathcal{E}(.)$ and decoder $\mathcal{D}(.)$ whose hidden layer is represented as a Gaussian distribution $\mathcal{N}(\mu, \sigma^2)$ with mean vector $\mu$ and standard deviation vector $\sigma$ that are obtained from the encoder $\mu, \sigma = \mathcal{E}(x)$, where $x \sim \mathbf{X}$. The decoder samples from this distribution and reconstructs the data $\hat{x} = \mathbb{E}_{z \sim \mathcal{N}(\mu, \sigma)}[\mathcal{D}(z)]$. Additionally, the aggregated distribution of $z$ over all data $\mathbf{X}$ is constrained to be $\mathcal{N}(0, \mathbf{I})$, allowing random vectors to be sampled from $\mathcal{N}(0, \mathbf{I})$ to be used for data generation with the decoder. [95]

The VAE is trained using the Evidence Lower-Bound (ELBO) loss

$$\mathcal{L} = \mathbb{E}_{x \sim \mathbf{X}}[\mathbb{E}_{z \sim \mathcal{N}(\mu, \sigma \mathbf{I})}[\|\mathcal{D}(z) - x\|_2^2] + \mathbb{KL}(\mathcal{N}(\mu, \sigma \mathbf{I})\|\mathcal{N}(0, \mathbf{I}))], \tag{16}$$

where the first term encourages the autoencoding part while the second term with the Kullback-Leibler (KL) divergence $\mathbb{KL}(p\|q)$ measures the difference between the probability distributions $p$ and $q$. [95]

Kingma et al. [96] applied the VAE as a generative model with MLP encoder and decoder to the MNIST and Frey Face datasets. They achieve faster convergence and a higher marginal likelihood than their reference, the wake-sleep algorithm, described in Section 2.7.1.

Gregor et al. [97] introduce the *Deep Recurrent Attentive Writer* (DRAW) for image generation. The authors follow the human drawing process, where rough outlines are iteratively refined until a realistic picture is generated. Similar to a VAE, DRAW consists of an encoder that learns a representation $\mathbf{z} \sim p(\mathbf{z})$ of the input and a decoder that reconstructs the input from $\mathbf{z}$, but both models are LSTMs that only handle regions of the full input defined by an attention mechanism at each time step. For evaluation, the authors use MNIST, the Street View House Numbers (SVHN) data set, and CIFAR-10, and they achieve highly realistic results on the first two while overfitting the last data set due to only 50,000 training samples.

Rezende et al. [98] propose the class of *sequential generative models*, which generate $T$ groups of $k$ latent variables sequentially instead of generating $K = kT$ latent variables at once and also sequentially reconstruct data to a modifiable canvas while incorporating inference and writing attention mechanisms. At the core of the model are one or more RNNs (both LSTMs and Gated Recurrent Units (GRUs)) and attention-based neural networks like spatial transformers [99] to write the RNN output to the canvas. The model performs well on unconditional sampling and *one-shot learning* tasks,

where a concept is only encountered once and compelling variations of the concept should be generated, or multiple concepts (e.g., letters of an alphabet) are provided during the sequential inference process, and a plausible new character is generated.

Higgins et al. [100] propose $\beta$-VAE, a reformulation of the original unsupervised VAE ($\beta = 1$) with an additional hyperparameter $\beta$ that encourages the model to learn better-disentangled representations of data, meaning that single hidden units encode single generative factors while being invariant to changes in others (e.g., skin color in face images), by enforcing more independence and less covariance of the latent variables and their distributions. $\beta$-VAE outperforms the previous unsupervised state-of-the-art model InfoGAN [101] and semi-supervised DC-IGN [102] in terms of disentanglement of factors qualitatively and quantitatively.

Tomczak et al. [103] implement the VAE's prior as a mixture distribution (e.g., a GMM) that can be trained with pseudo inputs and propose multiple layers of variables. They evaluate their model on multiple MNIST data sets, OMNIGLOT, Caltech 101 Silhouette, Frey Face, and Histopathology patches, resulting in state-of-the-art log-likelihood compared to a normal VAE and avoidance of its local optima problem.

Short overview of other usages of VAEs:

| Approach | Description | Year |
|---|---|---|
| [104] | Variational Recurrent Autoencoder (VRAE): A VAE with RNN encoder and decoder for modeling sequential data. The model is used for MIDI music generation and creates "medleys" of the training data. | 2014 |
| [105] | Denoising Variational Autoencoder (DVAE): Combination of the noise injection at the input like a DAE and the noise injection at the hidden layer of a VAE improves average log-likelihood results on MNIST and Frey Face data sets. | 2015 |
| [106] | Importance-weighted autoencoder (IWAE): A VAE with a tighter log-likelihood lower bound on $\log p(\mathbf{x})$ based on importance weighting using multiple samples $q(\mathbf{h}_i|\mathbf{x})$, that learns richer representations than VAEs, outperforming them on MNIST density estimation. | 2015 |
| [107] | A VAE with single-layer LSTM encoder and decoder for sentence generation. Interpolation between latent vectors of two sentences provides interesting results. | 2015 |
| [102] | Deep Convolutional Inverse Graphics Network (DC-IGN): A VAE built with a convolutional encoder and deconvolutional decoder. The model is further encouraged to assign certain transformations (e.g., lighting, rotation) in images to disentangled neuron groups unsupervised by training with mini-batches of transformed images. | 2015 |
| [108] | Composited spatially transformed VAE (CST-VAE): Layer-wise sequential image generation using pose and content encoder-decoder pairs on the partial results and a spatial transformer network [99] to output the next image layer front to back. | 2015 |
| [109] | VAE training with an additional similarity loss obtained from a jointly trained GAN discriminator for higher image quality and better representations. | 2016 |
| [110] | Attribute2Image: Using a VAE to learn disentangled latent representations of attributes (e.g., color, gender, viewpoint) from images and generate new images conditioned on attributes. | 2016 |
| [111] | A VAE with recurrent encoder and decoder (similar to [107]) is used to generate SMILES [112] text encodings of new valid molecules with desirable properties by using gradient-based optimization in the latent space. | 2016 |
| [113] | Proposal of the variational lossy autoencoder (VLAE), which allows the user to control what the latent variable can contain by limiting the receptive field of the autoregressive encoder and decoder. Combined with an autoregressive model modeling $p(z)$, state-of-the-art results are achieved on MNIST, OMNIGLOT, and Caltech-101 density estimation (competitive on CIFAR-10). | 2016 |
| [114] | PixelVAE: Combination of a VAE with a PixelCNN-based [115] autoregressive decoder that iteratively refines the image result. It performs comparably to PixelCNN with fewer autoregressive layers and a smaller latent variable than a normal VAE. | 2016 |
| [116] | GrammarVAE: To better generate discrete data, it is converted to a parse tree using a context-free grammar, and the rules used by the tree are then one-hot encoded in order, formatted as a matric and mapped to a latent space with a Convolutional Neural Network (CNN). A RNN decodes from this latent space back to valid rules. The model generates molecule structures and arithmetic expressions, outperforming text-based representations. | 2017 |

... Continuation

| [117] | Recurrent hierarchical VAE with BiLSTM encoder and 3-layer LSTM decoder for the creative reconstruction of short musical sequences with random sampling or interpolation in the latent space. | 2017 |
|---|---|---|
| [118] | A VAE with LSTM encoder and CNN decoder with dilated convolutions for text modeling and generation. Changes in the dilation configuration give control over the context size from previous words, and the convolutional decoder is less prone to ignore encoder information because its contextual capacity is lower than an LSTM's. | 2017 |
| [119] | Vector-Quantised VAE (VQ-VAE): The encoder CNN outputs $z$ are discrete to enforce more efficient usage of the latent space and prevent "posterior collapse", which is often caused by decoders ignoring $z$. The model is combined with an autoregressive decoder (PixelCNN for images, WaveNet for audio) and provides similar results to continuous VAEs. | 2017 |
| [120] | Character-level text generation VAE with convolutional encoder combined with a deconvolutional decoder with recurrent output layer. The CNNs make VAE training easier and an additional cost function encourages reliance of the decoder on the latent vector. Experiments on Twitter tweet generation show more diverse and coherent samples than a LSTM-based VAE. | 2017 |
| [121] | Sketch-RNN: Sketch-conditional and unconditional stroke-based sketch generation with a sequence-to-sequence VAE with bidirectional RNN encoder and RNN decoder. Possible applications also include latent space interpolation and sketch completion. | 2017 |
| [122] | First application of a recurrent VAE [107] to music generation, providing a good balance between local and global structures. | 2017 |
| [123] | Training a GAN to generate and modify the latent code $z$ of an unconditional VAE with $z' = G(z, y_{attr})$ to satisfy specific properties enforced by $D(z', y_{attr})$. Allows zero-shot conditional generation and identity-preserving transformation (e.g., same face with different hair color) of data from an unconditional VAE model. | 2017 |
| [124] | Application of the grammar VAE [116] to molecule generation with desired properties by randomly sampling $10^6$ samples from $p(\mathbf{z})$ and iteratively encoding and decoding them with the autoencoder many times before filtering the results with neural network regression functions to get the best samples for the desired property. | 2018 |
| [125] | GraphVAE (GVAE): Generation of probabilistic fully connected graphs from which can be sampled. The convolutional encoder gets as input a graph $G = (A, E, F)$ and graph label vector $\mathbf{y}$, where $A$ is the adjacency matrix, $E$ the edge attribute tensor, and $F$ a node attribute matrix, and computes a latent representation $\mathbf{z}$. The deterministic decoder MLP takes $\mathbf{z}$ and $\mathbf{y}$ and outputs $\tilde{G} = (\tilde{A}, \tilde{E}, \tilde{F})$, containing the independent node and edge, edge class and node class *probabilities* respectively for graphs with a fixed maximum number of $k$ nodes. The model is demonstrated on molecule generation tasks and is only suitable for slightly larger graphs than the provided input. | 2018 |
| [126] | MusicVAE: Recurrent VAE with two-layer bidirectional LSTM encoder and hierarchical RNN decoder, which consists of a two-layer unidirectional LSTM conductor that creates subsequence embeddings from the latent vector, and a two-layer LSTM decoder RNN that produces the final subsequence output inside these separate embeddings. The model achieves promising results on MIDI music reconstruction tasks. | 2018 |
| [127] | Junction tree VAE (JT-VAE): Encoding and decoding molecules using two representations: A fine-grained graph connectivity representation obtained from a message-passing network [128] and a representation of a junction tree (also from a message passing network), that models a molecule as a composition of valid subgraph components and avoids the node-by-node generation of invalid intermediaries. | 2018 |
| [129] | Constrained Graph VAE (CGVAE): Using gated graph neural networks as encoder and decoder for molecular graph generation. The encoder maps a graph with a maximum of $N$ nodes to $N$ latent codes $\mathbf{z}_v$ conforming to a Gaussian distribution. The decoder initializes a graph with $N$ nodes and a "stop node" from the latent codes. Then, a loop starts where new edges are added and labeled, and node representations are updated with messages from neighbors (see [128]) until an edge to the stop node is created. Correct atom valency is always enforced to guarantee valid molecules. | 2018 |

Continuation...

... Continuation

| | | |
|---|---|---|
| [130] | Syntax-directed VAE (SD-VAE): Applying syntax and semantic constraints to the decoder of a grammar-based VAE similar to GrammarVAE [116] to enforce syntactically valid and semantically reasonable reconstruction and optimization of molecule structures and program code. | 2018 |
| [131] | Differentially private autoencoder-based generative model (DP-AuGM) & variational autoencoder-based model (DP-VaeGM): Autoencoders are trained using stochastic gradient descent with clipped gradients and noise injection [132]. For DP-AuGM, only the encoder makes confidential data "private". The VAE version trains one model for each class, samples $z$ from each model, and merges the decoded samples, which is less stable than the former approach. | 2018 |
| [133] | IntroVAE: Integrating the concepts of VAE and GAN into a single model that is both a generator and a discriminator. This model self-evaluates the quality of generated images and improves itself accordingly, offering stable training and high-resolution image synthesis. The approach combines the advantages of VAE and GAN without needing extra discriminators, simplifying the architecture and improving training efficiency. | 2018 |
| [134] | A multilevel VAE architecture for generating coherent and long text sequences. The encoder consists of a lower and higher-level CNN producing separate latent representations where the lower latent vector is additionally conditioned on the upper. The lower representation is then fed to a hierarchical LSTM decoder network with a sentence-level and word-level LSTM. The model performs well on conditional (title-to-paragraph) and unconditional text generation. | 2019 |
| [135] | Topic-guided VAE (TGVAE): Modeling the latent space of a VAE using a GMM prior distribution parametrized by a neural topic module that is powered by the bag-of-words representation of the text. A GRU encoder and decoder processed the text inputs to outputs. The model is used for text generation and summarization. | 2019 |
| [136] | A VAE for molecule generation that uses a graph-convolutional encoder, a MLP to create a "bag-of-atoms" (counts for certain atoms in the target reconstruction), and a graph-convolutional decoder that takes the latent representation and the atom bag to create an edge probability matrix from which a beam search generates a discrete output. | 2019 |
| [137] | NeVAE: Generating molecular graphs with variable size from a convolution-style VAE with variable-length latent representations (one for each node). In addition to the atoms' types and their bonds, the model also predicts their positions. The decoder can be further optimized with a gradient-based algorithm to maximize specific properties of the generated molecules. | 2020 |
| [138] | Node-Edge Disentangled VAE (NED-VAE): Using three convolutional sub-encoders (node, edge, and graph) and two deconvolutional sub-decoders (node and edge), both with access to the graph representation, to reconstruct the node and edge attributes from which a graph can be recreated. The model also enforces the disentanglement of latent factors of nodes, edges, and joint patterns. | 2020 |
| [95] | Tabular VAE (TVAE): Tabular data generation with a VAE using probability distributions to encode discrete and continuous values. | 2020 |
| [139] | HI-VAE: A VAE that can handle incomplete (missing values at random) and heterogeneous (mixed continuous and discrete) data by learning the influence of input variables on the latent code individually, using special distributions for discrete variables, and only using observed variables for recognition (i.e., replace missing with zero). | 2020 |
| [140] | Message Passing Graph VAE (MPGVAE): Building on top of the GraphVAE [125], the authors use message passing neural networks [128] for both encoder and decoder, where edge and node representations are alternatingly updated multiple times based on messages from their neighbors. | 2020 |
| [95] | Besides a conditional GAN, the authors introduce a tabular VAE for generating synthetic tabular data using . This method addresses challenges in modeling tabular data, which often contains a mix of discrete and continuous columns and may exhibit imbalances and non-Gaussian distributions. It employs mode-specific normalization and reversible data transformations to generate synthetic data effectively. | 2020 |
| [141] | Treating the molecule as a sequence of valid SMILES-encoded [112] fragments/components, this approach uses GRUs to encode a molecule to a latent vector with Gaussian distribution and decode it back. The whole training process incorporates *low-frequency masking*, which masks rarely encountered fragments with a mask token that is replaced during sampling by any of the suitable masked fragments with uniform probability to improve uniqueness. | 2020 |

### 2.7.6 Deep Latent Gaussian Models

A Deep Latent Gaussian Model (DLGM) is a deep, directed generative model powered by Gaussian latent variables. A recognition model with $L$ layers encodes the training observations to provide the layer-wise parameters for the Gaussian distributions of the generation model. The generative process starts at the top latent layer ($L$) and draws mutually independent Gaussian variables $\mathcal{E}_l \sim \mathcal{N}(\mathcal{E}_l|\mathbf{0}, \mathbf{I})$. Each layer $\mathbf{h}_l = T_l(\mathbf{h}_{l+1}) + \mathbf{G}_l\mathcal{E}_l$ below the top layer $\mathbf{h}_L = \mathbf{G}_L\mathcal{E}_L$ depends on the layer above, where $T_l$ are MLP transformations and $\mathbf{G}_l$ are matrices. The visible data $\mathbf{v} = \pi(\mathbf{v}|T_0(\mathbf{h}_1))$ is generated from a distribution $\pi$. The model is trained using stochastic backpropagation, i.e., by computing gradients through random variables. [142]

The original DLGM authors Rezende et al. [142] evaluate the generative abilities of a three-layer DLGM on MNIST, CIFAR-10, the Frey faces data set, and NORB. They also propose an imputation use case, where the Gaussian model fills in missing data in the SVHN, Frey faces, and MNIST data sets. Their experiments achieve realistic-looking results comparable to other contemporary approaches such as the RBM and DBN.

### 2.7.7 Gated Autoencoders

A Gated Autoencoder (GAE) is a conditional bi-linear model that learns to represent a linear transformation encoded as *mapping units* $\mathbf{m}$ between two observations $\mathbf{x}^{(1)}$ and $\mathbf{x}^{(2)}$ using parameter matrices $\mathbf{U}$, $\mathbf{V}$ and $\mathbf{W}$ with

$$\mathbf{m} = \sigma(\mathbf{W}(\mathbf{U}\mathbf{x}^{(1)} \cdot \mathbf{V}\mathbf{x}^{(2)})). \tag{17}$$

The mappings, since the GAE is a symmentric model, can be used to reconstruct $\mathbf{x}^{(1)}$ or $\mathbf{x}^{(2)}$ respectively based on the other one [143]:

$$\tilde{\mathbf{x}}^{(2)} = \mathbf{V}^T(\mathbf{U}\mathbf{x}^{(1)} \cdot \mathbf{W}^T\mathbf{m}) \tag{18}$$

$$\tilde{\mathbf{x}}^{(1)} = \mathbf{U}^T(\mathbf{V}\mathbf{x}^{(2)} \cdot \mathbf{W}^T\mathbf{m}). \tag{19}$$

Training works similarly to a DAE: Input pairs are independently corrupted and concatenated. Then backpropagation and gradient-based optimization are used to minimize the loss function, for example, the symmetric reconstruction error [143, 144]:

$$\mathcal{L} = \|\mathbf{x}^{(1)} - \tilde{\mathbf{x}}^{(1)}\|^2 + \|\mathbf{x}^{(2)} - \tilde{\mathbf{x}}^{(2)}\|^2 \tag{20}$$

Michalski et al. [143] model a time series as a sequence of transformations applied to its elements. During the *predictive training*, a pyramid of stacked GAEs is used to learn basic and higher-order representations of transformations between observation pairs and predict future observations recurrently with a constant highest-order transformation (see Figure 18). The autoencoders are initialized by $k$ *seed frames* corresponding to the $k$ layers of the pyramid and optimized using backpropagation through time. The model is compared against a RNN and a RBM in the prediction of chirps (sinusoidal waves with changing frequencies) and video frames of bouncing balls and objects of the NORB data set, outperforming them in terms of mean squared error.

### 2.7.8 Masked Autoencoders

Masked autoencoders set weights in their input-to-hidden or hidden-to-output weight matrices to zero, meaning there is no computational path between certain input and output units, which is necessary, for example, for autoregressive tasks, where future inputs must not be seen. The masking approach also directly applies to deep architectures (see Figure 19). [145]

Germain et al. [145] propose the masked autoencoder for distribution estimation (MADE), which computes the joint probability distribution $p(\mathbf{x})$ of data $\mathbf{x}$ autoregressively as $p(\mathbf{x}) = \prod_{d=1}^{D} p(x_d|\mathbf{x}_{<d})$ and therefore needs to mask the future inputs $x_d, ..., x_D$ for the respective steps. The model can also be used for sampling according to the calculated probabilities demonstrated by generating binary MNIST images.

### 2.8 Neural Autoregressive Distribution Estimators

The Neural Autoregressive Distribution Estimator (NADE) is inspired by the RBM, whose joint probability estimation of an observation is intractable and can also be interpreted as an autoencoder that assigns probabilities to binary
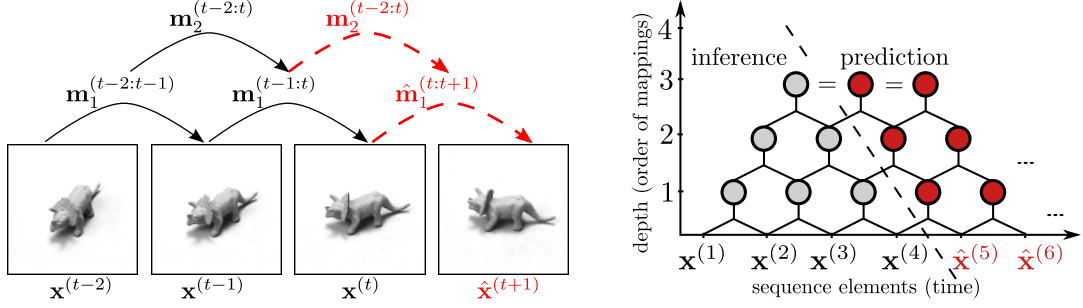
Figure 18: **Left:** A 2-layer pyramid model is used to predict the next transformation $\hat{\mathbf{m}}_1^{(t:t+1)} = \mathbf{V}_2^T(\mathbf{U}_2\mathbf{m}_1^{(t-1:t)} \cdot \mathbf{W}_2^T\mathbf{m}_2^{(t-2:t)})$ and the resulting observation $\hat{\mathbf{x}}^{(t+1)} = \mathbf{V}_1^T(\mathbf{U}_1\mathbf{x}^{(t)} \cdot \mathbf{W}_1^T\hat{\mathbf{m}}_1^{(t:t+1)})$ with $\mathbf{U}$, $\mathbf{V}$ and $\mathbf{W}$ being the filter matrices learned by the respective autoencoders. **Right:** Multi-step prediction with constant top-layer transformation in a 3-layer pyramid. (Source: [143])
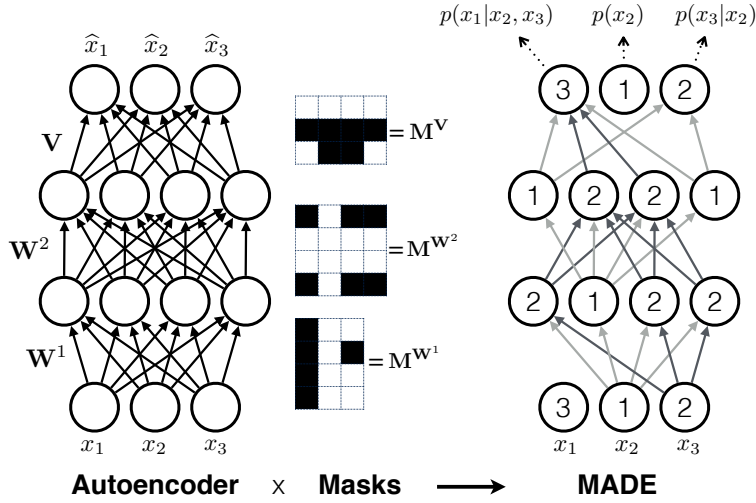


Figure 19: A deep masked autoencoder architecture. (Source: [145])

observations. NADE models a $D$-dimensional observation vector $\mathbf{x}$ as a product of the one-dimensional conditional distributions $p(\mathbf{x}) = \prod_{i=1}^{D} p(x_i|\mathbf{x}_{<i})$. The probability of $x_i$ is conditioned on all previously seen observations $\mathbf{x}_{<i}$, so the ordering of the variables matters. [146]

Each probability output $\hat{x}_i = p(x_i = 1|\mathbf{x}_{<i})$ depends on a $H$-dimensional hidden vector $\mathbf{h}_i$ that is computed recursively by

$$p(x_i = 1|\mathbf{x}_{<i}) = \sigma(\mathbf{V}_{.,i}\mathbf{h}_i + b_i) \tag{21}$$

$$\mathbf{h}_i = \sigma(\mathbf{W}_{.,<i}\mathbf{x}_{<i} + \mathbf{c}) \text{ and } h_1 = \sigma(\mathbf{c}) \tag{22}$$

with $\sigma$ being the sigmoid function and $\mathbf{V} \in \mathbb{R}^{D \times H}$, $\mathbf{b} \in \mathbb{R}^D$, $\mathbf{W} \in \mathbb{R}^{H \times D}$, and $\mathbf{c} \in \mathbb{R}^H$ being the parameters of the NADE model. This corresponds with each $\hat{x}_i$ being computed by a neural network and all neural networks having tied weights for each observation. [146]

NADE is trained using gradient descent on the Negative Log-Likelihood (NLL) given a training set $\mathbf{X}$ with size $T$ [146]:

$$\frac{1}{T}\sum_{t=1}^{T} -\log p(\mathbf{x}_t) = \frac{1}{T}\sum_{t=1}^{T}\sum_{i=1}^{D} -\log p(x_i|\mathbf{x}_{<i}) \text{ for } \mathbf{x}_i \in \mathbf{X} \tag{23}$$
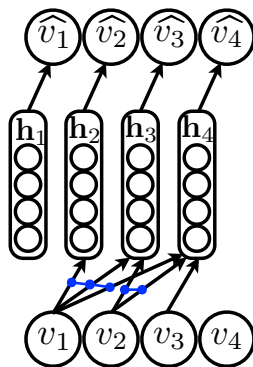
Figure 20: NADE architecture (blue lines are tied weights) with $v$ instead of $x$ notation. (Source: [146])

Uria et al. [147] propose real-valued NADE (RNADE), where real output values are computed using a GMM, so $p(x_i|\mathbf{x}_{<i}) = p_{GMM}(x_i|\theta_{\mathbf{i}})$ with $\theta_{\mathbf{i}}$ being the parameters of the GMM.

In [148], Uria et al. introduce an efficient procedure to train NADE and RNADE models for each possible variable ordering simultaneously by using shared weights and stochastic gradient descent to optimize the mean cost over all orderings. After that, the most suitable variable ordering for the data can be determined in constant time. They also introduce a deep NADE with multiple hidden layers that is efficient to train and often achieves better log-likelihood results on the test set than single-layer models.

Raiko et al. [149] propose NADE-$k$, which computed the density of an output as the $k$-th recurrent pass-through of the input $\mathbf{v}^{\langle 1 \rangle}$ through the neural network hidden layer, so $p(x_i = 1|\mathbf{x}_{\mathbf{obs}}) = v_i^{\langle k \rangle}$. The model outperforms previous NADE approaches, RBMs and DBNs in density modeling, also on masked inputs, and can generate binary MNIST digits and Caltech-101 silhouettes.

In [150], Uria et al. propose ConvNADE, which replaces the fully connected hidden layers with convolutional layers, allowing exploitation of the spatial structure, for example, of 2D images. They also combine the approach with the DeepNADE [148] architecture, which uses masking to improve the results of image modeling tasks.

### 2.9 Sparse Coding

Sparse coding is usually an optimization problem, where data is reconstructed by a weighted linear combination of as few as possible basis vectors (see Figure 21 for an example application). The reconstruction and sparsity costs of the linear combination representing the data have to be minimized. [151]
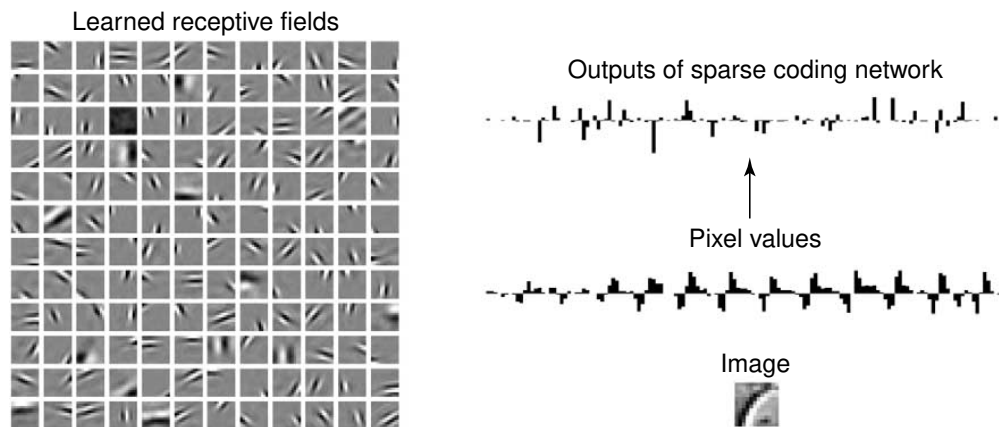


Figure 21: Example of an application of sparse coding to generate images. (Source: [152])

Wang et al. [153] propose using sparse representations of image patches for super-resolution. They assume that a low-resolution image patch and its closely related high-resolution pendant share the same sparse code $\alpha_{lowres} =$

$\alpha_{highres} = \alpha$ given properly defined reconstruction dictionaries $D_{lowres}$ and $D_{highres}$. They apply feedforward neural networks to compute approximate sparse codes.

Tonolini et al. [151] propose Variational Sparse Coding, which incorporates sparse coding at the inputs of a VAE recognition model to improve feature disentanglement in the latent code. The model is evaluated on the Fashion MNIST, celebA (celebrity faces), and UCI HAR (accelerometer and gyroscope time-series data of human activities) data sets to investigate the disentanglement of features in the latent space and provides promising results and visuals.

## 2.10 Recurrent Neural Networks

Recurrent Neural Networks (RNNs) are a superset of feedforward neural networks which include recurrent edges that incorporate hidden states of previous, and in some cases subsequent time steps. This enables the model to process sequential data of arbitrary length one at a time while maintaining a *memory* of the past. In the context of SDG, RNNs are especially useful for speech synthesis, music generation, or time series prediction. [154]
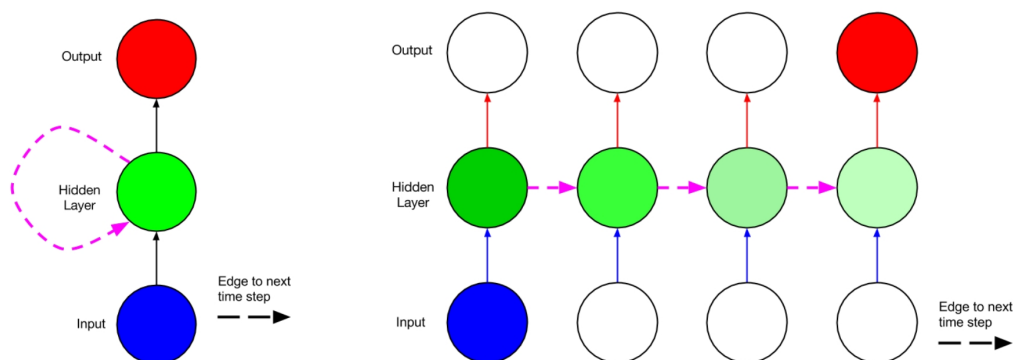


Figure 22: **Left:** A simple RNN network. **Right:** Depiction of the vanishing gradient problem, where through weights less than one, the influence of the first input will diminish over time. (Source: [154])

In a simple RNN (see Figure 22, left) at time step $t$, the current hidden state $\mathbf{h}^{(t)}$ depends on the current input example $\mathbf{x}^{(t)}$ and the previous state $\mathbf{h}^{(t-1)}$, resulting in

$$\mathbf{h}^{(t)} = \sigma(W_{hx}\mathbf{x}^{(\mathbf{t})} + W_{hh}\mathbf{h}^{(t-1)} + \mathbf{b}_h) \tag{24}$$

with the sigmoid activation function $\sigma$ and trainable weight matrices $W_{hx}$, $W_{hh}$ and bias $\mathbf{b}_h$. The predicted output $\hat{\mathbf{y}}^{(t)}$ is then computed from $\mathbf{h}^{(t)}$, so

$$\hat{\mathbf{y}}^{(t)} = \text{softmax}(W_{yh}\mathbf{h}^{(t)} + \mathbf{b}_y) \tag{25}$$

with bias $\mathbf{b}_y$ and trainable weight $W_{yh}$. The weights are usually trained using BPTT, which can suffer from the *vanishing gradient problem* (see Figure 22, right), which LSTMs, a special kind of RNN, and GRUs, a simplified version of LSTMs [155], aim to solve. [154]

Boulanger-Lewandowski et al. [156] generalize RTRBMs and introduce the *RNN-RBM*, which combines an RNN with distinct hidden units with an RBM, whose hidden units are related to the RNN's ones and visible units influence the next hidden state of the RNN. The RBM's ability to generate complex distributions for each time step allows the authors to model and generate polyphonic music in a binary matrix piano-roll representation but not observe long-term musical structure.

Graves et al. [157] propose a deep RNN with $N$ stacked recurrently connected LSTM hidden layers that, at each step, compute the prediction probability for the next word. The network always starts with a null vector as the first input, so all data is generated without prior information. First, they evaluate their model on one-hot-encoded discrete text data and then on online handwriting data, a sequence of pen tip locations, to generate random handwritten character sequences. Next, they combine the handwriting approach with a target character sequence to generate handwriting for a

given text. This works by providing a weighted window on the target text at each RNN time step to the hidden layers. The approach is capable of producing realistic results.

Ranzato et al. [30] introduce *rCNN*, an unsupervised recurrent convolutional neural network, to predict the next frame of a video. The model splits the video frame into $8 \times 8$ pixel patches and feeds a $9 \times 9$ patch of their quantized values into a RNN to create an embedding, which is then processed by two convolutional layers to predict the central patch. The patch-wise convolutional processing allows the rCNN to process videos of arbitrary frame size. The authors evaluate their model on the UCF-101 sport clip data set, resulting in better performance than n-grams and the neural net language model [31].

Vinyals et al. [158] generate captions for images using an end-to-end encoder-decoder architecture, fully trainable with stochastic gradient descent. The Neural Image Caption (NIC) model encodes images using the last hidden layer of a CNN pre-trained for image classification and feeding it to the LSTM decoder, which aims to maximize the likelihood of the sentence being a correct transcription of the image. The model achieves state-of-the-art BLEU, METEOR, and CIDER scores for the time (2015) on the described image data sets Pascal VOC 2008, Flickr8k, Flickr30k, MSCOCO, and SBU.

Donahue et al. [159] propose the Long-term Recurrent Convolutional Network (LRCN) for image interpretation tasks. Images are processed by a single CNN to extract visual features, and a LSTM encoder creates a total representation from these features. A LSTM takes the image representation and the previous word as inputs and generates a description word by word. The authors extract entities and their relations from videos for video description, and the LSTM decodes the entity collection into a meaningful sentence.

Srivastava et al. [160] use LSTMs to learn representations of video sequences and decode them to predict future frames or reconstruct the input sequence. They also propose a composite model performing both tasks simultaneously to overcome their shortcomings. The training is conducted using backpropagation, and the models are evaluated on the UFC-101, HMDB-51, Sports-1M, and moving MNIST digits data sets. The future prediction results are quite blurry, but the representations work well for action recognition when fed into a classifier.

Mansimov et al. [161] present *AlignDRAW*, which extends the DRAW model [97] to generate images given a caption. The caption is defined as a sequence of words and is encoded using a bidirectional RNN, which consists of a forward and backward LSTM whose representations at each time step are concatenated. The generative RNN works similarly to DRAW by including the respective caption representation at each time step and iteratively improving the quality of the existing image. The model also generates plausible results for previously unseen types of captions.

Jaques et al. [162] adopt LSTMs to music generation by training them to predict the next note on a large training corpus. They then use Reinforcement Learning (RL) to optimize their *Note-RNN* by combining a reward function based on rules of musical theory with the output of another fixed copy Note-RNN. They achieve more coherent results that comply with musical theory, avoiding the sometimes occurring randomness of RNNs. They also discuss the application of RL to other domains such as text generation, which could be used to enforce correct grammar.

Van den Oord et al. [115] advance two-dimensional RNNs to sequentially predict pixels (more specifically, their discrete RGB channel values) along the two spatial dimensions with the help of learned probability distributions and a softmax function. They propose two deep (up to 12 layers) LSTM architectures (see Figure 23), also called *PixelRNNs*: The row LSTM, which processes the image row by row from top to bottom, and the diagonal BiLSTM, which uses two LSTMs starting at each top corner and crossing the image diagonally to capture more context. The diagonal BiLSTM outperforms the row LSTM and other models like PixelCNN, DRAW, and DLGMs in image density estimation tasks.
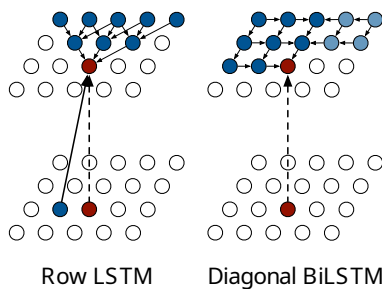


Figure 23: The row LSTM (left) and diagonal BiLSTM (right) PixelRNNs. (Source: [115])

Waite et al. [163] propose *Lookback* and *Attention* RNN, which try to improve the long-term structure modeling capabilities of recurrent networks for music generation. The basic LSTM takes the one-hot encoded vector of the

previous melody event as input. Lookback RNN takes two more previous event vectors, the current position in the music measure (e.g., $\frac{4}{4}$), and whether the last event repeats the event of the previous two events. The LSTM now has to label new vectors as "repeat-1-bar-ago", "repeat-2-bars-ago", or a new melody event. The Attention RNN uses an attention mechanism to look at the weighted sum of the previous $n$ outputs to generate the current step result. The model can be trained on MIDI files to create similar melodies.

Hadjeres et al. [164] introduce the *Anticipation-RNN*, which enables the interactive generation of music by allowing user-defined positional constraints. Because the incorporation of future constraints in a sequential probabilistic model would be computationally expensive, a backward *Constraint-RNN* going from constraint $N$ to 1 is proposed, whose step-wise outputs are combined with the input state of the forward *TokenRNN*, which generates the music tokens from 1 to $N$. The method is general and can be applied to other RNN-based approaches.

Oore et al. [165] train a LSTM on the Piano-e-Competition MIDI data set to generate natural-sounding piano performances in the MIDI event space with varying dynamics and velocity. The LSTM inputs are one-hot encodings over the MIDI event vocabulary (various NOTE-ON, NOTE-OFF, TIME-SHIFT, and VELOCITY events), and the model computes a softmax probability distribution for the output event conditioned on the input events from which it samples. The model can generate music from an initial starting sequence or from scratch (an empty sequence) that comes close to human improvisation. It does not just copy the training samples but can not decide on a coherent play style. The authors prove that powerful generative models for music are possible without defining rules or heuristics at all and without observing long-term relationships in the data.

Short overview of other usages of RNNs:

| Approach | Description | Year |
|---|---|---|
| [166] | RNN-NADE: Combining RNNs with a NADE [146] to output multivariate predictions (probabilities) for music data. | 2012 |
| [167] | Experiments with subword-level language models, resulting in smaller models than character or word-based approaches with similar performance and no out-of-vocabulary predictions. | 2012 |
| [168] | A reference melody and chaotic units are given to a LSTM as input to generate new melodies with a predefined *melodiouness*. | 2013 |
| [169] | Application of fast dropout to RNNs, which drops each incoming unit of a neuron with a certain probability, resulting in a respective zero value in the weighted sum for the neuron activation. Experiments on polyphonic music generation indicate that shallow RNNs perform better with dropout through better generalization. | 2013 |
| [170] | Experiments with multiple different deep RNNs on language modeling and polyphonic music generation tasks. They implement deep state transition and output functions and improve over conventional shallow models, except ones with fast dropout [169] or on character-wise generation, where the subword RNN [167] performs better. | 2013 |
| [171] | Video description by extracting features from each frame with a CNN, applying mean pooling across all frame embeddings and feeding the result to a LSTM which generates the description. | 2014 |
| [172] | Chinese poem generation from user-supplied keywords. Chinese symbols are one-hot encoded. A convolutional sentence model creates line embeddings fed to a recurrent context model that forwards a context vector to the recurrent generation model that creates word after word for the current line. | 2014 |
| [173] | Stochastic Recurrent Networks (STORNs): Two RNNs, the recognition (encoder) model $q(z_t \vert \mathbf{x}_{1:t-1})$ from which $\mathbf{z}_t$ is sampled, and the generating model $p(x_t \vert \mathbf{z}_{1:t})$, from which $\mathbf{x}_t$ is obtained, form a network easily trainable with stochastic gradient descent that can model multiple variables at each time step. Additionally, the latent representations $\mathbf{z}_t$ are conditioned on a prior $p(\mathbf{z})$ like a VAE. The model is evaluated on polyphonic music generation and motion capture data continuation, where it outperforms previous RNN approaches except RNN-NADE [166]. | 2014 |
| [174] | Clockwork RNN: Partition of hidden layers into separate modules $i$ with arbitrary periods $T_i$ that are active at time step $t$ only if $t \mod T_i = 0$. This reduces the number of parameters and increases the performance on sequence reconstruction tasks. | 2014 |
| [175] | LSTM training with resilient propagation [176] instead of BPTT for improved music composition represented as "binary" piano key presses. | 2014 |
| [177] | RNN-DBN: Very similar to the RNN-RBM [156], this model combines RNNs and DBNs, which are RBMs with multiple stacked hidden layers. It is also used for polyphonic music generation and improves upon the RNN-RBM results. | 2014 |

<div align="right">Continuation...</div>

... Continuation

| | | |
|---|---|---|
| [178] | Multimodal RNN (m-RNN): A RNN that models the next word probability distribution based on the previous word and the encoding of an image provided by a CNN. | 2014 |
| [179] | Recurrent image density estimator (RIDE): Combining spatial LSTMs [180] that compute the hidden state $\mathbf{h}_{ij}$ of the next pixel $x_{ij}$ based on the two axis-wise preceding states $c_{i,j-1}$ and $c_{i-1,j}$, and mixtures of conditional Gaussian scale mixture [181] that predict the state of the next pixel $p(x_{ij}|\mathbf{h}_{ij})$. | 2015 |
| [182] | Bidirectional RNNs as gap fillers in high-dimensional categorical and binary time series data (e.g., music), outperforming unidirectional RNNs and being applicable in more scenarios. | 2015 |
| [183] | Variational RNN (VRNN): A recurrent VAE for high-dimensional sequence generation using the hidden state $h_t$ of a RNN as the parameter for the distribution of the latent random variable $z$ of the VAE at each step. The model outperforms simpler configurations in unconditional natural speech and handwriting generation. | 2015 |
| [184] | Scheduled sampling: To bridge the gap between training, where usually the ground-truth previous token $x_{t-1}$ is taken for the next prediction, and inference (generation of new sequences) distributions, where the model uses its previous prediction $\hat{x}_{t-1}$, it is randomly decided during training for each token prediction, whether $x_{t-1}$ or $\hat{x}_{t-1}$ is taken. Improved results on image captioning (CNNencoder) and speech recognition compared to a baseline LSTM without scheduled sampling. | 2015 |
| [185] | Mind's Eye: Learning bi-directional mappings between visual features of images (obtained from a VGG [186]) and their text descriptions with RNNs. The model can be used to generate in both directions by first training a RNN to generate the text from the features and then a second RNN on top to reconstruct the visual features from the text. | 2015 |
| [187] | Alignment of visual (object detector CNN) and language (bidirectional LSTM) representations of image regions for full-frame and region-level image description with RNN decoding. | 2015 |
| [188] | DBN-LSTM: Improved version of the RNN-DBN [177], where the RNN is replaced with a LSTM for better performance in polyphonic music generation. | 2015 |
| [189] | LSTM-RTRBM: Replacement of some hidden units of a RTRBM with LSTM ones, which increases performance and learning speed. | 2015 |
| [190] | Image captioning with convolutional feature extraction and an attention-based LSTM that generates the caption word-by-word. | 2015 |
| [191] | Video description with CNN-encoded video frames as input for a two-layer LSTM that starts outputting words after the whole video sequence has been processed by the first LSTM layer. | 2015 |
| [192] | Video description with a 3D CNN encoder (width × height × timesteps of video) and LSTM decoder with temporal attention mechanism. | 2015 |
| [193] | A stochastic recurrent neural network (SRNN) with separate stochastic and deterministic layers propagates uncertainty in latent space through the network. The hidden representation $\mathbf{z}_t$ depends on a prior $p(\mathbf{z_z}|\mathbf{z_{t-1}})$, similar to a VAE, parameterized by a neural network. Achieves state-of-the-art performance on speech and polyphonic music modeling. | 2016 |
| [194] | Single-note melody generation and continuation with a multi-layer GRU trained on sequences of corresponding pitch and duration one-hot vectors. | 2016 |
| [195] | Improved LSTM training for music composition. Pitches and durations of notes are encoded together in one one-hot vector. Training is split into two parts, where first, the model is trained with real data, and new compositions are generated. These creations are filtered by the *grammar argumented method*, which only allows samples complying with defined musical rules to remain. These are then appended to the training data, and the model is retrained, resulting in the actual generative model. | 2016 |
| [196] | Training of a sequence prediction RNN *actor* with a second *critic* network that has access to the ground-truth data and computes the expected task-specific score. The RL-inspired actor-critic training approach fits training data faster than maximum-likelihood learning. It provides more coherent and accurate results on text prediction tasks (i.e., spelling correction and machine translation). | 2016 |
| [197] | Review Network: An encoder-decoder framework with reviewers in between that provide a discriminative loss. Additionally, an attention mechanism between the encoder and the review networks and the review networks and the decoder is implemented. The encoder can be a CNN or RNN, while the decoder is a LSTM. The model is used to caption images and generate comments for Java source code. | 2016 |

Continuation...

... Continuation

| [198] | Image captioning based on CNN encoder, visual attribute prediction (i.e., what happens in the picture), a LSTM for state progression, and two different attention mechanisms for input and output that determine the next word from the vocabulary. | 2016 |
|---|---|---|
| [155] | Video description in one or multiple sentences using a hierarchical RNNs with a sentence and paragraph generator. The network is based on GRUs and uses temporal and spatial attention mechanisms. Input features of the video are obtained from a pre-trained convolutional extractor [186]. | 2016 |
| [199] | Two text-based LSTMs (word-RNN and char-RNN) learn chord progressions from text representations of music for fully automatic music composition. | 2016 |
| [200] | Multi-track pop music generation with a hierarchical LSTM where each layer is responsible for predicting a certain aspect of the song at a time step: The bottom layer predicts the pressed key, the next layer the duration of the press, the third the chord and the fourth and last layer the drum beat. | 2016 |
| [201] | Folk-RNN: Large-scale generation of Celtic folk music transcriptions with LSTMs trained with a vocabulary of tokens on single transcriptions. Evaluation is performed on the large-scale distributions of real and generated data and the single transcription level. The LSTMs automatically learn to conform to structural constraints of folk music. | 2016 |
| [202] | BachBot: A three-layer stacked LSTM with optimized parameters generates Bach chorales or transfers their style to other melodies. It is trained with *teacher forcing* (always continue prediction with the correct previous token) on frame-based representations of the polyphonic piano roll data. | 2016 |
| [203] | Two-layer LSTM for token-level music generation from a short seed sequence. Each MIDI message or existing note combination from a piano roll representation is treated as a separate token. The results are comparable to RNN-NADE [166]. | 2016 |
| [204] | SampleRNN: An unsupervised and unconditional end-to-end model for raw audio waveform synthesis with hierarchical RNNs that cover different temporal ranges. | 2016 |
| [205] | Application of a (bidirectional) hierarchical recurrent encoder-decoder (HRED) to dialogue generation by utilizing question-answer pairs and pre-trained word embeddings. The model consists of a RNN encoder, whose final representation is fed into a context RNN, which maps the representation into the dialogue context. The context state is then fed to each step of the decoder RNN. | 2016 |
| [206] | The "Professor Forcing" algorithm is designed to align the behaviors of recurrent neural networks during training and sampling phases, addressing a common issue in RNN training. This method applies adversarial domain adaptation, where the network is trained to make its behavior indistinguishable between these two phases under the scrutiny of a discriminator. This approach helps in producing more coherent and structured outputs. The algorithm serves as a regularization technique, enhancing the network's ability to generalize from its training data, leading to improvements in performance on various tasks including language modeling and image synthesis. | 2016 |
| [207] | Char2Wav: End-to-end speech synthesis from text. A bidirectional RNN encodes text, and a RNN with attention decodes the representation at different time steps to produce features for a vocoder. The vocoder is a conditional SampleRNN [204] that produces raw waveform output from the vocoder features. | 2017 |
| [208] | Recurrent Highway Network: LSTMs with multiple highway layers [209] that allow deep step-to-step transition functions that are easily trainable. The model outperforms previous RNNs in character prediction tasks on Wikipedia texts. | 2017 |
| [210] | Deep Artificial Composer (DAC): Extension of [194] that uses LSTMs (duration and pitch RNN) to generate note transitions and is trained on a corpus with two different musical styles. It also emphasizes a new novelty measure (fraction of transitions found in a defined corpus) that is used to improve the creativity of the model. | 2017 |
| [211] | Improvement of image captioning with attention [190] by supervised training of the attention mechanism with ground truth text entity-image region mappings obtained from humans. | 2017 |
| [212] | Semantic Composition Network (SCN): Image captioning in multiple parts: A CNN extracts a feature vector, and a MLP computes probabilities of tags based on the most used words in the training descriptions. A LSTM then generates the description based on the feature vector and the tag probabilities. | 2017 |

... Continuation

| [213] | DeepBach: Generation or reharmonization of Bach chorales/MIDI data using deep LSTMs working in opposite directions and neural networks that merge the RNN results. Sampling from this dependency network is performed using pseudo-Gibbs sampling. | 2017 |
|---|---|---|
| [214] | Conditional drum rhythm generation with a combination of a two-layer stacked LSTM that learns drum sequences and a feedforward fully connected layer that processes metrical rhythm and a bass sequence as constraints. Their predictions are merged to predict the next drum event. | 2017 |
| [215] | Sequence Tutor: An improved fine-tuning approach for RNNs that first trains an RNN on data with maximum-likelihood estimation and uses its output as a policy (originally proposed in [162]) for a second RNN trained with RL for a specific domain. The effectiveness is demonstrated on music melody and molecule generation (i.e., SMILES [112] strings). | 2017 |
| [216] | Chord progression generation from monophonic melodies with BiLSTMs. | 2017 |
| [217] | Latent variable hierarchical recurrent encoder-decoder (VHRED): A stack of encoder, context, and decoder RNNs combined with a stochastic latent variable conditioned on all previously observed tokens that can capture the dependencies of sub-sequences of sequential data. The stochastic variable allows for diverse and coherent dialogues to be generated. | 2017 |
| [218] | Harmonic Improviser: A LSTM harmony agent trained on Jazz chord progressions and a rule-based melody agent manipulating provided melodies take turns improvising music in real-time and are rewarded for harmonic consistency and melodic flow. | 2017 |
| [219] | Performance RNN: MIDI music generation with a LSTM showing timing and dynamics, but lacking long-term coherence. | 2017 |
| [220] | TP-LSTM-NADE & BALSTM (biaxial LSTMs): Modification of an RNN-NADE [166] to model relative differences between nodes for transposition-invariant polyphonic music generation. | 2017 |
| [221] | Combination of a biaxial LSTM [220] for symbolic music generation and a conditional WaveNet-based audio generator for waveform music generation. | 2018 |
| [222] | Graph Neural Networks (GNNs): Computing and iteratively updating node and graph embeddings using fully connected neural networks, GRUs for nodes and LSTMs for edge modifications. MLPs are used to sequentially compute probabilities of adding new nodes (and their type) and edges based on these representations. The model can work conditionally and unconditionally and is demonstrated on molecule generation. | 2018 |
| [223] | WaveRNN: A sparse single-layer RNN mainly for text-to-speech synthesis where the majority of weights are pruned and subscaling is employed to fold long sequences into a batch (matrix) of short ones, which allows generating multiple samples (i.e., over multiple rows) at one step. This allows the model to produce results in real time on a mobile CPU. | 2018 |
| [224] | Relational RNN: A LSTM with a multi-slot memory matrix instead of the hidden state vector. Input is concatenated to the matrix as a new row, and multi-head dot product attention [18] is applied to generate the next memory state. The model achieves state-of-the-art results on language modeling tasks (next-word probability). | 2018 |
| [225] | DeepJ: A model built upon the biaxial LSTM [220] structure and incorporate dynamics (i.e., relative note volume) in note embeddings and global style and context (e.g., genre) conditioning in the network. They train three outputs simultaneously (play and replay probability and dynamics) for the predicted notes. | 2018 |
| [226] | GraphRNN: Generate large variable-length graphs without node ordering by treating the problem as a sequence of node and edge additions. At each step $i$, the graph-level GRU updates the graph state $h_i$ and adds a new node. The edge-level RNN then creates the adjacency vector for the new node to all old nodes or the end-of-sequence token from $h_i$. The model is trained on data obtained using breadth-first-search through any graph permutation with a random starting point. The graphs are evaluated by computing a Maximum Mean Discrepancy (MMD) score between the degree and clustering coefficient distributions and orbit count statistics between sets of graphs. | 2018 |
| [227] | Tacotron 2: A recurrent sequence-to-sequence text-to-speech network consisting of multiple LSTMs and CNNs that maps character embeddings to simplified *mel* spectrograms that are then converted to waveform audio with WaveNet [228]. | 2018 |

Continuation...

... Continuation

| [229] | Graph Recurrent Attention Network (GRAN): Improvement upon GraphRNN [226] using a GNN [222] with attention at each step to generate a block of new nodes and edges based on the already existing graph. Further, they propose training on adjacency matrices conforming to families of node orderings (e.g., nodes sorted by node degree, breadth/depth-first-search ordering from largest degree node, original data order) to improve model understanding. | 2019 |
|---|---|---|
| [230] | Generating undirected, fully connected graphs without self-loops as an ordered edge sequence using GRUs. The first GRU predicts the sequence of the first nodes of the edge pairs. The second GRU outputs the probabilities for the second node of the edge. The graph nodes are assigned a fixed order in advance and the RNNs are trained to maximize the node probabilities observed in the training data. This simple model performs similarly to the GraphRNN [226]. | 2019 |
| [231] | MolecularRNN: Extension of GraphRNN [226] that uses a *NodeRNN* to compute the next atom type and an *EdgeRNN* to compute the bond types to the previous atoms. The model enforces valid per-atom valency. It is first trained to reconstruct the training data and then fine-tuned using a RL critic that rewards molecules with certain properties. | 2019 |
| [232] | Deep Graph Distribution Learning (DeepGDL): Decomposing a graph into densely connected components with sparse connections between these communities. A GRU is then trained to learn the distribution of nodes and edges in these communities based on earlier node and edge observations. Synthetic communities are sampled from these predicted distributions and probabilistically connected to produce new large graphs (synthetic power grids) with similar properties to the real data. | 2019 |
| [233] | Generation of biomedical signals (electrocardiogram, etc.) for patients or specific events using bidirectional RNNs that are trained with real patient data. In the first optional stage, noise is injected into the real signals, and then the signal is segmented according to certain events or classes. Finally, the BiRNN generates new similar data based on the input, and a statistical stage evaluates the data quality. | 2020 |
| [234] | GraphGen: A domain-agnostic and scalable labeled graph generation method using a LSTM to sequentially append tuples containing the sequence and types of nodes and edges in a depth-first search order. Other models like GraphRNN [226] are outperformed in almost every evaluated criterion. | 2020 |
| [235] | Scramble: Music generation with pitch transitions generated by a Markov chain and a LSTM that imposes a learned style (velocity, rhythm, and beats per minute) on the pitch sequence, which the user can also tweak. | 2022 |

## 2.11 Convolutional Neural Networks

Convolutional Neural Networks (CNNs) are artificial neural networks based on matrix operations that can handle data of large sizes like images or speech with significantly fewer parameters to train than a fully connected neural network. The model can contain different types of layers [236]:

**Convolution** A *kernel* or *filter*, which is a lower-size matrix consisting of trainable weights, is applied to parts of the input matrix to extract local features, for example, to detect edges in an image. The kernel is moved over the input left-to-right, and top-to-bottom by a *stride*, which defines the number of units shifted at each step, to compute the output value at this kernel position through matrix multiplication of the weights with the respective values in the input at the kernel's position (see Figure 24a).

**Padding** To prevent loss of information at the border of the input or to preserve the input size, a zero-padding can be added around the matrix.

**Non-linearity** After a convolution layer, a non-linear function is used to modify or cut off the output. Typically, the rectified linear unit (ReLU) function $ReLU(x) = max(0, x)$ is used.

**Pooling** A pooling layer downsizes its input to reduce complexity for later model layers. Popular implementations are max-pooling or average-pooling, which return a region's maximum value or average, respectively. They are applied in the same manner as a kernel.

**Fully connected layer** Finally, a computationally expensive fully connected layer is applied to the significantly smaller input to, for example, classify an image or perform another task. An illustration of such a CNN architecture is provided in Figure 24b.

Lotter et al. [237] predict future frames of image sequences by first learning a representation of each single input image with a CNN. Then, a LSTM processes the images sequentially in order before the final output is forwarded to a

(a) Application of kernels or pooling functions to an input matrix with stride 1. (Source: [236])



(b) Example CNN architecture for an image classification task. (Source: [17])
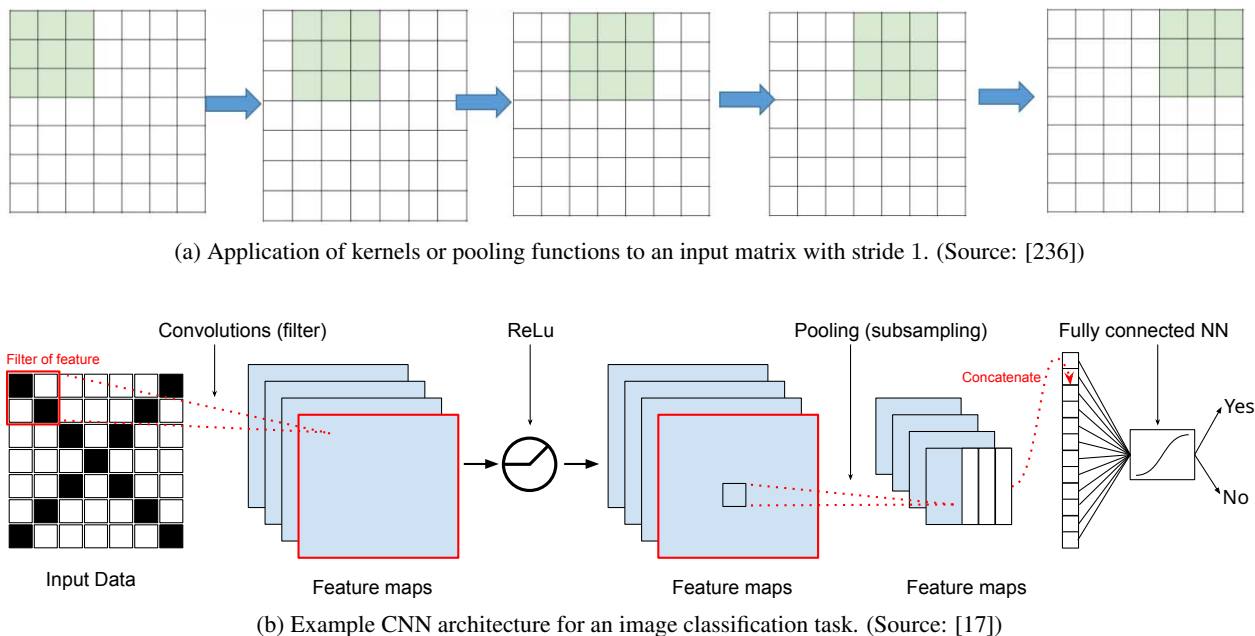
Figure 24: Illustrations of the structure of a CNN.

deconvolutional (i.e., reversed) CNN, which produces an image from the RNN prediction. They train the model using mean squared error but also experiment with adversarial loss [238] realized by a similar CNN-LSTM discriminator whose final prediction together with an encoding of the generator's result image or the real next frame is passed to a MLP. The model can predict movements and rotations in videos well, which supports the idea that prediction may enable the development of transformation-tolerant object representations.

Bruna et al. [239] propose CNNs for high-dimensional structured prediction problems such as image super-resolution. They model the conditional distribution $p(y|x)$ as a Gibbs density $p(y|x) \propto \exp(-\|\Phi(x) - \Psi(y)\|^2)$, where $\Phi : \mathbb{R}^N \to \mathbb{R}^P$ and $\Psi : \mathbb{R}^M \to \mathbb{R}^P$ are highly-informative non-linear mappings (sufficient statistics) obtained from deep CNNs that "minimize the uncertainty of $y$ given $x$". The model can provide solutions with spatial coherence. Still, the inference is computationally costly compared to GANs, and a trade-off between sharpness and stability must always be made.

Van den Oord et al. [115] propose the *PixelCNN*, which consists of multiple resolution-preserving convolutional layers with masks (see Figure 25) to ignore future pixels. Like the PixelRNNs, a softmax function is used to compute the discrete RGB values of pixels sequentially for image generation and completion tasks. The advantage of the PixelCNN over the PixelRNN is the possibility of parallelization during training and evaluation of test images, but the performance is worse. The model has been successfully adapted to video continuation as a decoder [240].

Van den Oord et al. [228] introduce *WaveNet*, which operates on the raw audio waveform and is similar to PixelCNN [115] in that it models the conditional probability distribution $p(x_t|x_1, ..., x_{t-1})$ with a stack of convolutional layers. It uses dilated causal convolutions (see Figure 26) to preserve ordering and process an area larger than the convolution length by skipping values by a step. The model can be trained in parallel but generates new audio sequentially. The model can also easily be transformed to incorporate additional input (e.g., text and speaker identity for text-to-speech) by appending it to the conditional distribution, resulting in a conditional WaveNet. WaveNet achieves state-of-the-art results in text-to-speech tasks and allows conditioning on different speakers. Further, the model can generate novel and realistic musical fragments.

Gatys et al. [242] propose to use a deep CNN [186] to separate style (texture) and content (object recognition) representations of an image for style transfer. The style transfer works by extracting multiple layers of style and content representations from style and content reference images, respectively. Then, a white noise image $\vec{x}$ is initialized and iteratively optimized using gradient descent with respect to the pixel values based on the combined style and content losses, which are the sum of squared errors between the respective representations (style image and $\vec{x}$, $\vec{x}$ and content image).

Kim et al. [243] use a recursive CNN for image super-resolution. First, an embedding network, similar to a MLP, encodes an image as a set of feature maps. Then the recursive CNN applies the same convolution followed by a ReLU
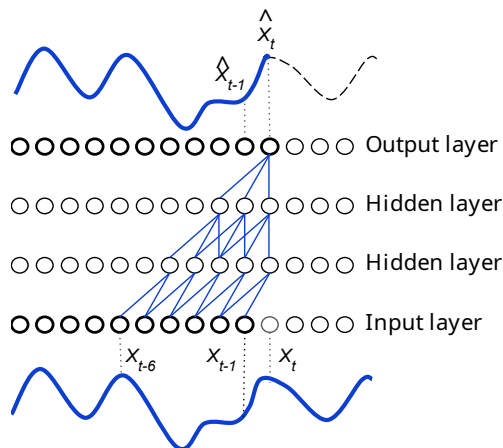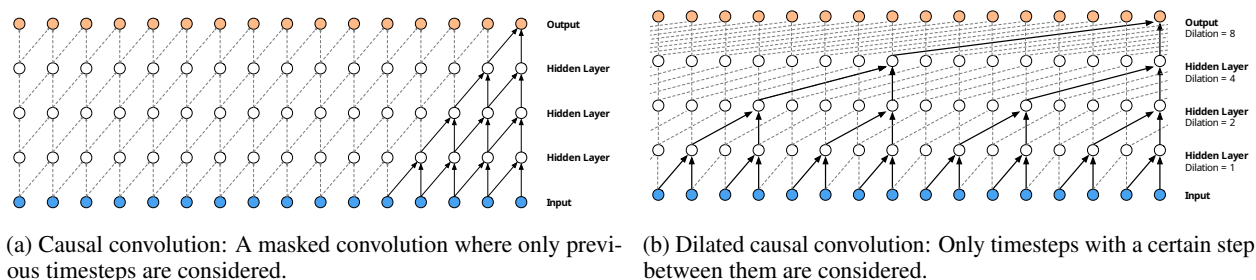
Figure 25: Auto-regressive 1D signal modeling with a masked CNN. (Source: [241])



(a) Causal convolution: A masked convolution where only previous timesteps are considered.

(b) Dilated causal convolution: Only timesteps with a certain step between them are considered.

Figure 26: Illustrations of a causal and dilated convolution. (Source: [228])

to the embedding to increase the observed range before a reconstruction network, also similar to a MLP, creates the final output. To solve the problems of vanishing/exploding gradients and finding the optimal amount of recursions during training, each recursive layer is reconstructed, and a weighted average of all predictions produces the final output (*recursive supervision*). A skip connection between the input image and the reconstruction network is established to improve results further. The method outperforms SRCNN [244, 245] and provides clearer images.

Salimans et al. [246] accelerate PixelCNN training and generate state-of-the-art results on class-conditional and unconditional image generation task CIFAR-10 with their improved implementation called *PixelCNN++*. They compute the RGB values of pixels assuming linear dependence and using continuous distributions that are rounded instead of 256-way softmax to reduce training cost. Further, they introduce downsampling and shortcut connections between layers to better capture the input structure and use dropout regularization to prevent overfitting.

Short overview of other usages of CNNs:

| Approach | Description | Year |
|---|---|---|
| [244] | SRCNN: A deep CNN learns an end-to-end mapping from low to high-resolution images for super-resolution. The first layer extracts feature maps from overlapping patches, the second maps these feature maps to higher-resolution maps, and the third combines the predictions in a neighboring area for the final result. | 2014 |
| [247, 248] | DeepDream: An image classification CNN that is reversed to (iteratively) amplify high-level features in random noise or real images, often resulting in dream-like over-interpretations. This approach allows a user to visually inspect what a model has learned about objects or concepts and can be used to create art-like images. | 2015 |
| [245] | Improvement of SRCNN [244] for simultaneous 3-channel color handling. Further, different model architectures (larger filters and more layers) and parameters are explored. | 2015 |

Continuation. . .

... Continuation

| [249] | Gated PixelCNN: Improved image quality over the original PixelCNN by utilizing gated convolutional layers, matching the PixelRNN. The model is suitable for class-conditional image generation and as a powerful decoder for an autoencoder. | 2016 |
|---|---|---|
| [241] | Proposal of the gated conditional PixelCNN with text (GRU encoder), segmentation map, and keypoint (i.e., annotated locations of certain human/bird body parts in the image) conditioning for image generation. | 2016 |
| [250] | Image super-resolution and style transfer like [242], but three orders of magnitude faster with qualitatively similar results by using a perceptual loss obtained from a pre-trained VGG network [186] instead of a per-pixel loss. | 2016 |
| [251] | ESPCN: Image and video super-resolution with sub-pixel CNN learning upscaling filters from low-resolution feature maps into high-resolution output. | 2016 |
| [252] | Deep Voice 3: A fully-convolutional encoder-decoder model with position-augmented attention mechanism for text-to-speech synthesis. The model can produce different parameters for various waveform synthesis models (e.g., WaveNet [228]) and incorporate a speaker representation to capture different speech styles. The model achieves state-of-the-art quality in human mean opinion score evaluations. | 2017 |
| [253] | PixelSNAIL: Application of SNAIL [254], a general purpose autoregressive meta-learning model using causal convolutions and self-attention to maximize its context size, to sequential image generation, resulting in state-of-the-art likelihood, but slow density estimation performance. | 2017 |
| [255] | Subscale Pixel Network (SPN): Images of size $N \times N$ are split into slices of size $\frac{N}{S} \times \frac{N}{S}$ that are interleaved. The network consists of a convolutional encoder that embeds previously processed slices and a convolutional decoder with masked convolution and self-attention that predicts the next slice given the embedding. The model is especially suitable for upscaling and generates coherent and exact samples. | 2018 |
| [256] | Joint training of an ensemble of shallow and deep CNNs to generate super-resolution images end-to-end. Optimization during training is alleviated by letting the shallow CNN restore the main structure of the image and the deep CNN fill in the details. The deep CNN extracts features, upscales them to the target factor, and uses multi-scale reconstruction to capture the context better and produce the output pixels. | 2019 |

## 2.12 Transformers

Transformers (see Figure 27c) are sequence-to-sequence transduction models with an encoder-decoder structure. They advance previous recurrent and convolutional encoder-decoder architectures because they allow for more parallelization and modeling of dependencies with arbitrary distance in the input and output sequences with a constant number of sequential operations. For that, the transformer utilizes a multi-headed *self-attention* mechanism (see Figure 27a and Figure 27b) instead of recurrent layers. Originally, the transformer was used for language translation tasks, where the original sentence was first encoded, and the decoder used attention to the encoder embeddings and the already generated output in the target language to compute the probabilities of the next token. [18]

Liu et al. [257] generate English Wikipedia articles by providing the target article title and summarizing multiple non-Wikipedia reference documents. The training data consists of Wikipedia articles, their citations, and the top 10 web search results for each article section title. The paragraphs of all reference documents are ranked by importance according to the target article title. Then, the first $L$ tokens of these ranked paragraphs are used as the input for the generative model. The generative model is a transformer without an encoder that concatenates the input sequence and the desired output sequence into a single vector for the decoder input and is trained to predict the next token based on the previous ones. The authors also modify the attention mechanism by splitting the tokens into blocks on which attention is independently applied and performing convolution on the key-value pairs to reduce memory requirements on long input sequences. The results show that the model can split articles into reasonable sections and fill in factual information from many different references.

Parmar et al. [258] propose the *Image Transformer* conditioned on a few class embeddings (decoder only) or low-resolution pictures (encoder-decoder architecture) to generate high-resolution images. The generation process is formulated as a sequence modeling problem, where the RGB channel values of the next pixel are predicted based on the other pixels' values in the local neighborhood to allow for larger image sizes. The authors propose two different attention mechanisms, 1D and 2D local attention (see Figure 28). Experiments on CIFAR-10, ImageNet, and celebA

(a) The scaled dot-product attention maps weights to values of key-value pairs $(K, V)$ according to the correspondence between a query $Q$ and the key.

(b) Multi-head attention uses multiple attention layers and concatenates the results.

(c) Encoder-decoder architecture of the transformer. Both parts consist of $N = 6$ layers.
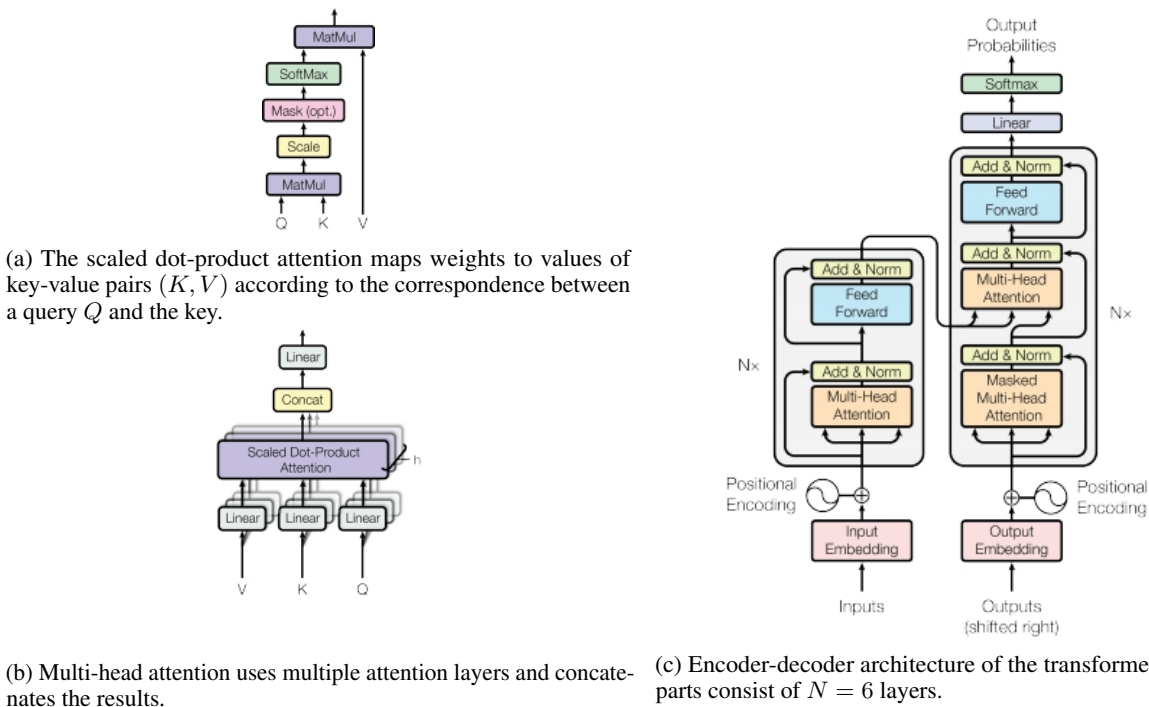
Figure 27: Structure of the transformer itself and essential parts. (Source: [18])

data sets show that the transformer architecture outperforms previous state-of-the-art architectures like RNNs, CNNs, and GANs in terms of processible image size and quality.
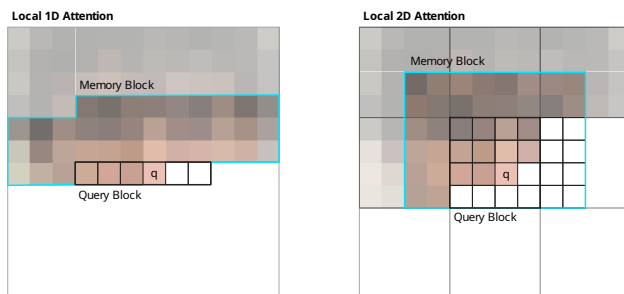


Figure 28: Illustration of the 1D and 2D local attention mechanisms of the *Image Transformer*. The 2D attention performs slightly better in terms of perceptual image quality evaluated by humans. (Source: [258])

Huang et al. [259] use a transformer to generate symbolic music represented as a sequence of discrete tokens. Since pieces often repeat and modify previous motifs or sections, relations between such sections are explicitly modeled. Therefore, the authors adopt a relation-aware self-attention [260] that creates relative position embeddings and optimizes the memory consumption to allow long sequences to be autoregressively modeled. They achieve the best NLL scores and most win in a human "musicality" comparison test when compared against the PerformanceRNN and LookBackRNN [163] LSTM models and a baseline transformer when continuing a music sequence they were initialized on.

Child et al. [261] introduce the *sparse transformer*, which uses sparse factorizations of the attention matrix to reduce the time and memory requirements from $O(n^2)$ to $O(n\sqrt{n})$ for the sequence length $n$ without performance loss. This works by splitting the attention operation into multiple faster operations that only access a subset of all previous positions and combining their results to approximate the full attention. The model can generate large unconditional sequence samples in various domains such as natural images (CIFAR-10, ImageNet64) and raw audio data of classical music and achieves state-of-the-art results in density modeling tasks compared to contemporary models.

Sun et al. [262] build *VideoBERT*, a joint visual-linguistic model in the style of BERT [263], which uses masked language model and next sentence prediction tasks to train the language understanding of a transformer. VideoBERT

pairs representations of videos, consisting of spatiotemporal features extracted with pre-trained video classification models and their text description, obtained through an automatic speech recognition system, and trains the transformer on filling in masked tokens in both data types (representations of frames, not raw image data) or deciding whether the text matches the video features. The trained model is then used for action classification, video captioning, future video token forecasting, and prediction of video tokens for text descriptions (see Figure 29), achieving coherent results and state-of-the-art captions.
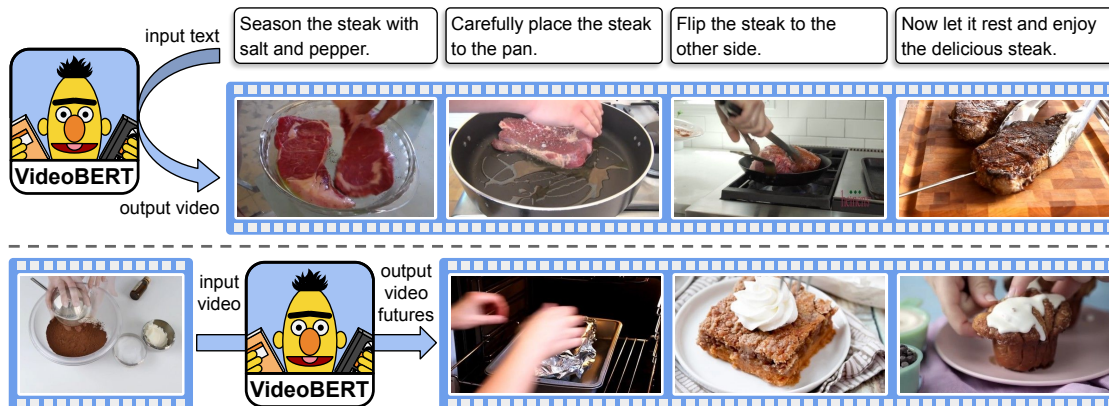


Figure 29: Text-to-video token generation and future token prediction with VideoBERT. The images depicted from the training data have the most similar token representation to the prediction. (Source: [262])

Liu et al. [264] propose a *graph transformer* that replaces the edge-output network of the GraphRNN [226] with a transformer decoder with self-attention layers and attention layers referring to the hidden graph state of the node RNN. The results are competitive or better than GraphRNN on a variety of metrics.

## 2.13 Generative Adversarial Networks

Generative Adversarial Nets (GANs) are frameworks consisting of a generator $G$ creating synthetic data from random noise and a discriminator $D$ determining whether a provided sample came from $G$ or the training data. The authors describe their system as a "minimax two-player game" [238], where the generator tries to deceive the discriminator. [265]
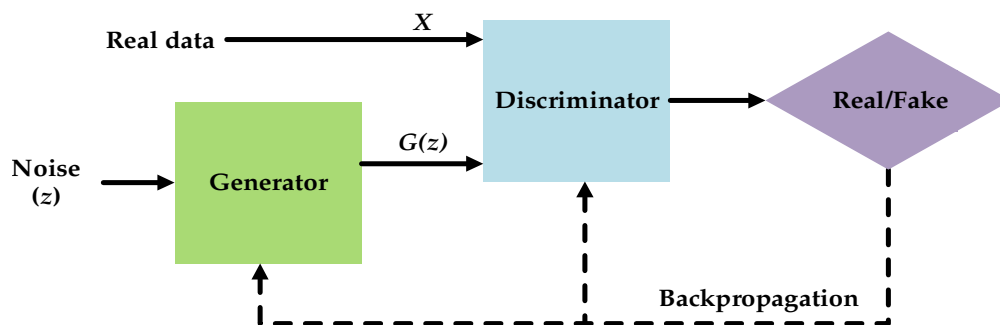


Figure 30: The original GAN implementation. (Source: [265])

As seen in Figure 30, both $D$ and $G$, originally implemented as MLPs to provide a network model with non-linear mapping, are learning from the results of $D$. $G$ learns the probability distribution of the real data $p_{data}(x)$ and $D$ the distribution of the random noise $p_z(z)$. The goal of the optimization process in Algorithm 2 is reaching the *Nash equilibrium* between $D$ and $G$, so both distributions become indistinguishable and the probability of $D$ classifying a sample as either fake or real approaches 50%. [265]

In practice, minimizing the cost of the discriminator and generator jointly is difficult and often leads to instability because minimizing one cost function often means increasing the other. A GAN may fail to converge. Further, GANs often lack diversity due to the *mode collapse problem* induced by the effort of the generator to deceive the discriminator, not represent a realistic data distribution. This often leads to only certain "easy" data types being generated and repeated. [4]

---

**Algorithm 2** Training algorithm for GANs. $\theta_d$ and $\theta_g$ are the parameters of the respective MLPs $D$ and $G$ that are updated. In the original experiments, $k = 1$ is used, but a higher $k$ is recommended to keep $D$ close to the optimal solution and prevent $G$ from overfitting. (Source: [238])

---

**for** number of training iterations **do**
    **for** $k$ steps **do**
        Sample minibatch of $m$ noise samples $\{z^{(1)}, ..., z^{(m)}\}$ from $p_z(z)$.
        Sample minibatch of $m$ examples $\{x^{(1)}, ..., x^{(m)}\}$ from $p_{data}(x)$.
        Update the discriminator by ascending its stochastic gradient:
$$\nabla_{\theta_d} \frac{1}{m} \sum_{i=1}^{m} [\log D(x^{(i)}) + \log(1 - D(G(z^{(i)})))].$$
    **end for**
    Sample minibatch of $m$ noise samples $\{z^{(1)}, ..., z^{(m)}\}$ from $p_z(z)$.
    Update the generator by descending its stochastic gradient:
$$\nabla_{\theta_g} \frac{1}{m} \sum_{i=1}^{m} \log(1 - D(G(z^{(i)}))).$$
**end for**
The gradient-based updates can use any standard gradient-based learning rule. We used momentum in our experiments.

---

The authors of the original GAN framework, Goodfellow et al. [238], train a setup of two MLPs on three datasets: MNIST, TFD and CIFAR-10. They achieve first and second-place results on the MNIST and TFD datasets, respectively, when compared against a DBN, a deep GSN, and a stacked CAE in terms of log-likelihood estimates [266].

Radford et al. [267] develop the Deep Convolutional GAN (DCGAN) architecture, which aims to adopt CNNs to unsupervised learning tasks by using the GAN framework. DCGANs replace the pooling functions of CNNs with stridden convolutions for the discriminator and fractional-stridden convolutions for the generator and remove fully connected layers entirely. Further, the neural networks now utilize ReLU activation functions and batch normalization for all parts. An example of a DCGAN architecture for image modeling is depicted in Figure 31. Frid-Adar et al. [268] apply the DCGAN to generate labeled liver lesion images.
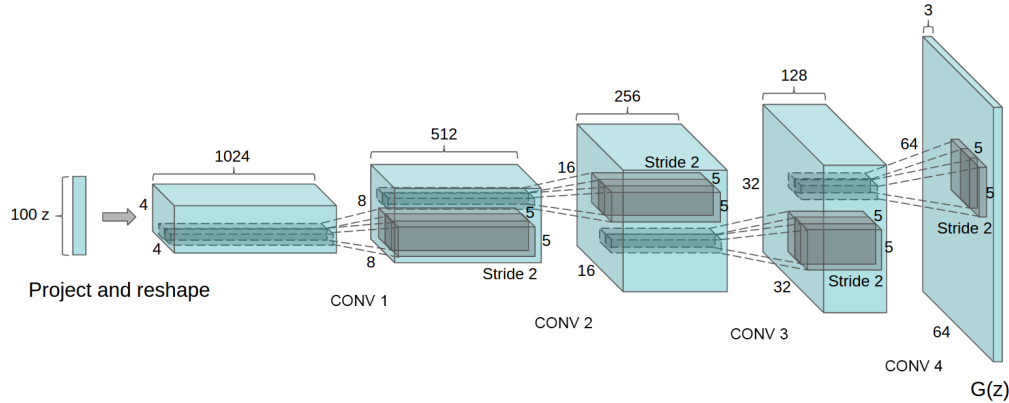


Figure 31: The DCGAN generator architecture used to generate bedroom scene images trained on the LSUN data set. (Source: [267])

Metz et al. [269] unroll the parameters of GAN discriminators $\theta_D^0 = \theta_D$ for $K$ future steps

$$\theta_D^{k+1} = \theta_D^k + \eta^k \frac{df(\theta_G, \theta_D^k)}{d\theta_D^k} \tag{26}$$

and update the generator parameters

$$\theta_G \leftarrow \theta_G - \eta \frac{df(\theta_G, \theta_D^K)}{d\theta_G} \tag{27}$$

with learning rate $\eta$ and objective function $f$, accordingly to stabilize the training of the generator at the expense of increased computational cost during training. By letting the generator "see into the future", the next discriminator step becomes less effective, which balances the two models better and improves convergence.

Elgammal et al. [270] train a DCGAN [267] to produce art images. They achieve this by changing the training objectives of both the generator and discriminator to encourage learning about styles and art separately. The Creative

Adversarial Network (CAN) generator aims to deviate as much as possible from learned styles while keeping the learned art aspects. The discriminator, on the other hand, decides whether an input image is art or not and classifies the art style. During training, the discriminator is trained with art images with style labels to refine the decisions and style classifications. The generator receives the art/not art decision and the style ambiguity of the discriminator as a loss. An evaluation with human subjects indicates that generated art is indistinguishable from real art.

Donahue et al. [271] introduce *WaveGAN*, which synthesizes one-second clips of raw-waveform audio unsupervised. WaveGAN can generate words, bird chirping, and instruments by capturing periodic patterns in the sampled waveforms with convolutions. The approach is based on DCGAN [267] but has a modified one-dimensional convolution kernel with a higher stride. Further, *phase shuffle* is used to shift the activations of each layer's activations by a random integer $\in [-n, n]$ to prevent the discriminator from learning trivial policies. They compare WaveGAN to another of their models based on DCGAN, *SpecGAN*, which works on spectrograms, and find that SpecGAN has a higher Inception Score (IS) and label accuracy through humans, but WaveGAN produces samples of higher sound quality.

Short overview of other usages of GANs:

| Approach | Description | Year |
|---|---|---|
| [272] | Generative Multi-Adversarial Networks (GMANs): The first introduction of DCGAN [267] with multiple discriminators demonstrated on image generation tasks, resulting in higher-quality images and robustness to mode collapse. | 2016 |
| [273] | Mode regularized GANs for stable training and reduced risk of mode collapse: Train an encoder $E(x) : X \rightarrow Z$ together with generator $G(z) : Z \rightarrow X$ and add a similarity loss $distance(x, G(E(x)))$ for stable training gradients. Further, a mode regularizer objective $D_1(G(E(x)))$ on the training data $x$ is employed to force the generator to cover the whole data space. After training with $x$, a second discriminator $D_2$ discriminates $G(z)$ and $G(E(x))$ to bring both distributions to the same manifold efficiently. | 2016 |
| [274] | 3D-GAN: A CNN generator creates objects in 3D voxel space from a random vector, and the discriminator, which is a reversed generator, judges whether the objects are real. The model learns mappings between low-dimensional latent vectors and 3D objects, which can be used, apart from generative purposes, for object recognition and description. They further introduce 3D-VAE-GAN, which consists of a convolutional image encoder and allows one-shot image-conditional 3D model generation. | 2016 |
| [275] | TextGAN: Text generation with a LSTM generator and a CNN discriminator/encoder. Instead of the standard GAN objective, the generator is trained to minimize the covariance matrices of the real and synthetic sentence feature vectors obtained from the CNN encoder. The discriminator is trained with the standard adversarial loss, latent code reconstruction loss, and generator loss. The generator is pre-trained in an autoencoder LSTM setting, while the discriminator is pre-trained to classify sentences with swapped words from true sentences. In [276], they go into further detail and explore a MMD-based feature matching loss instead of covariance. | 2016 |
| [277] | Coupled GAN (CoGAN): Multiple generator-discriminator pairs with shared weights in the higher abstraction layers allow learning of image relations without tuples of corresponding images. The model learns to generate corresponding images for different domains from the same noise vector **z**, making the model suitable for unsupervised domain adaptation and image transformation (e.g., learning the relationship between depth and color in RGBD images). | 2016 |
| [278] | Energy-based GAN (EBGAN): Defining the discriminator as an energy function [279] that assigns low energies to data points near the data manifold and vice versa. The discriminator is implemented as an autoencoder, and the reconstruction error is the energy function. During training, the EBGAN is more stable than regular GANs and generates high-resolution images. | 2016 |
| [280] | C-RNN-GAN: A model with deep LSTM generator and discriminator to create continuous sequential data with one or more values at each step. The adversarially trained model generates polyphonic music with a sense of timing and variation but is distinguishable from real music. | 2016 |
| [281] | Introspective adversarial network (IAN): A hybrid of a VAE for representation learning and reconstruction and a GAN to improve VAE performance. The network combines GAN discriminator and VAE encoder and GAN generator and VAE decoder. Used for realistic photo editing/interpolation. | 2016 |
| [101] | InfoGAN: An unsupervised GAN that learns disentangled representations with a mutual information objective on latent variable subsets. The model can separate writing styles from digits on MNIST and hairstyles, emotions, and eyeglasses from faces on the CelebA data set. | 2016 |

Continuation...

... Continuation

| [282] | f-GAN: Training of generative neural samplers (probabilistic feedforward neural networks), which efficiently convert a random input vector to a sample, using an auxiliary discriminative neural network or other $f$-divergences (e.g., Kullback-Leibler, Jensen-Shannon). Results, also with DCGAN [267], show that the discriminator approach does not necessarily perform better. | 2016 |
|---|---|---|
| [283] | Generative Recurrent Adversarial Network (GRAN): A recurrent attention-based generator/decoder applies sequential changes to a canvas based on the previous hidden state and random noise at each time step, while an attention-based encoder/discriminator provides a feature-based loss. An evaluation method involving a "battle" between separately trained generators and discriminators is proposed, where GRAN outperforms DRAW [97] and the denoising VAE [105]. | 2016 |
| [284] | Offering a suite of methodologies for improved training of GANs. These techniques, such as feature matching, minibatch discrimination, and historical averaging, address the challenges of GAN training, particularly in achieving convergence. The authors further contributed to stabilizing GAN training and proposed a novel evaluation metric now widely used for GAN performance evaluation, the IS, to assess sample quality. | 2016 |
| [285] | Wasserstein GAN (WGAN): Replacement of the discriminator with a *critic* that approximates the Earth-Mover distance (Wasserstein-1) between the real data distribution and the generator data distribution. Since the Wasserstein distance is continuous and differentiable under most circumstances, the generator can be trained with gradient descent, solving the training instability and mode collapse problems of GANs and providing meaningful learning curves, as demonstrated on image generation tasks. In [286], the training is improved with a gradient penalty with respect to the critic input. | 2017 |
| [287] | Boundary Equilibrium GAN (BEGAN): Similar to EBGAN [278], an autoencoder is used as a discriminator, but instead of using the reconstruction loss of samples directly, the Wasserstein distance between error distributions of real and generated samples is computed. Further, the equilibrium condition between the expected errors of generated and real samples is relaxed with a hyperparameter $\gamma \in [0, 1]$ that influences the diversity of the generated images. BEGAN outperforms previous models in image quality at higher resolutions. | 2017 |
| [288] | Least Squares GAN (LSGAN): Using the least squares loss for the discriminator instead of the sigmoid cross-entropy loss, resulting in higher quality images generated and a more stable training procedure compared to regular GANs. The new loss penalizes correct decisions far beyond the decision boundary (i.e., too much certainty by the discriminator) to prevent vanishing gradients and move the generator towards the decision boundary, which converges to the real data manifold. | 2017 |
| [289] | RankGAN: Instead of binary judgments of individual data samples by the discriminator, generator outputs are mixed with real data and ranked by the discriminator according to a reference. The more detailed feedback allows the generator to learn better what makes data realistic. The model generates natural language sentences but could be extended for image generation and captioning. | 2017 |
| [290] | Chekhov GAN: Treating GAN training as a zero-sum game that can be solved by finding a mixed strategy. Using ideas from online learning, where a player aims to minimize a sequentially-revealed cumulative loss function, a no-regret strategy for both generator and discriminator that incorporates the history of the model's actions is employed. The model improves stability and mode collapse problems and is guaranteed to converge to equilibrium for specific GAN architectures. The model is evaluated on image generation and density estimation tasks. | 2017 |
| [291] | Adversarially Regularized Autoencoder (ARAE): A framework that trains a Wasserstein GAN [285] to produce latent codes of a simultaneously trained autoencoder to create a GAN for discrete data (e.g., images, text). | 2017 |
| [292] | Mean and covariance feature matching GAN (McGAN): Training GANs by matching statistics (e.g., embedded mean or covariance of features) of real and fake data instead of classifying individual samples. The approach is adapted to DCGAN [267] and used to generate images. | 2017 |
| [293] | Fisher GAN: Instead of penalizing the gradients of the Wasserstein GAN discriminator [286], a constraint is imposed on its second-order moments, inspired by the *Fisher Discriminant Analysis* method. The approach is applied to a DCGAN [267] architecture, outperforming other models in unconditional image generation. | 2017 |

Continuation...

... Continuation

| [294] | $\alpha$-GAN: Combining an autoencoder with a GAN, where the GAN generator is trained to reconstruct latent representations of an encoder and the latent space is encouraged to conform to a Gaussian distribution. Further, a discriminator is used to evaluate the realism of the generated or reconstructed samples. The combination of autoencoders and GANs (DCGAN [267] in particular) solves blurriness and mode collapse issues for image generation. | 2017 |
|---|---|---|
| [295] | SeqGAN: Application of the GAN architecture to the generation of sequences of discrete tokens. Since the default GAN generator generates continuous values and receives instant feedback from the discriminator, and the discriminator can only evaluate complete sequences, the authors model the generator as a LSTM and RL agent with a stochastic policy whose intermediate actions are rewarded after the complete sequence has been judged by the CNN discriminator. Lee et al. [296] applied SeqGAN to polyphonic music generation using efficient representations of polyphonic MIDI files. | 2017 |
| [297] | Temporal GAN (TGAN): The video generator consists of a temporal generator CNN that generates $T$ latent variables $z_1^t$ from initial noise $z_0$ and an image generator that generates $T$ video frames from $z_0$ and the respective $z_1^t$. The discriminator receives all frames and evaluates, with multiple convolutional layers, whether the video is real or generated. Image generator and discriminator are very similar to DCGAN [267], and training also incorporates the WGAN [285] objective. | 2017 |
| [298] | VEEGAN: Using a reconstructor network $F(x)$ to map the true data distribution $p(x)$ back to the Gaussian distribution $p(z)$ of the generator $G(z)$'s latent code. The Kullback-Leibler divergence serves as an expressive loss function between distributions $F(G(p(z)))$ and $p(z)$ because they should be identical. This autoencoder-styled structure aims to solve the mode collapse problem of GANs by checking if this whole distribution is covered. The model is demonstrated on density estimation and image generation tasks, where it is less prone to mode collapse issues than other GAN approaches. | 2017 |
| [299] | Boundary-seeking GAN (BGAN): Using the estimated difference measure of a discriminator as a differentiable policy gradient for the generator to allow GANs to generate discrete data. The approach is also suitable for continuous data and is demonstrated in natural language and image generation. | 2017 |
| [300] | medGAN: Generation of synthetic Electronic Health Records (EHRs) for privacy-preserving data sharing for medical research. An autoencoder provides discrete real data $\mathbf{x}$ or discrete data from noise $\mathbf{z}$ to a discriminator $D$ in the form $Dec(Enc(\mathbf{x}))$ or $Dec(G(\mathbf{z}))$ respectively, where generator $G$ and $D$ are feedforward neural networks. | 2017 |
| [301] | MidiNet: A generator CNN creates symbolic MIDI melodies from random noise and prior knowledge (e.g., chord progression, previous bars, priming melody) from a conditioner CNN with *transposed convolutions* [302] and a discriminator CNN predicts whether an input score is real or fake. | 2017 |
| [303] | Progressive GAN: Accelerating and stabilizing GAN training by progressively adding convolutional layers to the generator and discriminator, increasing the resolution for image generation iteratively. | 2017 |
| [304] | DRAGAN: Treating GAN training as a zero-sum game solved using a regret minimization technique. Mode collapse is interpreted as a problem caused by local equilibria in the training "game". The model is stable and outperforms Wasserstein GAN [285, 286] significantly. | 2017 |
| [305] | Gang of GANs (GoGAN): Improvement over Wasserstein GAN [285] by using a margin-based discriminator loss that enforces a certain distance $\epsilon$ between discriminator scores of fake and real samples and progressively training multiple GANs that are evaluated against each other with a maximum margin ranking loss that ensures that later GANs perform significantly better than earlier ones. The model is used for image completion. | 2017 |
| [306] | LeakGAN: The CNN discriminator leaks its high-level abstracted features of the currently evaluated sentence to the LSTM generator that takes the additional input as guidance for next word generation for a text. | 2017 |
| [307] | Combining a classic GAN architecture with an inference network. This allows the discriminator to distinguish between joint latent/data-space samples from the generative network and joint samples from the inference network (which receives a data sample as input and outputs synthetic data). The setup enhances the performance of state-of-the-art tasks. | 2017 |

Continuation...

... Continuation

| [308] | LR-GAN: Generating images by creating backgrounds and foregrounds separately and then stitching them together. It employs a recursive approach where each layer (background or foreground) is generated step-by-step, each with its own shape and pose. This method allows for a contextually relevant composition of images, leading to more natural and realistic image generation. | 2017 |
|---|---|---|
| [309] | Enhancing GANs using denoising feature matching. This technique guides the generator towards more probable configurations of abstract discriminator features, generating more object-like samples. The approach uses a denoising auto-encoder to estimate and track the distribution of these features derived from real data. This is combined with the original GAN loss, and the augmented training procedure is shown to improve its stability. | 2017 |
| [286] | Applying Wasserstein GAN for data generation, using a gradient penalty to enforce a Lipschitz constraint on the discriminator. This addresses the training instability associated with weight clipping in Wasserstein GAN. The method enables stable training across various GAN architectures, significantly reducing the need for hyperparameter tuning. | 2017 |
| [310] | In contrast to GAN training training, the authors update the discriminator and generator at different rates. This method facilitates more stable convergence and enables the networks to reach a local Nash equilibrium effectively. Additionally, the introduction of the Fréchet Inception Distance (FID) provides a more reliable and sensitive metric for GAN performance evaluation when compared to the IS. The effectiveness was validated using several GAN architectures, improving stability and image quality in the generated samples. | 2017 |
| [311] | Motion and content decomposed GAN (MoCoGAN): Video generation using a fixed content vector and a sequence of motion vectors that are converted to a state using a RNN. An image generator CNN then creates a frame from both vectors, and CNN-based video and image discriminators provide feedback. | 2018 |
| [312] | MuseGAN: A WGAN [285] utilizing deep CNNs that generate multi-track music sequences with piano-roll representations. It combines a jamming model that improvises one track with a composer model that creates multiple accompanying tracks at once, so multiple hybrid generators with inter-track random vector $\mathbf{z}$ and intra-track random vectors $\mathbf{z}_i$ are paired with one discriminator. | 2018 |
| [313] | CapsuleGAN: Replaces the CNN discriminator in a DCGAN [267] with a capsule network (CapsNet) [314] classifier to improve accuracy and generator performance. | 2018 |
| [315] | MolGAN: Training a MLP generator to produce adjacency and node annotation matrices for small molecules with a graph-convolutional discriminator and reward network. The reward network, like in RL, learns to score non-differentiable metrics (e.g., solubility in water) with the help of external software and to guide the molecule generation towards a specific target. | 2018 |
| [316] | NetGAN: Generating graphs as a set of random walks (node sequences) with a LSTM generator outputting vertex after vertex of the sequence and a discriminator LSTM that judges whether the sequence belongs to the real graph or is fake. Training samples (random walks) are obtained from a real graph, and synthetic graphs are obtained by merging the random walks. | 2018 |
| [317] | table-GAN: Application of the DCGAN [267] architecture to the generation of private and useful tabular data. Original table entries are fed as square zero-padded matrices to the GAN for training, with a neural network classifier $C$ besides the generator and discriminator that enforces data consistency learned from the original table. | 2018 |
| [318] | corrGAN: Generation of correlated discrete data (e.g., table entries, binary images) with an autoencoder that learns mappings between discrete input/output and continuous latent space while considering the correlations between subsets of the discrete variables. The decoder then serves as the output layer of a GAN generator modeling the continuous latent space. | 2018 |
| [319] | Tabular GAN (TGAN): A LSTM generates discrete and continuous values encoded with probability distributions column by column for each entry, and a MLP discriminator scores the likelihood and diversity of the data. | 2018 |
| [320] | Consistency Term (CT) GAN: Improved training of improved Wasserstein GANs [285, 286] by additionally training the discriminator on once and twice perturbed real data and evaluating the responses. Combined with other small optimizations, this approach achieves state-of-the-art generative results and is better suitable for semi-supervised learning than other GANs. | 2018 |

... Continuation

| [321] | AmbientGAN: Training a GAN on lossy measurements by passing the generator output through a simulated random measurement function that corrupts it. The discriminator then tries to distinguish the lossy real and fake measurements. The model can successfully inpaint lossy images. | 2018 |
|---|---|---|
| [322] | Protein structure generation and completion by encoding them as pair-wise distances between $\alpha$-carbons in a matrix using a DCGAN [267] and recovering the 3D structure using the "alternating directional method of multipliers" (ADMM) algorithm. | 2018 |
| [323] | Human trajectory generation as a sequence of stays with a location $(x, y)$, a start time $t$, and a duration $d$. A GAN made of CNNs creates and evaluates "maps" with the aforementioned times as coordinate values. | 2018 |
| [324] | MGAN: Tackling the mode collapse in GANs by using multiple generators. This method employs a mixture of generators with a shared classifier and discriminator, aiming to produce diverse outputs that cover different data modes. The authors demonstrate that their approach effectively minimizes the Jensen-Shannon Divergence between the mixed generator distributions and the real data distribution. Empirical results showed the model's ability to generate diverse and recognizable images, indicating a significant improvement over single-generator GAN. | 2018 |
| [325] | RelGAN: Text generation with a configurable trade-off between sample quality and diversity. The recurrent generator incorporates relational memory [224] (multiple memory slots with self-attention [18]) to model long-range dependencies, and the CNN discriminator creates multiple representations for each sentence to evaluate the sentence from different aspects and provide better feedback. | 2018 |
| [326] | Dist-GAN: A GAN enhanced with distance constraints. It addresses two key issues in GAN training: gradient vanishing and mode collapse. The novel approach includes integrating an autoencoder with the generator and implementing two distance constraints, one in the latent space and another based on discriminator scores. This stabilizes the GAN training while retaining competitive scores in classification tasks. | 2018 |
| [327] | Introducing spectral normalization, a technique for stabilizing GAN training. It normalizes the spectral norm of the weight matrices in the discriminator, using a Lipschitz constant as the only hyper-parameter to be tuned. This method is simple, computationally efficient, and stabilizes GAN performance, particularly in image generation tasks on datasets. | 2018 |
| [328] | Bayesian modeling is employed to tackle mode collapse in GANs, a well-documented and actively researched problem. The authors propose learning the distributions of generators, as it mirrors a common approach of training a GAN with multiple generators to alleviate the mode collapse. | 2018 |
| [329] | Using the Sinkhorn divergence to train generative models based on regularized optimal transport with an entropy penalty. This method addresses the computational and statistical challenges of using optimal transport metrics in training generative models. The Sinkhorn divergence interpolates between Wasserstein and MMD losses, leveraging the geometrical properties of the optimum transport and the favorable high-dimensional sample complexity of MMD. | 2018 |
| [330] | StyleGAN: A style-transfer related progressive GAN [303] capable of unsupervised learning of disentangled high-level attributes of images in latent space and generation of realistic images from that latent space and random noise. The model can also interpolate (style mixing) in the latent space. | 2019 |
| [331] | PATE-GAN: Differentially private data generation with the original GAN utilizing the private aggregation of teacher ensembles (PATE) [332] framework which allows to bound the influence of individual samples on the generator. A set of *teacher* discriminators is trained on equal amounts of generator output and disjoint training data. A *student* discriminator is trained on the output labels of the teachers, and the generator is trained on the student's output. | 2019 |
| [333] | Improved text generation with GANs by segmenting sentences into sub-sequences and providing simultaneous discriminator feedback for all sub-sequences and the entire sentence instead of the sentence alone. Applying this approach to previous state-of-the-art GANs [295, 306, 325] significantly improves their results. | 2019 |

Continuation...

... Continuation

| | | |
|---|---|---|
| [334] | SemGAN: Generation of pixel-level semantic images from latent vectors with prior distribution to speed up convergence. Class probabilities for each pixel are forwarded as a softmax distribution to the discriminator, allowing for detailed feedback. The SemGAN produces significantly cleaner results than a default GAN working with RGB semantic mappings. Using the image-to-image GAN from [335], the creation of natural images from these artificial semantic images is demonstrated. | 2019 |
| [336] | Medical Wasserstein GAN (medWGAN) & Medical Boundary-seeking GAN (medBGAN): Modified versions of medGAN [300] with better performance. medBGAN outperforms all other approaches on synthetic EHR generation. | 2019 |
| [337] | Using a default GAN to generate images of skin cancer, a domain for which only a few labeled samples are available. The synthetic data is used to boost the performance of a cancer detection CNN. | 2019 |
| [338] | AutoGAN: An automated approach of solving neural architecture search specifically tailored to GAN architectures. To do so, an additional RNN controller is incorporated for the search process, while the IS is used as the success reward of the reinforcement. | 2019 |
| [339] | DoppelGANger: A generative model for creating synthetic time series data tackling several issues about GANs and data privacy. This model effectively captures and generates complex correlations between time series and its attributes while addressing challenges such as mode collapse and variable data lengths. | 2019 |
| [340] | TimeGAN: Generating realistic time-series data by combining supervised and adversarial training, addressing the challenge of preserving temporal dynamics in generated sequences. It features an embedding network for dimensional reduction, enhancing the generative model's learning of temporal relationships. | 2019 |
| [341] | StyleGAN2: Improvement of StyleGAN [330] by simplifying and restructuring the generator architecture. They further employ lazy regularization and remove the progressive growing in favor of skip connections between up- and downsampling layers. They also developed a method to project images back to latent space where style could be changed/interpolated. | 2020 |
| [342] | HealthGAN: A feedforward neural network based GAN that combines ideas from medGAN [300] with the Wasserstein GAN gradient penalty [286] and data transformations to enable the use of continuous and categorical data. The model is purposefully small to prevent memorization and used to generate private EHR records. | 2020 |
| [343] | Using a progressively growing GAN (similar to [303]) with Wasserstein gradient penalty loss [286] to generate magnetic resonance imaging of brains that can be used for Attention Deficit Hyperactivity Disorder prediction. | 2020 |
| [344] | Using DCGAN [267] to generate positron emission tomography brain images for three stages of Alzheimer's disease to build an automated diagnosis model. | 2020 |
| [345] | SynSigGAN: Generation of fixed-length labeled biomedical signals (electrocardiogram and other time series data) using a bidirectional grid LSTM and a CNN discriminator with a small amount of training data. The data can be used for automatical medical diagnosis or training medical students and achieves state-of-the-art performance when used as training data for a classifier. | 2020 |
| [346] | By analyzing the loss of the generator and the discriminator during the training process, overarching calculations such as IS or FID can be eliminated. This reduces the computing time greatly and optimizes the network at the same time, especially compared to either a manual GAN configuration or similar automated approaches based on the IS oder FID. | 2020 |
| [347] | COT-GAN: A generative model for sequential data employing causal optimal transport for training implicit generative models, integrating classic optimal transport methods with a temporal causality constraint. The COT-GAN framework is adept at generating low- and high-dimensional time series data. | 2020 |
| [348] | Using differentiable augmentations applied to real and fake samples during training, this approach effectively prevents overfitting in GAN and reduces the training size. The method shows notable improvements across various GAN architectures and achieves state-of-the-art results while notably only using 20% of the training data. | 2020 |
| [349] | Using DCGAN [267] to boost the amount of training data for traffic sign recognition, resulting in increased accuracy and reduced detection time. | 2021 |

... Continuation

| [350] | Combining the boundary-seeking GAN [299] with *noise addition* directly on the data to generate differentially private smart health care data sets of populations for medical research. Training data is obtained from Fitbit smartwatches. | 2021 |
|---|---|---|
| [351] | Introducing a method for generating synthetic Electrocardiograms (ECGs) using GANs to address privacy concerns in medical data sharing. The authors presented two GAN models, WaveGAN and Pulse2Pulse, trained on real normal ECGs to produce synthetic, plausible ECGs. The presented approach allows for generating an arbitrary amount of normally very sensitive patient data and open-sources the over 100,000 normal ECGs. | 2021 |
| [352] | TTS-GAN: Generating synthetic time-series data using transformer architecture. It employs transformer encoders in both generator and discriminator networks, overcoming limitations of RNN-based GAN in handling long sequences. It showcases improved performance in generating realistic sequences across multiple datasets. | 2022 |

### 2.13.1 Conditional GANs

Normal GANs have no control over the type of data the generator outputs. Conditional GANs solve this problem by conditioning generator and discriminator on additional information $\mathbf{y}$, which could be class labels, for example, that are provided at the input layer (see Figure 32). The objective function of the GAN is modified as follows [353]:

$$\min_G \max_D V(D, G) = \mathbb{E}_{\mathbf{x} \sim p_{data}(\mathbf{x})}[\log D(\mathbf{x}|\mathbf{y})] + \mathbb{E}_{\mathbf{z} \sim p_z(\mathbf{z})}[\log(1 - D(G(\mathbf{z}|\mathbf{y})))]. \tag{28}$$
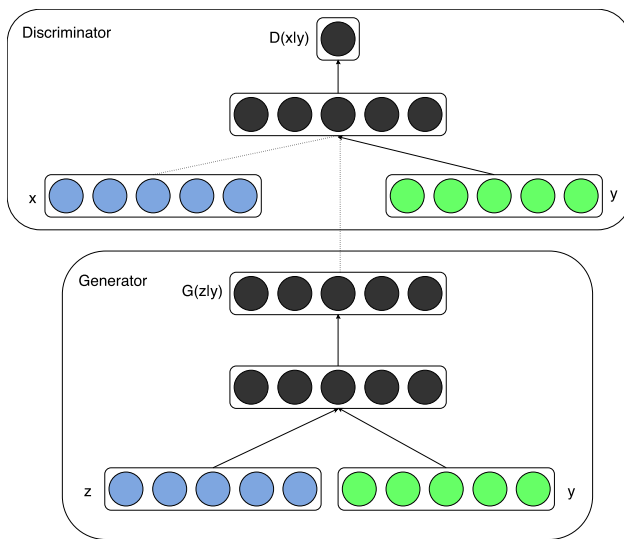


Figure 32: A conditional GAN conditioned on additional input $\mathbf{y}$. (Source: [353])

Mathieu et al. [354] adopt the GAN architecture to predict the next video sequence frames and introduce a loss function that improves the sharpness of the image predictions. The generator and discriminator use multi-scale networks, which are CNNs that up- or down-scale the resolution of a prediction multiple times, respectively, until the target image size or a single scalar output is reached. The model is trained by providing a sequence of video frames to both models and the real or generated additional frames to the discriminator. The model parameters are updated using stochastic gradient descent. They achieve state-of-the-art sharpness and similarity on a collection of sports clips from the Sports1m and UCF-101 data sets.

Zhu et al. [355] propose using the DCGAN [267] architecture to learn the approximate natural image manifold by training the generator to produce realistic images given a 100-dimensional random vector $z$. Then, real images are projected to such a latent representation $z_0$, and manipulating operations (coloring, sketching, or warping executed with brush tools) are applied to that vector, resulting in $z_i$. The changes in the artificially generated images $G(z_i)$ are then sequentially transferred to the original photo using optical flow methods to maintain image quality. Besides interactive image manipulation, the *iGAN* model can also generate objects from drawn sketches.

Reed et al. [356] generate images from textual descriptions (sentences) using a DCGAN [267] architecture. First, a text embedding encoder model (hybrid character-level convolutional recurrent neural network) is pre-trained by comparing its embeddings to the ones encoded by a corresponding image encoder (deep CNN). The GAN consists of a deep convolutional generator that takes as input the text embedding and random noise and a deep convolutional discriminator that decides whether the image is real or fake based on the text embedding and a provided image. They also test a modified discriminator with additional real training samples with mismatched text to condition the model on matching texts in addition to image realism. A further addition to the generator training objective is the interpolation of text representations to ensure that gaps in the training data also correspond to the data manifold:

$$\mathbb{E}_{t_1, t2 \sim p_{data}} [\log(1 - D(G(z, \beta t_1 + (1 - \beta)t_2)))] \tag{29}$$

with noise sample $z$, text embeddings $t_1$ and $t_2$ and $\beta = 0.5$, which works well as long as the discriminator can recognize the matching image and text pairs. Further, the author developed a style encoder, which, in combination with the trained generator, allows them to combine the style of one image and a description to generate a same-style image with properties matching the text (e.g., an image of an eagle with the description of a red bird results in a red eagle).

In [357], Reed et al. introduce the Generative Adversarial What-Where Network (GAWWN) as an extension of [356], that allows specification of locations (bounding-boxes) of objects and their parts (keypoints) for text-to-image synthesis. They present two models, the bounding-box-conditional and the keypoint-conditional text-to-image model, and apply a GAN architecture to generate missing keypoints given some user-defined keypoints and the text description or text alone, which they demonstrate also works.

Isola et al. [335] perform image-to-image translation using a conditional DCGAN [267] with a generator $G : \{x, z\} \to y$ and a discriminator $D : \{x, y\} \to [\text{real}, \text{fake}]$, where $x$ is the "concept" of an image (e.g., an edge drawing or blueprint), $y$ is the real or generated image and $z$ is random noise, which is later replaced by dropout on several layers because $G$ tends to ignore $z$. The generator is a deep CNN which first down- and then up-samples $x$ to obtain essential features and create a fitting image in the target "style". The discriminator *PatchGAN* is a convolutional model that classifies if each $N \times N$ patch of an image is real or fake and averages all outputs. The loss function combines the discriminator output and the L1 loss (mean absolute error) of the generated image and the ground truth. The model is then trained using alternating gradient descent steps on $D$ and $G$.

Choi et al. [358] introduce *StarGAN* based on [359] for multi-domain image-to-image translation using one generator and multiple discriminators, one for every pair of image domains (see Figure 33a). This allows StarGAN to be trained on multiple data sets with different labels simultaneously and combine the information learned into one generative model.



(a) Overall architecture of StarGAN with one generator and multiple discriminators.

(b) **Right:** Discriminator training. **Left:** Generator training with reconstruction loss **(c)** and adversarial loss by the discriminator **(d)**.
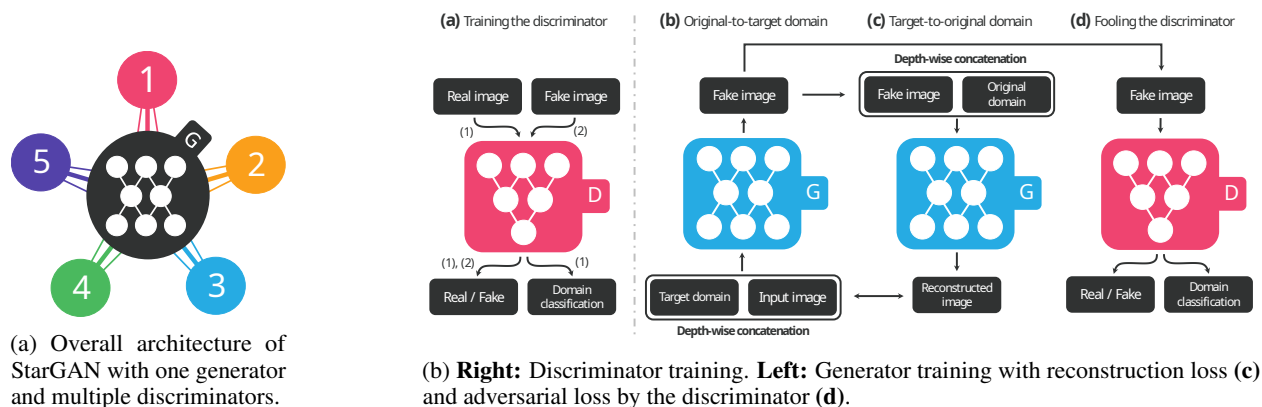
Figure 33: Illustrations of StarGAN. (Source: [358])

Kocaoglu et al. [360] propose *CausalGAN* (Figure 34a), whose generator is structured according to a given causal graph, similar to a BN (see Figure 34b). Given binary labels and some real observations for the discriminator, the model can be conditioned to generate data, for example, face images, according to the labels (e.g., sex, hair color). First, the *causal controller*, implemented as a Wasserstein GAN [285], generates the labels. Then, the generator produces data according to the labels and noise. Finally, the generator competes with three adversaries to produce realistic samples (discriminator) with correct labels (labeler) while avoiding easy-to-label unrealistic image distributions (anti-labeler).

(a) Overall architecture of CausalGAN.

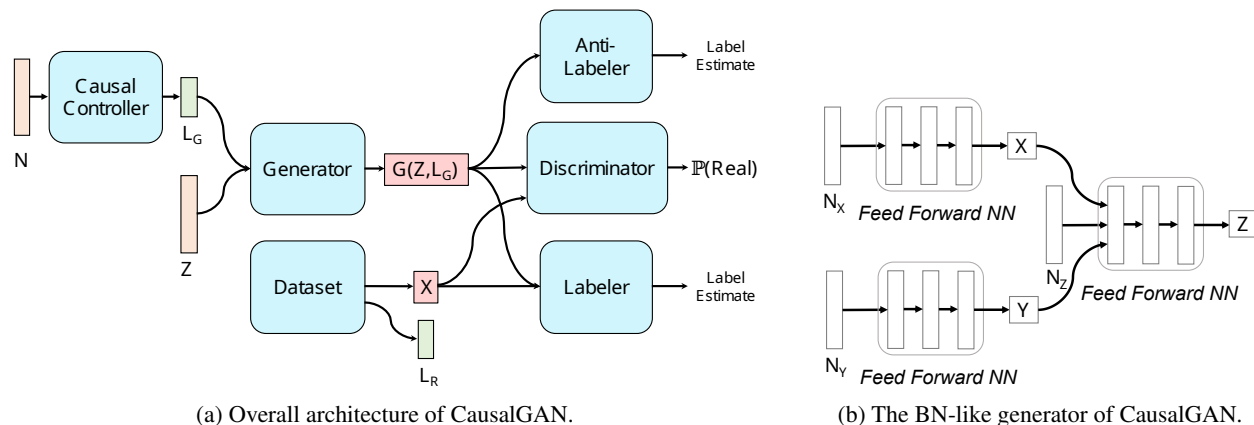(b) The BN-like generator of CausalGAN.

Figure 34: Illustrations of CausalGAN. (Source: [360])

Zhang et al. [361] develop the self-attention GAN (SAGAN) to combine the computational and statistical efficiency of a convolutional GAN for class-conditional image generation with the long-range dependency modeling capability of self-attention mechanisms (see Section 2.12). SAGAN significantly outperforms previous state-of-the-art models [362] by increasing the IS from 36.8 to 52.52 and decreasing the FID from 27.62 to 18.65.

Krishna et al. [363] propose a conditional GAN to merge the style of one Computed Tomography (CT) image with the content of another to solve the problem of scarcity and privacy issues in clinical training data. They introduce a convolutional encoder-decoder generator trained to map the style of a CT image to the segmentation map of other images per organ influenced by random noise. The training requires only a small data set and combines style and content loss to guide the generator and the convolutional discriminator.

Alonso et al. [364] extend a GAN to generate images of handwritten words. They use a 4-layer bidirectional LSTM to create an embedding of the target character sequence fed to the generator. Further, an auxiliary text recognition network consisting of a CNN encoder and LSTM decoder evaluates the generated image in addition to the discriminator to encourage faithful recreation of the target word. The resulting GAN produces realistic images of French and Arabic words that improve text recognition results of a neural network when used as additional training data.

Brock et al. [365] build upon SAGAN [361], increase the batch size by factor 8 and the width of each layer by 50%, doubling the number of parameters in the generator and the discriminator. Their *BigGAN* achieves new state-of-the-art results in class-conditional image synthesis with high resolutions up to $512 \times 512$ and a controllable trade-off between detail and variety.

Lučić et al. [366] aim to train the BigGAN [365] with unlabeled data for conditional image generation. They experiment with unsupervised clustering and linear classifiers (pre-trained or co-trained with the discriminator) on top of representations of the images to substitute large parts of the labels (up to 90%), achieving similar or better inception and FID scores compared to the normal BigGAN.

Feng et al. [265] introduce CA-GAN, a symmetric and convolutional GAN that implements collaborative learning between the generator and discriminator to provide real sample information to the generator and an attention mechanism for the generator to remove spurious features. The generator and discriminator interact at multiple convolution steps with the attention mechanism. CA-GAN is used to generate label-conditional high-quality Hyperspectral Images (HSIs), which are pictures with multiple layers of features per pixel, with the help of a limited sample set and classify them.

Short overview of other usages of conditional GANs:

| Approach | Description | Year |
|---|---|---|
| [353] | First introduction of conditional GANs on MNIST digit and image tags generation. | 2014 |
| [367] | Conditional face image generation with the conditional GAN [353] architecture with deconvolutional generator and convolutional discriminator. The conditional vector allows age specification and other attributes found in the training data. To avoid the reproduction of the training data by the generator, the conditional vectors are not directly used but are sampled from a kernel density estimate of the training data. | 2014 |

Continuation...

... Continuation

| [368] | Utilizing an objective function balancing mutual information between observed examples and predicted class distributions against the robustness to an GAN. This approach extends the regularized information maximization to robust classification against an optimal adversary. The study includes empirical evaluations of synthetic data and image classification tasks, demonstrating the robustness of learned classifiers and the fidelity of samples generated by the adversarial generator. | 2015 |
|---|---|---|
| [369] | Invertible conditional GAN (IcGAN): Using two encoders to create an inverse mapping from images to a latent representation $z$ and conditional attributes, allowing modifications of images with a conditional DCGAN [267] whose generator also functions as the decoder. The model is evaluated with different configurations, enabling realistic and complex image modifications. | 2016 |
| [370] | Attribute-Layout Conditioned GAN (AL-CGAN): Generating images of outdoor scenes from semantic layouts and scene attributes (e.g., weather) with a DCGAN [267]. | 2016 |
| [371] | GAN for video (VGAN): A spatio-temporal DCGAN [267] architecture for unconditional and conditional video generation (e.g., future prediction from static images) with moving objects and a static camera. Generates one-second-long videos that resemble real videos with separate foreground and background generators. | 2016 |
| [372] | SimGAN: Use a GAN to refine synthetic data from a simulator and make it realistic while preserving the simulator's annotations. An additional loss is employed to minimize the applied per-pixel difference of the CNN refiner. | 2017 |
| [373] | SeGAN: Segmentation and painting of occluded parts of an image. A segmentation CNN creates the complete segment mask from a partial mask obtained by another pre-trained model [374] and the RGB image. A conditional GAN (similar to [335]) then paints the invisible parts of the object based on the masks and the image. | 2017 |
| [375] | Training a GAN generator to learn a pixel-level image transformation from one domain to another (e.g., depth for an RGB image) using convolutions and fully connected layers. The discriminator is structured similarly and outputs the probability that an image was sampled from the target domain. For some tasks, the generator is further conditioned using a content similarity loss that penalizes large differences between the input and output of the generator. | 2017 |
| [376] | SegAN: Training a CNN to produce segmentation masks for medical images (CT images) and an adversarial critic network compares the generated mask to the ground truth by evaluating features extracted with a CNN. | 2017 |
| [377] | DualGAN: Training two convolutional GANs [335] to translate images between two domains in both directions from unpaired data. Also, since both directions' models are available, a reconstruction loss is employed to improve training. Instead of a dedicated noise vector $z$, dropout is used in multiple layers. | 2017 |
| [378] | Steganographic GAN: Using a DCGAN [267] generator to encode encrypted messages in natural-looking images, a discriminator to ensure data realism, and a steganalyzer network to retrieve the message using a shared key. | 2017 |
| [379] | SRGAN: Generator and discriminator consisting of multiple convolutional layers generate super-resolution (factor $4\times$) images. The GAN is additionally trained on a perceptual loss obtained from feature maps of a VGG network [186], a deep CNN for image classification. | 2017 |
| [380] | Variational autoencoding Wasserstein GAN (VAW-GAN): Voice conversion with a conditional VAE encoding phonetic content from speech parameters (spectral frames) and decoding it depending on speaker identity. The VAE decoder is then treated as the generator of a Wasserstein GAN [285] and trained to generate clearer speech. | 2017 |
| [381] | Triple-GAN: Jointly train a generator, discriminator, and classifier for semi-supervised learning, resulting in state-of-the-art classification results and smooth class-conditional generation and interpolation in the latent space. | 2017 |
| [382] | Recurrent GAN (RGAN) and Recurrent Conditional GAN (RCGAN) for real-valued (medical) time series generation. The conditional version has additional inputs besides noise and data for the generator and discriminator LSTMs, respectively. The models are evaluated using an evaluation scheme, trained on synthetic data, and evaluated on real data. | 2017 |
| [383] | Multi-Agent Diverse GAN (MAD-GAN): Using multiple generators and one discriminator that also has to identify the generator that created the sample to generate images with different classes in a conditional and unconditional setting. The discriminator pushes certain generators to different classes (modes) and improves the overall performance. | 2017 |

Continuation...

... Continuation

| [384] | Variational InfoGAN (ViGAN): Combining VAEs with InfoGAN [101], using the VAE encoder to generate a latent representation $z$ of an observation that combined with a set of interpretable variables $c$ allows to generate/decode modified images $\tilde{x} \sim P(x|z,c)$. These images are then evaluated by a discriminator and a recognizer that reconstructs $c$. | 2017 |
|---|---|---|
| [362] | Auxiliary Classifier GAN (AC-GAN): A class-conditional deconvolutional generator produces images, and the convolutional discriminator outputs the real/fake probabilities and a probability distribution over all classes for input images. The objective function encourages both models to maximize the class likelihood, while the generator aims to minimize the correct real/fake judgments that the discriminator tries to maximize. The model works better on higher-resolution images and achieves a state-of-the-art IS on CIFAR-10. | 2017 |
| [385] | Structure Correcting Adversarial Network (SCAN): A CNN generator learns segmentation masks of chest x-rays, and a critic network (CNN and fully connected network for classification) learns to discriminate between real and generated masks for an image. To mitigate mode collapse, the generator is pre-trained with pixel-wise loss. | 2017 |
| [386] | Image manipulation with text descriptions using a generator that consists of a CNN image encoder, a pre-trained RNN text encoder, a residual transformation unit that jointly encodes text and image embeddings further and a decoder that reconstructs the image with upscaling layers. A CNN discriminator evaluates the probabilities of the generated image and the text description matching. | 2017 |
| [359] | CycleGAN: Unpaired image-to-image translation (i.e., two unrelated sets of images with no relation information provided) by learning mapping $G : X \to Y$ and inverse mapping $F : Y \to X$ and adding a cycle consistency loss $F(G(X)) \approx X$ next to the adversarial loss $G(X) \approx Y$. The model is implemented as CNNs and used for style transfer, object transfiguration, season transfer, and photo enhancement. | 2017 |
| [387] | BicycleGAN: Image translation (e.g., night-to-day image conversion) with invertible latent codes to prevent many-to-one mappings (i.e., mode collapse). For that purpose, a generator/decoder is conditioned on an input image and a latent vector to produce a suitable output image, and a VAE encoder is trained to map the output back to the same latent code and a predefined distribution. | 2017 |
| [388] | Unsupervised image-to-image translation with coupled GANs utilizing a *shared latent space assumption*. Two encoders $E_1$ and $E_2$ translate images from two domains to the same latent space $z$. The generators $G_1$ and $G_2$ serve as decoders of a VAE and recreate images from $z$ suitable for the domain. Their performance is evaluated by discriminators $D_1$ and $D_2$ respectively. The high-level connection weights between the encoders and the generators are tied to account for the shared latent space assumption. | 2017 |
| [389] | Conditional CycleGAN: Extending CycleGAN [359] with an additional attribute condition vector (e.g., hair color, gender, smiling) for face image super-resolution and attribute-conditional translation. | 2017 |
| [390] | Showcasing that GANs can produce higher quality samples by using a conditional setup, more precisely, when class labels are provided. The approach proposed augments class labels by clustering the representation space learned by the model itself. The method is, however, based on the more computationally costly Wasserstein GAN approach. | 2017 |
| [391] | Using the *pix2pix* [335] approach to generate infrared images and videos with tracking labels for entities from labeled RGB video frames. | 2018 |
| [392] | Perceptual Adversarial Network (PAN): An image-to-image transformation CNN $T$ learns mappings between domains and a discriminator CNN $D$ tries to find discrepancies between transformed images and ground truth. | 2018 |
| [393] | Text-adaptive GAN (TAGAN): A GAN that allows the modification of images using natural language descriptions while preserving text-irrelevant features of the original image, for example, the form of an object. The generator is an encoder-decoder architecture derived from [386] that combines the image embedding with the text representation obtained from a bidirectional GRU working on pre-trained fastText [394] word vectors. The text-adaptive discriminator classifies each attribute independently using word-level discriminators. | 2018 |

Continuation...

... Continuation

| [395] | vid2vid: Video-to-video synthesis from segmentation masks or other sources to photorealistic output with GANs and a spatiotemporal adversarial objective. A recursively applied feedforward neural network generates the next frame based on the current source frame and the past source and generated frames. A conditional image discriminator compares the current source and the generated image. Additionally, a conditional video discriminator estimates the optical flow of the generated sequence to ensure plausible temporal dynamics. | 2018 |
|---|---|---|
| [396] | Training a progressive GAN [303] to produce face images conditioned on the identity embedding of a person. The identity is built with discrete labels and converted to a continuous representation that follows a simple distribution. The model can be used for conditional and unconditional image generation and improves and alleviates the training of face recognition models by boosting the training data size with interpolated image reconstructions. | 2018 |
| [397] | Produce identity-preserving photorealistic face images from renderings of a 3D morphable model showing different poses, expressions, and illuminations with a conditional GAN. The GAN uses semi-supervised learning, so it trains with a few pairs of real and rendered face images and has a large amount of unpaired data available. | 2018 |
| [398] | Graph-translation GAN (GT-GAN): A GAN consisting of a (de)convolutional graph translator with an encoder-decoder structure pitted against a conditional graph discriminator that takes the translated/real graph pairs as input and uses the same CNN as the encoder to classify whether they are real or fake. | 2018 |
| [399] | MaskGAN: Training a GAN for text generation. The generator consists of a LSTM encoder that generates a context representation of text with masked tokens and a LSTM decoder that takes the context and the masked text to autoregressively fill in the gaps and reconstruct the sequence probabilistically. The discriminator is a LSTM similar to the decoder and outputs the probability of a token in the sequence to be real. On top of the discriminator, a critic network computes a reward function based on the predicted token likelihoods used to train the generator. The model is also used for unconditional text generation by masking the entire input. | 2018 |
| [400] | Using the *pix2pix* [335] GAN approach to generate images and segmentation maps of brains and tumors. | 2018 |
| [401] | Making the improved Wasserstein GAN with gradient penalty [286] differentially private by adding Gaussian noise to the activations of the second-to-last layer of the discriminator. The model is used to produce label-conditioned images. | 2018 |
| [402] | GANsynth: High-fidelity and locally-coherent audio (music notes) synthesis using high-resolution frequency spectrograms, several orders of magnitude faster than autoregressive models. The generator is additionally conditioned on a one-hot pitch label that the discriminator tries to predict with an auxiliary classifier. | 2019 |
| [403] | Labeled-Graph GAN (LGGAN): An improved Wasserstein GAN [320] approach that uses a MLP generator to produce adjacency and one-hot node label matrices for graphs and a graph convolutional network [84] with residual connections [374] as the discriminator. Both conditional and unconditional configurations are presented. | 2019 |
| [404] | DP-CGAN: Differentially private data generation using a conditional GAN that injects Gaussian noise and clips the discriminator gradients to limit the amount of information from the training data transferred to the generator. The model is trained until the privacy budget monitored by a Rényi differential privacy accountant [405] is spent. The model improves on differentially private results on MNIST over a simpler baseline conditional GAN model. | 2019 |
| [406] | Misc-GAN: Translating graphs from a source to a target domain (e.g., for anonymization) by extracting coarse (partial) structures from real graphs using clustering methods, permutating them, generating new coarse graphs from these templates, and combining the generated coarse graphs to produce a full synthetic graph. | 2019 |
| [407] | COCO-GAN: A conditional GAN only trained using so-called micropatches of image datasets. The generator and discriminator learn only parts of the image via conditional formatting; however, they can generate full images during the inference phase. | 2019 |
| [95] | Conditional Tabular GAN (CTGAN): A GAN to effectively generate tabular data with mixed discrete and continuous columns. The generator is conditioned on a one-hot vector determining a randomly selected category and samples a row (each column independently) from the learned marginal distributions. The critic then scores the generated sample against a random sample from the training data matching the same criterion. | 2020 |

... Continuation

| [408] | MedGAN: Image-to-image translation in the medical domain using a *CasNet* generator that progressively refines the image via encoder-decoder pairs. This feature-extracting discriminator evaluates the presence of desired modalities and non-adversarial losses that evaluate the generated image's style (e.g., structure and texture application). | 2020 |
|---|---|---|
| [409] | Mol-CycleGAN: Using CycleGAN [359] to enable the bidirectional translation of molecule embeddings between molecules with and without an optimized property. The GAN operates in the latent space of a junction tree VAE [127], simplifying training because the similarity measure used for this molecule representation is differentiable. | 2020 |
| [410] | SMOOTH-GAN: EHR generation conditioned on diagnosis codes (binary vector of present diseases) built upon the conditional GAN [353] with the Wasserstein GAN gradient penalty loss [285, 286]. | 2020 |
| [411] | WG$^2$GAN: Image-to-image translation to generate wounds from segmentation maps with a conditional GAN. | 2021 |
| [412] | Using a modified DCGAN [267] to generate lung CT images for COVID-19 classification. | 2021 |

### 2.13.2 Deep/Stacked GANs

Deep GANs contain multiple layers of generators and discriminators or share partial results (e.g., images at different resolutions and the respective discriminative feedback) at specific steps with each other. This allows the discriminator(s) to evaluate the final product, assess higher-order features, and provide more detailed feedback to the generator(s), often resulting in better results.

Denton et al. [413] introduce LAPGAN, a Laplacian pyramid framework having conditional CNNs trained with GANs at each layer to generate coarse-to-fine images. For this purpose, the generator upsamples the image from the previous stage to double the width and height and then applies a convolution based on the upscaled image and random noise. During the training process, a high-resolution image $I$ is downscaled, blurred, and upscaled again. This modified version $l$ is then used to either generate a real high-pass image $h = I - l$ or synthetic sample $G(z, l)$ that is given to the discriminator together with $l$ to determine whether it is real or generated. Log-likelihood evaluation with Parzen window estimates and human evaluation show that LAPGAN produces more realistic results than a standard GAN.

Zhang et al. [414] propose *StackGAN*, where the text-to-image task is divided into two stages. The first stage GAN draws the shape and color of the object described by the text. In stage two, a second GAN refines the sketch given the text description to produce photo-realistic images. StackGAN outperforms previous state-of-the-art conditional GAN models GAN-INT-CLS [356] and GAWWN [357] in terms of resolution ($256 \times 256$ pixels) and realism.

Short overview of other usages of deep/stacked GANs:

| Approach | Description | Year |
|---|---|---|
| [415] | Stacked GANs: A stack of symmetric encoder-decoder layers with layer-wise representation pair discriminators and a Q-Net that evaluates the diversity of the output of the generator/decoder at each layer. Each generator layer has its independent noise input and previous layer input, which allows conditioning on the output of the encoder. The resulting label-conditioned images are of higher quality than compared shallow GANs. | 2017 |
| [416] | A two-stage GAN approach for retina images. In the first stage GAN, vessel tree segmentation masks are generated with the help of a small human data set. In stage two, a conditional DCGAN [267] learns to create the corresponding retina fundus image for the vessel tree. | 2017 |
| [417] | TransGAN: A deep GAN consisting of a purely transformer-based architecture, excluding conventional convolutional layers. TransGAN consists of a memory-friendly transformer-based generator that progressively increases feature resolution and a multi-scale discriminator that captures both semantic contexts and low-level textures. The study includes a new grid self-attention module for high-resolution image generation and a unique training procedure to address the instability issues associated with TransGAN. | 2021 |

### 2.13.3 Bidirectional GANs

A Bidirectional GAN (BiGAN) [418], also introduced as Adversarially Learned Inference (ALI) [419], lets the discriminator evaluate data-representation pairs $(x, z)$ of either a generator $x = G(z)$ creating data from a representation or an encoder $z = E(x)$ inversely generating representations for data and decide, whether $x$ is real or fake (see Figure 35). The result is a model that learns meaningful data representations for detection or classification tasks despite the encoder

and generator being unable to communicate. It can also be used similarly to an autoencoder $\tilde{x} = G(E(x))$ for generative tasks, for example, image generation. [418]
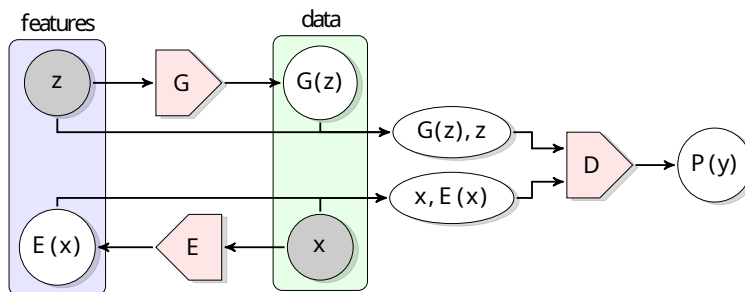


Figure 35: The BiGAN architecture. (Source: [418])

Donahue et al. [420] propose *BigBiGAN*, which combines the BigGAN [365] with the encoder of a BiGAN. They further modify the BigGAN discriminator to compute a score for $x$ and $z$, respectively, and a joint score. BigBiGAN achieves state-of-the-art results in unsupervised representation learning and unconditional image generation tasks, outperforming BiGAN and BigGAN, respectively.

### 2.13.4 Adversarial Autoencoders

The Adversarial Autoencoder (AAE) is a probabilistic autoencoder whose latent representations $\mathbf{z} \sim q(\mathbf{z})$ of the encoded data are forwarded to a discriminator who tries to distinguish the samples $q(\mathbf{z})$ from ones of a prior user-defined distribution $p(\mathbf{z})$ and provides an adversarial loss, as seen in Figure 36. By using the variational inference technique similar to a VAE (see Section 2.7.5), the autoencoder learns to map data to and generate meaningful samples from the whole space defined by the prior $p(\mathbf{z})$, which is configurable by the user and allows learning of powerful representations (see Figure 37 for example) for classification, disentangling of style and content, unsupervised clustering or data visualization. [421]
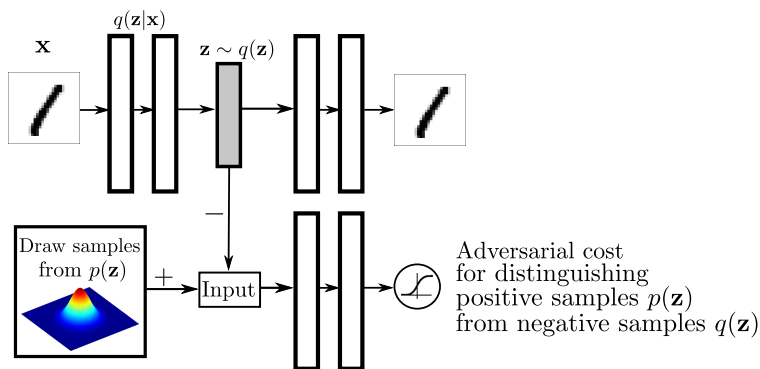


Figure 36: Architecture of an Adversarial Autoencoder. (Source: [421])

Makhzani et al. [421] demonstrate the generative capabilities of the AAE on MNIST, SVHN and TFD, achieving state-of-the-art log-likelihood on the test data compared to the standard GAN [238], Generative Moment Matching Network (GMMN) [422], DBN [65] and GSN [94] and allowing the combination of disentangled styles (writing style or font) and contents (numbers) on MNIST and TFD.

Tolstikhin et al. [423] propose the Wasserstein Autoencoder (WAE) as a generalization of the AAE, inspired by the Wasserstein GAN [285] adversary that uses the probability distribution discrepancy between real and generated data as the adversarial loss. Two approaches to compute the discrepancy between $q$ and $p$ are presented: The WAE-GAN incorporates an adversary that approximates the Jensen-Shannon (JS) divergence as loss and is trained together with the autoencoder, while the WAE-MMD used an adversary-free MMD-based loss, similar to GMMNs (see Section 2.14). WAE-GAN is less stable but generates higher quality samples than WAE-MMD and a normal VAE.
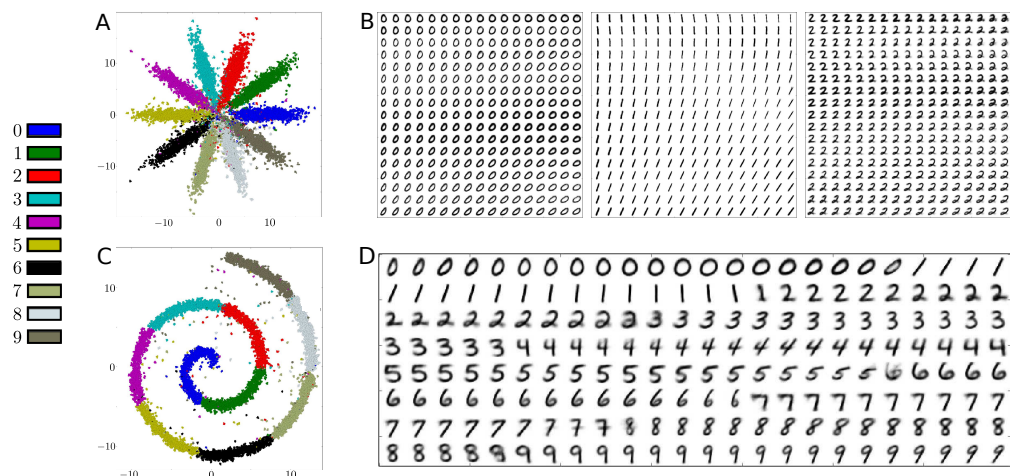
Figure 37: By forcing the latent space distribution of an Adversarial Autoencoder into specific regions, for example, to separate different numbers of MNIST, the quality and interpretability of representations can be greatly improved. (Source: [421])

Short overview of other usages of AAEs:

| Approach | Description | Year |
|---|---|---|
| [424] | An AAE is used to generate suitable molecule fingerprints for cancer treatment (trained on 6252 samples), and these probabilistic fingerprints are used to search for similar real molecules in the PubChem database with 72 million entries. | 2017 |
| [425] | druGAN: An improved AAE trained on larger data sets than Kadurin et al. [424] that demonstrates superior reconstruction results over a comparable VAE with the same sampling variability (i.e., coverage), which they identify is a trade-off necessary for generative autoencoders. | 2017 |
| [426] | End-to-end retinal image generation with an AAE generating the vessel network and a conditional GAN creating the corresponding fundus image. The GAN discriminator then evaluates the vessel-fundus pairs in terms of realism. | 2017 |
| [427] | An AAE with an additional classification step for biomedical time series generation. This allows the authors to generate labeled time series data using a semi-supervised approach. A dimensionality reduction is further introduced, transforming three-dimensional biomedical breathing data into one-dimensional time series and back. | 2021 |

## 2.14 Generative Moment Matching Networks

Generative Moment Matching Networks (GMMNs) use a neural network to learn deterministic mappings from an easy-to-sample data distribution to the real data distribution, similar to a GAN generator. The model starts with a top hidden layer $\mathbf{h} \in \mathbb{R}^H$ whose elements are usually independently sampled from a uniform distribution so that $p(\mathbf{h}) = \prod_{j=1}^{H} U(h_j)$. Then, $\mathbf{h}$ is passed through multiple neural network layers until the final output is returned as a data sample. [422]

The GMMN is trained, unlike a GAN generator, minimizing the MMD criterion instead of adversarial training with a discriminator. The MMD criterion is defined as

$$\mathcal{L}_{MMD^2} = \frac{1}{N^2} \sum_{i=1}^{N} \sum_{i'=1}^{N} k(x_i, x_{i'}) - \frac{2}{NM} \sum_{i=1}^{N} \sum_{j=1}^{M} k(x_i, y_j) + \frac{1}{M^2} \sum_{j=1}^{M} \sum_{j'=1}^{M} k(y_j, k_{j'}) \tag{30}$$

and computes, whether the generating distributions for two sets of samples $X = \{x_i\}_{i=1}^{N}$ and $Y = \{y_j\}_{j=1}^{M}$ are the same using a Gaussian kernel $k(x, x') = \exp(-\frac{1}{2\sigma}|x - x'|^2)$ with $\sigma$ being the bandwidth parameter. The gradient is differentiable; that is, the gradient can be backpropagated to update the parameters of the neural network. The larger the dimensionality of the data, the more training data is required to estimate the data distribution. [422]

Li et al. [422] use a GMMN with four intermediate non-linear ReLU layers and a logistic sigmoid output layer to model the MNIST and TFD data distribution. Further, they improve their generative model by training their GMMN on top of the latent representation of an autoencoder to reduce the data dimensionality and, consequently, the amount of training data needed.

Li et al. [428] propose *MMD GAN*, which replaces the Gaussian kernel of a GMMN with an adversarial kernel learning technique to reduce the amount of required training data and improve the model accuracy and efficiency. The Gaussian kernel is modified with injective functions $f_\phi$ with optimizable parameters $\phi$, resulting in the new kernel function $\tilde{k}(x, x') = \exp(-\|f_\phi(x) - f_\phi(x')\|^2)$. The GMMN model (generator) parameters $\theta$ and MMD (discriminator) kernel parameters $\phi$ are then adversarially optimized in a minimax game $\min_\theta \max_\phi \mathcal{L}_\phi(X, Y_\theta)$, like in a GAN, minimizing the generator loss for the worst possible discriminator result.

Short overview of other usages of GMMNs:

| Approach | Description | Year |
|---|---|---|
| [429] | MMD nets: Proposal of training/initializing generative neural networks with the inexpensive MMD statistic instead of a GAN discriminator at the same time as Li et al. [422]. | 2015 |
| [430] | Conditional GMMN: Feeding a sample drawn from a simple distribution and the conditional variables to the network that generates the target sample. A conditional MMD criterion is developed to learn the parameters. The model is used for conditional generation of MNIST and face images, combined with an autoencoder like Li et al. [422]. | 2016 |
| [431] | Speech synthesis by sampling speech parameters from a GMMN-learned distribution to induce variation in the generated speech and make it more natural. | 2017 |
| [432] | Introducing a repulsive loss function to address the limitations of existing MMD loss functions that may hinder learning fine details in data. The repulsive loss function enhances learning by emphasizing differences among real data samples. Additionally, the study proposes a bounded Gaussian kernel to stabilize MMD-GAN training. The methods are applied to unsupervised image generation tasks on datasets, showing significant improvements over the traditional MMD loss without additional computational costs. The paper also explores regularization techniques for MMD and the discriminator, contributing to the stability and effectiveness of the training process. | 2018 |
| [433] | Generation of high-dimensional load curves (cooling, heating, power) of integrated energy systems with a CNN generator and transformation of real and generated samples to latent space with an autoencoder where the distributions are compared using an MMD loss. | 2022 |

### 2.15 Plug & Play Generative Networks

Plug & Play Generative Networks (PPGNs) consist of a pre-trained generator network $G$ with latent variable space, usually obtained from a GAN, and a replaceable ("plug and play") pre-trained condition network $C$, which can be a classifier or image captioning network for example and "tells the generator what to draw". The class probabilities of the classifier are used to perform gradient ascent on the latent space of the generator, iteratively improving the generated results. $G$ and $C$ can even be trained on different data sets or domains, allowing for at least a limited form of domain transfer. [434, 435]

Nguyen et al. [434] introduce the deep generator network with activation maximization, short DGN-AM (see Figure 38), which is the origin of PPGNs. It aims to synthesize data from a pre-trained generator, which maximizes the activation of a specific neuron of a classifier. The architecture produces images of state-of-the-art quality but little diversity [435].

In another work [435], Nguyen et al. improve upon the lack of diversity of DGN-AM by learning a prior distribution for the latent variable of the generator using a DAE, which serves as an additional optimization criterion. They further define the generalized class of PPGNs and experiment with other building blocks like image captioning networks for the condition network. They achieve higher sample quality and diversity than DGN-AM.

### 2.16 Copulas

A copula $C$ models and decomposes the joint probability distribution of a continuous random vector $\mathbf{X}$ into a product of the marginal distributions and a representation of the marginal variables' dependence structure. [436]

The creation of a copula model relies heavily on estimation: First, the marginal cumulative distribution functions $F_i$ for individual random variables $\mathbf{X}_i$ are approximated with the help of training data and converted to a uniform distribution using probability integral transform. Then, the variables' dependencies and correlations can be estimated in various
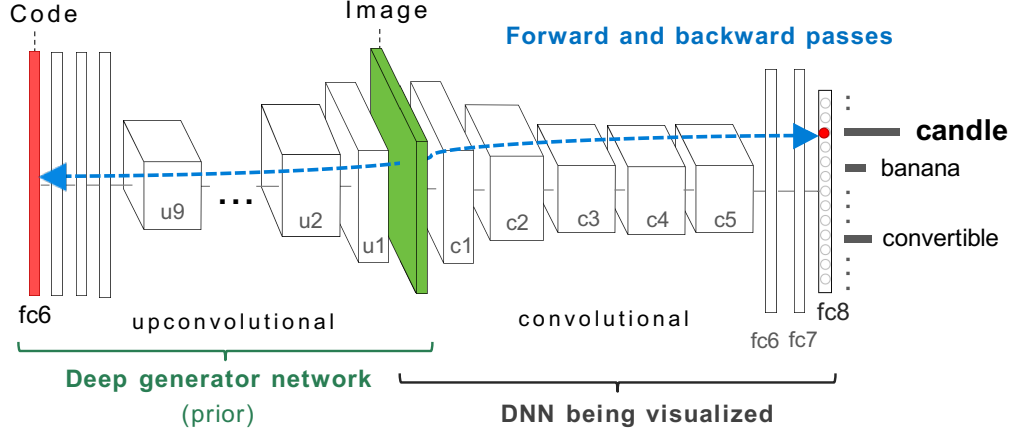
Figure 38: Example of the first PPGN, the DGN-AM. (Source: [434])

ways, for example, nonparametrically via kernel estimation [437, 436] or decomposition into trees of pair-copulas called *vines* [436].

Li et al. [438] introduce *DPCopula*, a technique to synthesize differentially private multi-dimensional data with copula functions efficiently. They propose two metrics, maximum likelihood estimation, and Kendall's $\tau$ correlation, to estimate the parameters of the Gaussian copula function. Both methods are analyzed in terms of privacy guarantees and computational complexity on census data sets, outperforming previous state-of-the-art approaches like Private Spatial Decomposition (PSD) and P-HP (lossy compression), especially on high-dimensional data sets.

Kulkarni et al. [437] use vine copulas to iteratively generate mobility trajectories, which are defined as a temporally ordered sequence $T = \langle (l_1, t_1), ..., (l_n, t_n) \rangle$ with locations $l_i$ (cell IDs) and timestamps $t_i$. They find that the copulas model observed statistical and semantic/geographic similarities particularly well at a fraction of the computational cost of neural network approaches while also incorporating long-range dependencies. Less accurate models like RNNs and GANs perform better in their privacy tests.

Tagaskova et al. [436] propose Vine Copula Autoencoder (VCAE), which utilize a CNN autoencoder based on DCGAN [267] to extract lower-dimensional representations from data and fit a nonparametric vine copula to learn the representations' distribution. By sampling from the trained copula, the decoder of the autoencoder can be used as a generative model. The VCAE is tested on three real-world image data sets in terms of MMD score and Classifier Two Sample Test (C2ST) accuracy: MNIST, SVHN, and CelebA. VCAE performs similar to DCGAN and outperforms the VAE.

### 2.17 Normalizing Flow Models

Normalizing flows are deterministic invertible transformations $f : \mathcal{E} \to \mathcal{Z}$ with parameters $\theta$ between a base distribution $\mathcal{E}$ (e.g., Gaussian distribution) and observational space $\mathcal{Z}$. The transformation can be used to calculate the exact density (probability) of an observation by using $f^{-1} : \mathcal{Z} \to \mathcal{E}$ or sample new observations by sampling from the simpler distribution $\epsilon \in \mathcal{E}$ and transforming it to observation space $z \in \mathcal{Z}$. [439]

Autoregressive Flows (AFs) are a variant of normalizing flows that models an observation $z \in \mathbb{R}^D$ as

$$p(z_d|z_{1:d-1}) = \mathcal{N}(z_d|\mu_d, (\alpha_d)^2), \text{ with } \mu_d = g_\mu(z_{1:d-1}; \theta), \alpha_d = g_\alpha(z_{1:d-1}; \theta), \quad (31)$$

where $g_\mu$ and $g_\alpha$ are unconstrained positive scalar functions, usually implemented as neural networks, that compute the mean and deviation for a normal distribution $\mathcal{N}$. The transformations are defined as

$$f_\theta(\epsilon_d) = z_d = \mu_d + \alpha_d \cdot \epsilon_d; \; f_\theta^{-1}(z_d) = \epsilon_d = \frac{z_d - \mu_d}{\alpha_d}. \quad (32)$$

To sample $z$ from an AF, first $\epsilon \in \mathbb{R}^D$ is sampled, then $z_1$ is computed using Equation 32. Finally, each subsequent $z_d$ can be computed using Equation 31. [439]

Shi et al. [439] propose *GraphAF* for graph-based molecule generation with parallel training, which uses AFs to generate graphs sequentially. At each step, the type of the next node (made continuous with a dequantization technique) is predicted before all edges to existing nodes are determined. The node and edge distribution parameters respectively $\mu_i^X$ and $\alpha_i^X$ and $\mu_i^A$ and $\alpha_i^A$ are computed with MLPs that are conditioned on the node embeddings $H_i$ that are obtained from a relational graph convolutional network [440]. The authors additionally employ *valency checking* [231] to reject invalid edges and implement a RL approach for goal-directed molecule generation.

Short overview of other usages of normalizing flow models:

| Approach | Description | Year |
|---|---|---|
| [441] | Non-linear independent components estimation (NICE): Learning of a non-linear deterministic and easily invertible transformation (deep neural network) $f$ that maps an input $x$ to a hidden representation $h = f(x)$ and back, so $x = f^{-1}(h)$. The training criterion is the exact log-likelihood. Since the transformations are deterministic, noise is injected at $h$ for generation. Results in image generation achieve high log-likelihood but often do not look realistic. | 2014 |
| [442] | Real-valued non-volume preserving (real NVP) transformations: Efficient invertible mapping of images to latent variables. | 2016 |
| [443] | Glow: Using an invertible $1 \times 1$ convolution to generate realistic and large images. | 2018 |
| [444] | Neural Ordinary Differential Equation (ODE): A continuous-time mapping $z(t)$ from latent variables to data defined by ODEs, also called continuous normal flow [445]. The model is used to accurately model and extrapolate time series or replace discrete hidden layers of a neural network. | 2018 |
| [445] | FFJORD: Combining the continuous normal flow procedure from [444] with an estimator of the log density for training instead of using maximum likelihood. This significantly reduces the computational cost and allows unrestricted architectures. The model outperforms previous flow-based methods on density estimation and also image generation. | 2018 |
| [446] | GraphNVP: The first flow-based invertible graph generation model that can handle fixed-size node type assignments and adjacency matrices together. The model maps the node feature and annotation matrices to latent representations, which can also be randomly sampled for generative purposes. They also train a simple linear regressor on the latent space for property-targeted molecule generation. | 2019 |
| [447] | Fourier Flows: This approach operates in the frequency domain, using discrete Fourier transformation to handle variable-length time-series with varying sampling rates and leveraging the more computationally efficient convolutions in the frequency domain. Fourier Flows applies data-dependent spectral filters to the frequency-transformed data, enabling efficient Jacobian determinants and inverse mapping computation. This method shows competitive performance compared to state-of-the-art models. | 2021 |

## 2.18 Reinforcement Learning

In Reinforcement Learning (RL), an agent starts in a state $s_0 \in \mathcal{S}$ within its environment and obtains an initial observation $\omega_0 \in \Omega$. The agent then has to decide on an action $a_t \in \mathcal{A}$ at each time step $t$. After performing the action, the agent receives a reward $r_t \in \mathcal{R}$, the state transitions to $s_{t+1} \in \mathcal{S}$ and the agent gets a new observation $\omega_{t+1} \in \Omega$, as seen in Figure 39. The goal of RL is to learn and optimize a policy $\pi$ so that the actions taken maximize the cumulative reward. Therefore, the agent uses a value function $V$ to predict the reward for an action. The policy can be *deterministic*, so it can be defined as $\pi(s) : \mathcal{S} \to \mathcal{A}$, or *stochastic*, assigning a probability $\pi(s, a) : \mathcal{S} \times \mathcal{A} \to [0, 1]$ to each action. A significant advantage of RL is that the agent does not need complete knowledge or control of the environment, which often makes it more computationally efficient than classic supervised and unsupervised ML methods. [448]
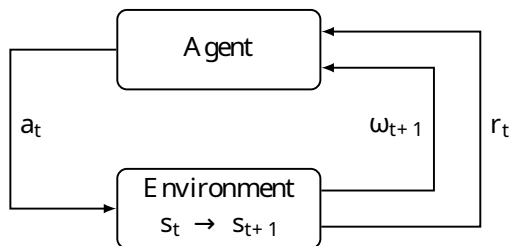


Figure 39: Agent-environment interaction in RL. (Source: [448])

Oh et al. [449] predict the next frames of Atari video games conditioned on previous frames and player actions. Input frames are encoded with a CNN (and optionally a RNN) to extract spatio-temporal features and then combined with a 1-hot encoded action variable in a transformation layer to obtain a high-level prediction of the next frame. A CNN takes this prediction and uses upsampling to generate a full-size frame. The model is trained on emulator recordings with the corresponding user inputs using stochastic gradient descent with backpropagation through time. The trained model can generate future frames for arbitrary input sequences. The evaluation shows realistic 100-step future frames for a variety of Atari games.

Jia et al. [450] train a painting agent's policy with RL to brush strokes step-by-step on a canvas guided by a reference image. The painting process of *Paintbot* consists of multiple steps:

1. A random stroke starting point $p$ on the canvas is selected.

2. Image patches centered around $p$ from the reference image and the canvas act as the observation $o$.

3. While the predicted reward $V_\pi(o)$ is positive, actions (strokes) are performed, consisting of continuous values for angle, length, color, and width, the canvas is updated, and $p$ and $o$ are both updated to the new position.

This process is repeated until a specified similarity threshold between the reference image and painting is reached. The training process consists of an additional loss function used to train the deep neural network for reward prediction $V$ (see Figure 40).
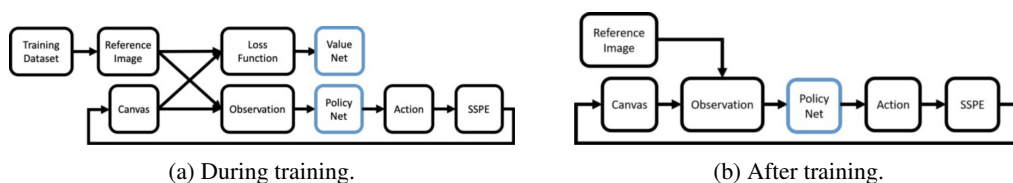


(a) During training.    (b) After training.

Figure 40: The Paintbot painting process during and after training. (Source: [450])

Krishna et al. [451] use a RL approach to synthesize full-resolution CT images. They create anatomically correct semantic masks and use their existing conditional GAN style transfer network [363] to fill the generated masked areas with correct textures. Semantic masks are represented as vectors via b-splines and principal component analysis for which the agent learns a policy. An image classifier CNN is used as the reward predictor and trained with human feedback on the agent's generated masks through an interface.

Short overview of other usages of RL:

| Approach | Description | Year |
|---|---|---|
| [452] | Dialogue text generation with two agents, which are encoder-decoder LSTMs. Agents are rewarded for displaying three conversational properties: Informativity, coherence, and ease of answering. | 2016 |
| [453] | Extending the idea from [452] for dialogue generation, but using an adversarial discriminator to distinguish between human-generated and synthetic dialogues and using the probabilities put out by the discriminator as the reward function. | 2017 |
| [454] | ORGAN: Combination of a SeqGAN [295] with RL where the GAN generator is trained with a tunable reward function consisting of the discriminator classification result, a repetition penalty, and domain-specific objective functions. The network is evaluated on SMILES [112] molecule and musical melody generation. | 2017 |
| [455] | Generation of SMILES [112] representations of molecules with specific properties using a prior RNN trained on a SMILES database and a user-defined scoring function as rewards for an RL agent RNN initialized by the prior RNN. | 2017 |
| [456] | Application of inverse RL to text generation, where a reward function is alternately learned on training data and an agent learns an optimal policy to maximize the total reward. In the implementation, a reward approximator MLP aims to maximize the log-likelihood of the training set samples, and the text generator LSTM is trained with a policy gradient [457] technique based on the reward and entropy regularization to encourage more diverse results. | 2018 |

Continuation...

... Continuation

| [458] | Goal-directed molecular graph generation with a graph convolutional policy network (GCPN): Molecule generation as a Markov decision process where the next graph state only depends on the previous one. At each step, a new subgraph, in this case, predefined as all single-node graphs of all allowed atom types, is connected to an existing node, or two existing nodes are connected by the GCPN. The GCPN is rewarded by a sum of domain-specific and adversarial rewards to ensure the molecule's utility, realism, and validity. | 2018 |
|---|---|---|
| [459] | Amadeus: Train a LSTM on a representation of multiple monophonic note streams that provide a polyphonic piece of music to simplify the learning process. Then RL is applied to select high-level LSTM configurations that produce the desired outputs instead of modifying weights or outputs. | 2019 |
| [460] | Addressing the compounding errors in sequential generation by combining contrastive imitation and an energy model. The model aims to capture both the step-wise transitions and the overall trajectory distribution, balancing the local and global properties of time-series data. | 2021 |

## 2.19 Diffusion Models

Diffusion models are Markov chains that iteratively add Gaussian noise to data $\mathbf{x}_0$ in a *forward process* over $T$ steps and also learn the *reverse process* that iteratively maps the noise input back to the data distribution (see Figure 41). [461]
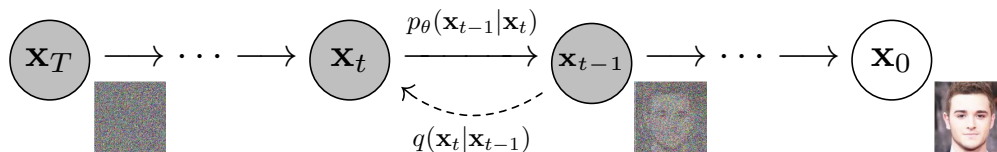


Figure 41: Diffusion forward and reverse processes. (Source: [461])

The forward process is defined as

$$q(\mathbf{x}_{1:T}|\mathbf{x}_0) = \prod_{t=1}^{T} q(\mathbf{x}_t|\mathbf{x}_{t-1}), q(\mathbf{x}_t|\mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1-\beta_t}\mathbf{x}_{t-1}, \beta_t \mathbf{I}), \tag{33}$$

where the variance schedule $\beta_1, ..., \beta_T$ can be constant or learned hyperparameters. The reverse process is defined as

$$p_\theta(\mathbf{x}_{0:T}) = p(\mathbf{x}_T) \prod_{t=1}^{T} p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t), p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \mu_\theta(\mathbf{x}_t, t), \mathbf{\Sigma}_\theta(\mathbf{x}_t, t)), \tag{34}$$

where $\mu_\theta(\mathbf{x}_t, t)$ and $\mathbf{\Sigma}_\theta(\mathbf{x}_t, t)$ are functions that provide the mean and covariance for the Gaussian and are defined using MLPs. [462, 461]

The model is trained by maximizing the variational lower bound on the NLL (like a VAE)

$$\mathbb{E}[-\log p_\theta(\mathbf{x}_0)] \leq \mathbb{E}_q[-\log \frac{p_\theta(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T}|\mathbf{x}_0)}] = L, \tag{35}$$

which is done by optimizing the aforementioned MLPs. [462, 461]

Sohl-Dickstein et al. [462] provide the first implementation of diffusion models and apply it to the generation and inpainting of images and binary sequences, performing worse than GANs, but better than DBNs, GSNs and CAEs in terms of log-likelihood.

Ho et al. [461] propose the Denoising Diffusion Probabilistic Model (DDPM), which simplifies the training process of [462] by replacing $\mathbf{\Sigma}_\theta(\mathbf{x}_t, t)$ with untrained time-dependent constants $\sigma_t^2 \mathbf{I}$ in the reverse process and training an estimator $\epsilon_\theta(\mathbf{x}_t, t)$ to predict the noise $\epsilon$ added to $\mathbf{x}_t$ instead of predicting the mean $\mu_\theta(\mathbf{x}_t, t)$. The variational lower bound is simplified to optimize the difference between actual and predicted error. For the estimator $\epsilon_\theta$, PixelCNN++ [246] is adopted with self-attention at the small feature maps.

Nichol et al. [463] introduce the improved DDPM for better log-likelihood results. It reintroduces $\Sigma_\theta(\mathbf{x}_t, t)$ as a trainable model and replaces the linear noise schedule for $\beta_t$ with a cosine schedule, which spreads the noise addition in the forward process more evenly. They also reduce the number of diffusion steps to improve sampling speed with very little quality loss by scaling the schedule parameters, and increasing the model size also increases performance.

Dhariwal et al. [464] propose the ablated diffusion model (ADM) with classifier guidance (ADM-G), which improves upon [463] by using more attention at different scales, class conditioning, a deeper model architecture, and a classifier to guide the generation process more precisely.

Ramesh et al. [465] add text embeddings produced by a decoder-only transformer to the existing timestep embedding of a diffusion model to produce text-conditional images. Images can be manipulated based on text by reconstructing images from $\mathbf{x}_T$ with the diffusion model conditioned on new/interpolated text embeddings.

## 2.20  Virtual Environments

Virtual environments are computer-simulated "graphic and real-like models of real-life objects" [5] that are simple to annotate by nature since object locations, classes, and other properties are apparent in the simulation software. More realistic depictions of virtual objects have become possible because of the increasing computational processing power of Graphics Processing Units (GPUs) in recent years. They allowed virtual environments to be adopted for various tasks, such as autonomous driving or gesture recognition. Based on [5], we derive three categories for virtual environment usages:

**Graphic Models**  In this simple scenario, single 3D Computer Aided Design (CAD) models or compositions are used to alleviate ML model training for, e.g., gesture recognition and object recognition from different viewpoints or building 3D models from images.

**Virtual Worlds**  These models are computer-simulated and emulate a complex real-world environment. They are populated with many objects that can sometimes move or interact with the world or each other. They are especially popular for generating annotated training data for autonomous driving systems.

**Interactive Environments**  These interactive virtual worlds, usually video games or simulators with an inherent goal and well-defined rules, are especially suitable for the training and benchmarking/competition of RL agents because they allow user inputs and directly present a result. The advantages are their high availability for various scenarios (e.g., driving simulations, open-world, and strategy games) and their often simple adaptability to research tasks through mod support and editor software.

For optimal results with rendered synthetic data for computer vision tasks, Mayer et al. [466] present several findings: Multistage training on different data sets works better than mixing or training on one data set alone, complex and more realistic lighting does not necessarily help, and incorporating flaws of a real camera during training improves model performance.

### 2.20.1  Graphic Models

Butler et al. [467] introduce *MPI-Sintel*, an optimized version of the open-source 3D animated short film Sintel [468] for optical flow evaluation. As ground truth, motion vectors for each pixel are computed using a modified version of Blender's [469] motion blur pipeline. The advantages over other data sets for optical flow evaluation are the long sequences and motions and the ability to render the film with various effects like motion blur, specular reflections, and atmospheric effects enabled or disabled.

Handa et al. [470] create the first RGB-D (RGB plus depth) image collection with camera trajectory and full 3D scene ground truth attached to each frame. They provide raytraced renderings of two rooms, the office room and the living room, via the POVRay [471] raytracing software. They also separately apply artificial noise to the RGB and depth values to make the images more realistic. The data set is then used to benchmark algorithms for visual odometry, 3D reconstruction, and Simultaneous Localization and Mapping (SLAM).

Su et al. [472] utilize 3D models rendered on top of real background images to train CNNs for viewpoint estimation. The models are rotated and inserted at different positions in the picture for that purpose. Their "render for CNN" approach achieves state-of-the-art performance on a benchmark data set at negligible human cost using existing 3D model repositories.

Peng et al. [473] train deep CNN object detectors with synthetic images rendered from non-photorealistic 3D CAD models that are freely available on the Internet to detect novel object categories not available in the real training data. They evaluate the models on real images, showing better detection performance than models trained on real data from a different domain.

Handa et al. [474] introduce a framework to randomly generate realistic and automatically annotated 3D indoor environments using 3D objects from public databases. The *SceneNet* uses a hierarchical scene generator that learns relationships (co-occurrence frequency) between objects from prior indoor scene data sets. In another work [475], they demonstrate their model's utility by using its renderings from random perspectives with added noise to train a deep model for depth-based semantic per-pixel segmentation.

Short overview of other usages of graphical models:

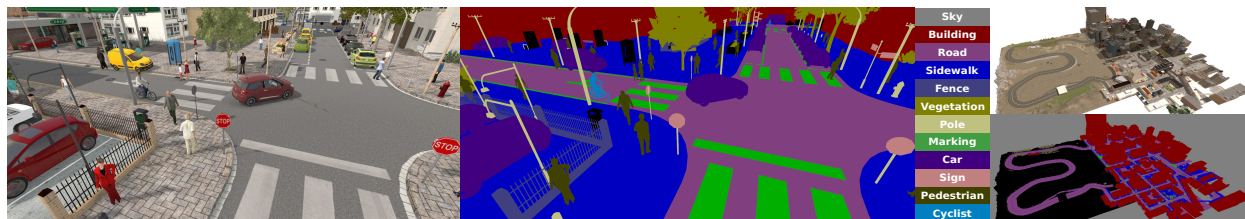| Approach | Description | Year |
|---|---|---|
| [476] | Creation of the Tsukuba CG Stereo data set, which is a collection of photo-realistic renderings of the *head and lamp* scene under varying conditions and camera perspectives for stereo matching task training and evaluation. | 2012 |
| [477] | Kinematic hand model with random initial global rotation is used to render hand gesture sequences for human-computer interfaces. | 2014 |
| [478] | Training of object detectors with 3D model renderings combined with domain adaptation using discriminative decorrelation performs comparably to models trained on real data (ImageNet). | 2014 |
| [479] | (Re-)Synthesis of natural images from different viewpoints aided by structural information from 3D models of the same object class aligned to an image. | 2014 |
| [480] | Simultaneous class, pose, and location prediction of possible objects using a deep CNN trained on RGB-D renderings of randomly generated 3D rooms built with intersection detection and plausibility checks (e.g., sofas are usually near walls). | 2015 |
| [481] | A collection of *Flying Chairs* image sequences combining natural background images with renderings of 3D chair models is used to train the *FlowNet* CNN for optical flow estimation. | 2015 |
| [482] | Using a 3D morphable face model to create synthetic training data for face recognition systems reduces the amount of real data needed significantly and improves performance. | 2018 |
| [483] | Rendering video sequences of neurosurgical instrument movements from a microscope's perspective in Blender [469] for optical flow estimation benchmarks in this domain. | 2021 |
| [484] | Using Blender [469] to render steel pieces with defects (e.g., cracks in the surface texture) and masks for a steel defect detection task. | 2021 |

### 2.20.2 Virtual Worlds



Figure 42: The SYNTHIA data set: Image rendered from the virtual world (left), ground truth segmentation map (middle), and city overview (right). (Source: [485])

Haltakov et al. [486] propose a framework built on top of the open-source driving simulator *VDrift* [487]. The modified software allows them to render realistic synthetic images with pixel-wise object annotations, depth, and optical flow maps (movements) in various scenarios, perspectives, and driving styles. The framework is then applied to create a large image training set for a multi-class image segmentation task.

Richter et al. [488] utilize the *detouring* technique to generate semantic label maps for images from closed-source modern computer games. They evaluate the communication between the game and graphics hardware using a graphics API wrapper library, hashing rendering resources like textures or geometry and creating persistent object signatures to which labels for urban scene understanding are applied. The game training data, obtained from *Grand Theft Auto V*, boosts the accuracy of segmentation models and reduces the need for expensive labeled real-world images.

Johnson-Robertson et al. [489] use open-source plugins and GPU buffer data to extract annotated and realistic images from *Grand Theft Auto V* for vehicle detection tasks. They achieve state-of-the-art performance on real data with a model trained only on simulated images.

Ros et al. [485] introduce the *SYNTHIA* data set (see Figure 42), which provides a semantically segmented collection of images of urban scenes obtained from the Unity game engine [490]. Since SYNTHIA is intended for autonomous

driving tasks, it not only contains images from various perspectives in the virtual world but also four video sequences of a virtual car driving through the urban simulated landscape, one for each season. The virtual car consists of two multi-cameras 0.8 meters apart, each consisting of four monocular cameras with depth sensors, a common center, and 90-degree rotation between them. In combination with real data, SYNTHIA can significantly improve segmentation model accuracy.

Short overview of other usages of virtual worlds:

| Approach | Description | Year |
|---|---|---|
| [491] | Assessment of the effectiveness of models trained with computer games on real-world data in image segmentation and depth estimation tasks. The results show similar or better performance than models trained on real data. | 2016 |
| [492] | Virtual KITTI: A labeled video data set containing computer-rendered sequences of driving through realistic virtual worlds obtained from the Unity [490] game engine. The virtual worlds are created with positional information from the original KITTI [493] data set and human optimization. | 2016 |
| [494] | Visual perception benchmark (VIPER): More than 250,000 high-resolution video frames with annotations, for instance, segmentation, optical flow, object detection, tracking, visual odometry, and object-level 3D scene layout tasks obtained by moving through the world of the video game *Grand Theft Auto V*. | 2017 |
| [495] | Meta-Sim: Using probabilistic scene grammars to create and render valid virtual environments. Performance improvements of the method are demonstrated by training a task network on the synthetic data and comparing it to a model trained on real data. | 2019 |

### 2.20.3  Interactive Environments

Bellemare et al. [496] build the Arcade Learning Environment (ALE) on top of an Atari 2600 emulator to evaluate AI techniques like RL and planning algorithms on arbitrary Atari games, which are usually split into a training and testing set. ALE provides an interface to the game controls, screen information, RAM, and registers for an AI agent to control or read. The reward function for an agent is defined per game based on the score difference between frames.

Kempka et al. [497] build upon the idea of [496], where 2D Atari 2600 games are used as an evaluation platform for RL agents, and propose *VizDoom*, a more realistic 3D game based on the first-person shooter Doom, as a research platform for visual reinforcement learning. They train competent bots using convolutional deep neural networks and RL on various tasks and scenarios.

Sadeghi et al. [498] make deep reinforcement learning applicable to safety-critical domains such as autonomous flights by training in virtual environments built entirely with CAD models. The CAD$^2$RL method trains a RL agent (deep CNN) on RGB images of a monocular camera mounted to a drone in the virtual environment to output velocity commands that avoid collisions. The authors find that by using highly randomized rendering settings, the agent's policy can be trained to generalize well to real-world applications, which they demonstrate by letting the trained agent fly a real drone through indoor environments.

Vinyals et al. [499] introduce the *StarCraft II Learning Environment* (SC2LE), which combines a Python-based interface for the game engine with specifications for possible observations, actions, and rewards. As a complex multi-agent problem with incomplete information, long-term strategies, and large action space, StarCraft provides a difficult class of problems to evaluate RL models on. The authors test agents on various mini-games, resulting in agent behavior similar to a novice player and on the main game, where the agents cannot progress noticeably.

Short overview of other usages of interactive environments:

| Approach | Description | Year |
|---|---|---|
| [500] | *DeepMind Lab*: First-person 3D game framework based on the *Quake III* game engine for easy task and AI design. | 2016 |
| [501] | *TorchCraft*: Providing an interface between the *Torch* ML framework and real-time strategy game "StarCraft: Brood War". | 2016 |
| [502] | Project Malmo: An AI experimentation platform for complex navigation, survival, collaboration, and problem-solving tasks built on top of Minecraft, a popular game mimicking the real world as a collection of blocks and friendly and hostile entities (e.g., animals, zombies). | 2016 |

# 3 Classification of Generative Models

In this section, we classify the models we presented in Section 2 by criteria presented in Section 3.1. In Section 3.2, we perform a trend analysis before finally presenting a guideline for model selection in Section 3.13.

## 3.1 Criteria for Classification

For each model presented in this work, we collect certain information we find to be comparable and useful for our trend analysis and guidelines:

**Metadata** We collect name, release year, model (sub-)category (according to Section 2) and citations (according to Google scholar). Optional entries are predecessors (other models that are the foundation for this model) and combinations, which describe the categories a proposed model combines (e.g., a GAN can combine a CNN generator with a RNN discriminator).

**Data Structure** Determines if the size/amount of samples of the data generated by the model is limited, for example, an image with a static resolution or a fixed-length video, or (theoretically) infinite, which is often required for processing arbitrarily-sized sequences.

**Data Type** Type of the data generated by the model, for example, natural language text, time series, music, or images.

**Sampling Requirements** The data generation process of a model can be unconditional or conditional, which means an input is required based on which synthetic data is generated, or both (see Figure 43). If an input is accepted, we also capture the types of input that the model accepts.

**Sampling Process** Describes whether a sample is generated in "one go" or iteratively refined by the model.

**Training Process** Describes how the generative model is trained and is described by two aspects, inspired by [503]:

> **Loss Type** The loss function(s) that is/are used to optimize the model.
> **Optimization** Additional penalization of the model or modification of the loss function to improve results.

**Data Sets** The data sets used to train and evaluate the models.

**Model Performance** Comparing different models is complicated because no commonly used performance measure exists, and many proposed measures only work for specific data types or domains (e.g., music can not be evaluated in the same way as natural language). To work around this issue, we collect the *performance predecessors*, that is, the list of outperformed models, from the evaluation section of the respective paper, if available.

**Privacy** Shows whether the generated data is considered differentially private or private by another criterion defined by the respective authors.



this small bird has a pink breast and crown, and black primaries and secondaries.

this magnificent fellow is almost all black with a red crest, and white cheek patch.

(a) Image-to-image translation (Source: [335])     (b) Text-to-image translation (Source: [3])

Figure 43: Examples of GANs being applied to translative tasks.

## 3.2 Data Evaluation and Trend Analysis

This section presents the findings from our survey of existing SDG literature. In the following subsections that build on one another, we investigate the criteria proposed in Section 3.1 before concluding.

## 3.3 Metadata

First, we look at the metadata from our 417 models. In Figure 44, we show the number of papers grouped by model category we evaluated for each year. We focus primarily on models proposed in the last ten years, so the data starting in
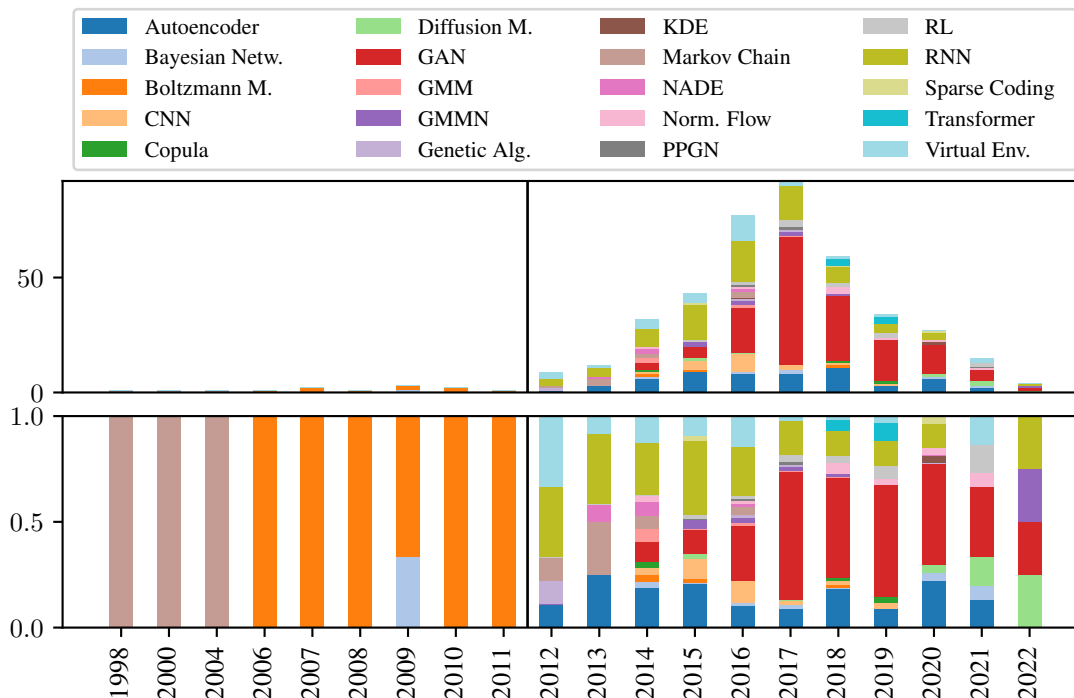
Figure 44: Total (top) and normalized (bottom) amount of models we present in our work grouped by year and model category.

2012 is primarily relevant. GANs, RNNs, autoencoders (especially VAEs), virtual environments, and CNNs experience high usage throughout the years in the literature we evaluated, with GANs quickly surpassing the other approaches since their first proposal in 2014 [238]. The usage of Markov chain models and Boltzmann machines declined over the years, while RL and diffusion models slightly gained popularity.

Next, in Figure 45, we evaluate the amount of Google Scholar citations per model category. Compared to other models covered by our work, GANs have received the highest amount of citations over the last ten years. In addition, RNNs, CNNs, and autoencoders often receive comparable attention.

In Figure 46, we evaluate which architectures and concepts our models borrow from other model types. The models often used as submodels according to Figure 46a are CNNs, followed by RNNs and autoencoders. In Figure 46b, we also show that some apparent assumptions are confirmed:

- GANs and the very similar GMMNs often use autoencoders, especially their decoders, as generative networks.
- 70% of our RL models use RNNs, which are suitable for guidance by reward functions due to their sequential data generation process. RL is also applied to GAN and CNN training.
- Autoencoders often incorporate RNNs and CNNs as their encoders and decoders.
- Diffusion models, GANs and normalizing flow models heavily utilize CNNs. This is because many GAN models are based on the CNN-based DCGAN [267], and diffusion models are also mostly related to the CNN-based DDPM [461], which is illustrated in Figure 47.

Finally, we discuss the *predecessors* class of our metadata. In Figure 47, we build a graph of all predecessor connections for notable model categories. We observe that especially GANs and RNNs are strongly interconnected. For autoencoders, we only found strong relations between VAE models. Four of the five diffusion models we presented significantly depend on each other, while the fifth model is extended for text-conditional tasks and, therefore, considerably changed.

## 3.4 Data Structure

In Figure 48, we discuss the data modeling capabilities of different model types. As expected, sequence-based models such as Markov chains, RNNs, and transformers are mainly applied to generate arbitrarily-sized data. Autoencoders and
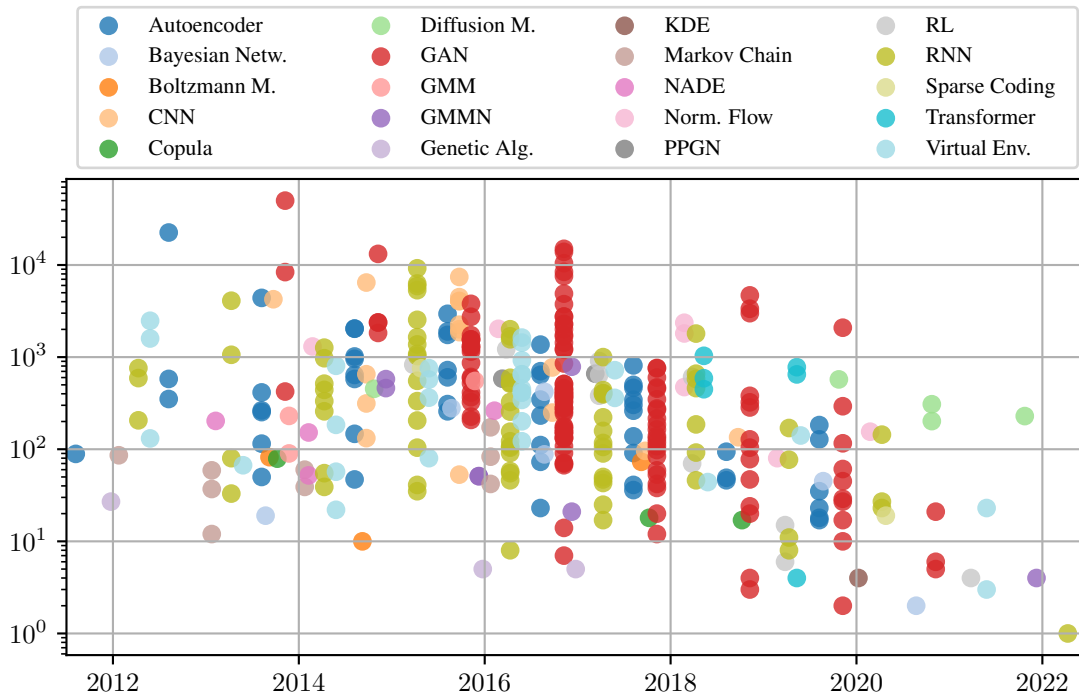
Figure 45: Google Scholar citations for different model categories over the last ten years on a logarithmic scale. For better readability, we apply a $[-0.4, 0.4]$ offset to the points' year value in the order of their legend appearance (top to bottom, left to right). The numbers were last obtained Date.

GANs are based initially on neural network architectures with fixed input and output sizes but can process sequences by adopting RNNs. These correlations are easy to observe by comparing the RNN connections in Figure 46b to our findings in this section. RL, which also heavily relies on RNNs, works well for unlimited data lengths as well, also due to the reward function being a powerful tool to learn and improve during the generation process. Modern Boltzmann machines like the TRBM are also strongly related to RNNs, making them noticeable in our data structure evaluation.

### 3.5 Data Types

We identify several data types in our research, of which multiple are often used by a single model. We categorize them as follows:

**Audio** Represents raw waveform audio. Due to the high sampling rate required to produce natural-sounding audio, this data type usually requires complex models and thorough training to produce good results. We further differentiate between music and speech generation due to their difference in complexity (e.g., note and instrument vs. text, speaker, and character transitions).

**Image** Describes classic two-dimensional bitmaps, usually with RGB or grayscale pixel values to which we refer as natural images. Binary images are simplified versions where pixels can either be on or off (black and white). Segmentation masks describe images where pixels are of a particular class instead of color. Images with more information encapsulate other pixel values besides color, for example, depth in RGB-D images or HSIs.

**Text** This class describes a sequence of characters. We differentiate between natural language, as spoken by humans, and text representations that encode other data types, for example, SMILES strings [112] for molecules.

**Time Series** Sequences of one (univariate) or more (multivariate) variable values that have to be determined for each step. The more potentially interdependent variables have to be specified, the more complex the task becomes. We additionally define symbolic music as an additional category because of the considerable interest in the topic and the complex constraints provided by music theory that must be considered. Depending on the task definition, symbolic music can be regarded as a univariate or multivariate time series.

(a) Total amount bottom models borrowed from upper ones.



(b) Normalized (in %, rounded down) by the number of models of bottom model categories that use the upper models.
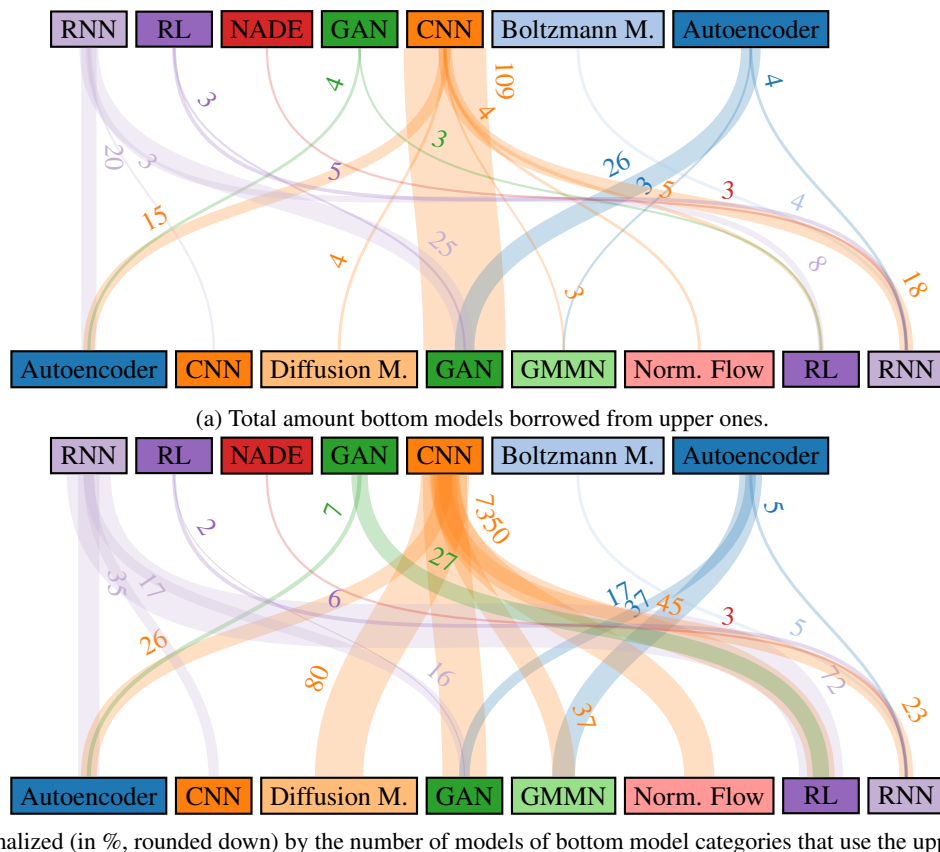
Figure 46: Weighted dependency graph for model compositions. The bottom models rely on submodels (e.g., encoders and decoders) from the upper categories. The edges' size and labels denote the amount or percentage of models that use the respective submodel. We omitted edges with less than three total usages for readability.

**Graphs and Molecules** Graphs are collections of nodes with connections (edges) between them. Nodes and edges can have additional properties and values assigned to them. They are usually represented by adjacency matrices or time series of node and edge creations. Many presented works consider molecules as a subset of graphs, where atoms are the nodes connected in a specific way, and they are often represented as SMILES [112] text representations in addition to the representation above types.

**Tabular Data** We consider tabular data to be tuples, sequences of tuples, or matrices containing categorical and numerical values. Besides generating tables, it is often used to provide additional information to another data type. Virtual environments use it to provide, for example, emulator or game information.

**Video** A sequence of images as defined above.

Figure 49 shows that natural image generation is the data type of most overall importance. RNNs are predominantly applied to sequential arbitrary-length domains like natural language text and symbolic music. Due to their flexible architecture, autoencoders and GANs are applied to almost all data types.

Figure 50 illustrates that some model types are limited to specific data types: Boltzmann machines are exclusively applied to time series, video, and image data. CNNs mainly focus on images and other high-dimensional formats like waveform audio. Virtual environments are exclusively applied to visual data and provide additional information. RL puts itself forward for domains suited for sequential generation.

## 3.6 Sampling Requirements

In Figure 51, we compare the models' reliance on additional information. Genetic algorithms, PPGNs, sparse coding models, and CNNs that we covered can always be conditioned on additional information to guide the generation process. Boltzmann machines, RNNs, transformers, RL, Markov chains, diffusion models, GANs, and autoencoders

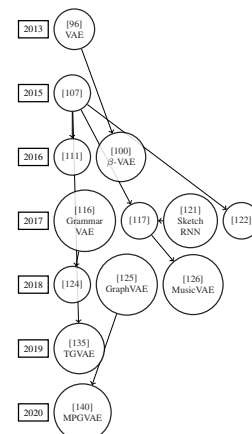(a) GANs.



(b) RNNs.



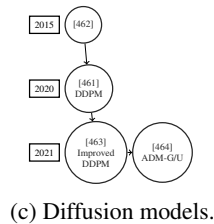(c) Diffusion models.



(d) Autoencoders.

Figure 47: Inheritance graphs of the model types for which we acquired a significant amount of data in the *predecessors* section of the metadata. We omit models without documented edges, and some years are spread across multiple rows to accommodate width limitations.

Figure 48: Fraction of data structures different types of models put out.



Figure 49: Fraction of data type usages per model category.

are also very flexible and often support conditional generation. The other model types are mainly used for unconditional generation; they only require random noise as input or during the generation process.

Figure 50: Heatmap showcasing the data output provided by different types of models.



Figure 51: Fraction of models conditioned on input data per model category.

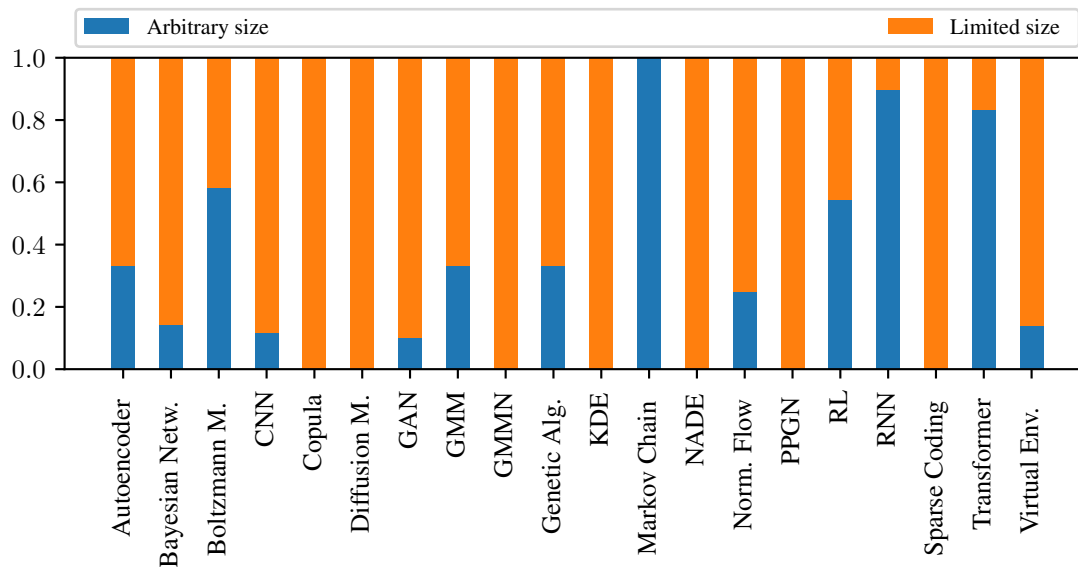In Figure 52, we provide a more detailed overview of the reliance of specific model types on certain conditioning input types. We extend the data types proposed in Section 3.5 with three new entries:

**Model Constraints** Constraints imposed on the model by the user, for example, static generation rules from music theory or positional constraints that require a specific step in a sequence to have a particular value.

**Class Labels** A one-hot vector or embedding that specifies certain aspects the generated data should have. This could be hair or skin color for human face image generation.

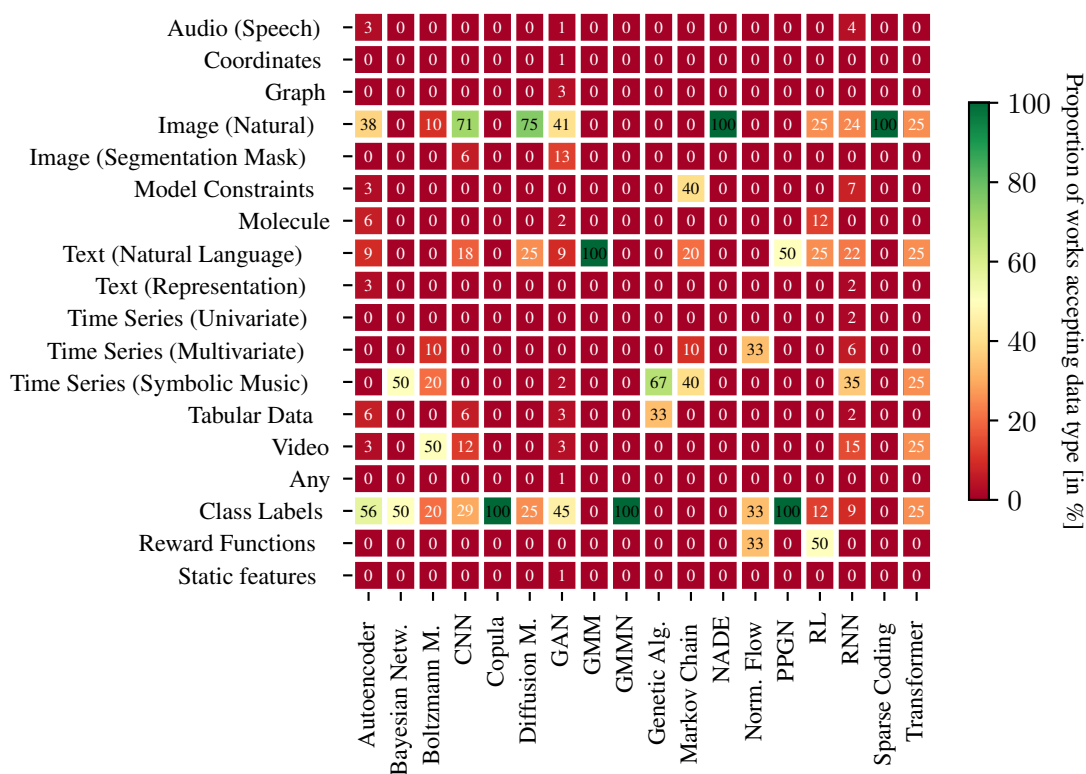| | Autoencoder | Bayesian Netw. | Boltzmann M. | CNN | Copula | Diffusion M. | GAN | GMM | GMMN | Genetic Alg. | Markov Chain | NADE | Norm. Flow | PPGN | RL | RNN | Sparse Coding | Transformer |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Audio (Speech) | 3 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 4 | 0 | 0 |
| Coordinates | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| Graph | 0 | 0 | 0 | 0 | 0 | 0 | 3 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| Image (Natural) | 38 | 0 | 10 | 71 | 0 | 75 | 41 | 0 | 0 | 0 | 0 | 100 | 0 | 0 | 25 | 24 | 100 | 25 |
| Image (Segmentation Mask) | 0 | 0 | 0 | 6 | 0 | 0 | 13 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| Model Constraints | 3 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 40 | 0 | 0 | 0 | 0 | 7 | 0 | 0 |
| Molecule | 6 | 0 | 0 | 0 | 0 | 0 | 2 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 12 | 0 | 0 |
| Text (Natural Language) | 9 | 0 | 0 | 18 | 0 | 25 | 9 | 0 | 100 | 0 | 20 | 0 | 0 | 50 | 25 | 22 | 0 | 25 |
| Text (Representation) | 3 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 2 | 0 | 0 |
| Time Series (Univariate) | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 2 | 0 | 0 |
| Time Series (Multivariate) | 0 | 0 | 10 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 10 | 0 | 33 | 0 | 0 | 6 | 0 | 0 |
| Time Series (Symbolic Music) | 0 | 50 | 20 | 0 | 0 | 0 | 2 | 0 | 0 | 67 | 40 | 0 | 0 | 0 | 0 | 35 | 0 | 25 |
| Tabular Data | 6 | 0 | 0 | 6 | 0 | 0 | 3 | 0 | 0 | 33 | 0 | 0 | 0 | 0 | 0 | 2 | 0 | 0 |
| Video | 3 | 0 | 50 | 12 | 0 | 0 | 3 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 15 | 0 | 25 |
| Any | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| Class Labels | 56 | 50 | 20 | 29 | 100 | 25 | 45 | 0 | 100 | 0 | 0 | 33 | 0 | 100 | 12 | 9 | 0 | 25 |
| Reward Functions | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 33 | 0 | 50 | 0 | 0 | 0 |
| Static features | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

*Proportion of works accepting data type [in %]*

Figure 52: Heatmap showcasing the data types accepted by different types of models.

**Reward Functions** User-defined functions that provide a specific value to be optimized by the model in addition to its default targets. For example, a generated molecule should have a particular chemical property.

The most often used sampling requirements across all model types are images and class labels, followed by natural language text and music. Due to their simple structure, Markov chains are suitable for applying model constraints. In contrast, RL models are predestined for using reward functions due to their flexible policy learning architecture. In the presented literature, we find segmentation masks to be exclusively used by CNNs and GANs, which also often use CNNs as generators (see Section 3.3).

## 3.7 Sampling Process

We found two types of sampling processes to be relevant: Determining the values of a sample or sequence of samples "in one go" (one shot) and iteratively refining the sample values a specific number of steps or until a criterion is reached. We present our findings in Figure 53, where we show that most models use one-shot sampling. PPGNs, genetic algorithms and diffusion models always iteratively sample data because of their architecture. Other models like autoencoders, CNNs, Markov chains (especially HMMs), NADEs, RL, and RNNs are also suitable to be applied multiple times or as deep architectures consisting of multiple submodels to the data to refine the results, but less often used in that way.

## 3.8 Training Process

In Figure 54, we provide a combined overview of different losses and optimization techniques used to train and improve generative models:

**MMD Loss/Data Distribution Matching** Comparison of overall data statistics between sets of real and generated samples. GMMNs use the MMD metric as their main training objective. Many GANs use the Wasserstein distance in addition to the adversarial loss to prevent mode collapse.

**Adversarial Loss/Classification Error/Auxiliary Classifier** Describes the competition of a generative model against a classifier model that judges whether its input is real or fake (adversarial loss, the foundation of GANs) or
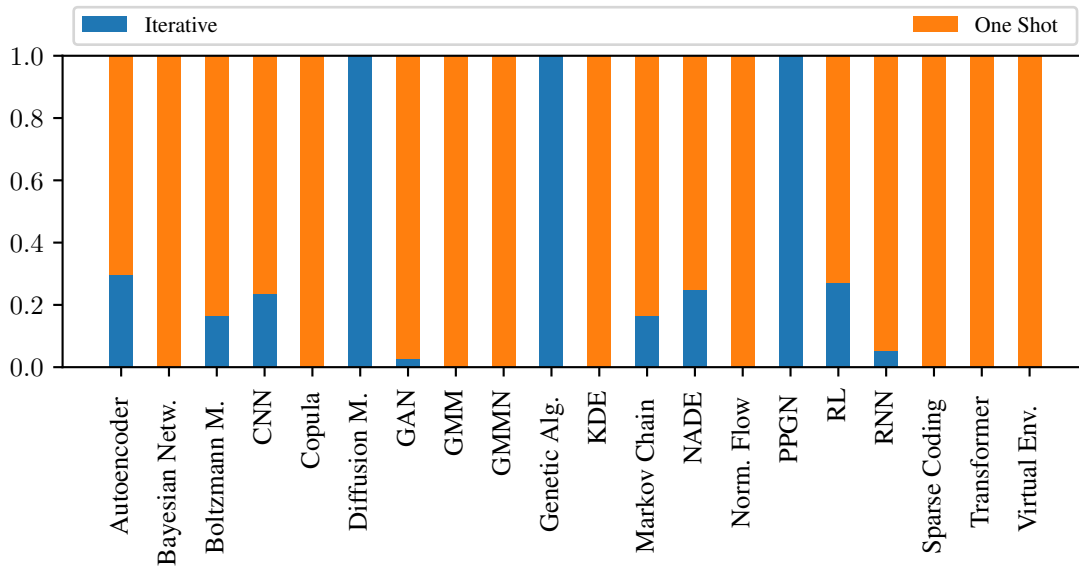
Figure 53: Fraction of models utilizing a specific sampling process.

the probabilities that the input belongs to a specific class. Classifiers are often used besides other loss types (auxiliary) or in PPGNs as the primary training objective, which can be considered a more powerful GAN discriminator.

**Data Likelihood** A simplistic evaluation approach of the model performance is often used by older models, letting the model assign a generation probability to real test data.

**Reconstruction Error/Cycle Consistency Loss** Reconstruction error is the foundation of autoencoders, Boltzmann machines, and the more general sparse coding model that takes data as input, converts it to a small representation, and then aims to reconstruct the data as accurately as possible. The accuracy is measured by the error metric, which is often a mean absolute or mean squared error. Cycle consistency loss, introduced by CycleGAN [359], allows a GAN to train based on reconstruction loss by training conditional generators and discriminators for both directions in an unpaired image translation task.

**Reward/Score Function and Content Feature Evaluation** Evaluation of specific aspects of the generated data for model training and guidance. This is especially relevant for policy-learning models like RL.

**Rule-Based Loss** Hard constraints imposed on the generated data and model, implemented by humans. It is used to force RNNs to comply with the basic rules of music theory.

**Limitations and Adaptation** The training process of many models is significantly modified. Model training with gradient descent is often improved and accelerated by using a dynamically adapted learning rate or restricting the gradient itself (e.g., gradient clipping, normalization, penalization). A model's network weights or parameters can also be heavily restricted by normalization, freezing, decay, clipping, or connection of some of them (i.e., weight sharing). Models that work with latent codes representing data (mostly autoencoders) also often impose distribution constraints on them, usually to simplify the sampling of new data.

**Ranking of Samples** A ranking of fake samples among real ones provides more detailed feedback to the generative model.

**Dropout** Set a random fraction of neurons of a neural network to zero at each training step to prevent neurons from learning the same features.

**Early Stopping** Stop optimization on the training data when model performance stops increasing on the evaluation data to prevent overfitting.

**Mode Regularization** Offering an incentive to the model or penalizing it for its coverage of data classes or the data distribution. Makes models more resilient to imbalanced data sets or mode collapse in the case of GANs.

**Noise Injection** Adding noise to the data or at specific layers in a neural network. Forces the model to generalize more.
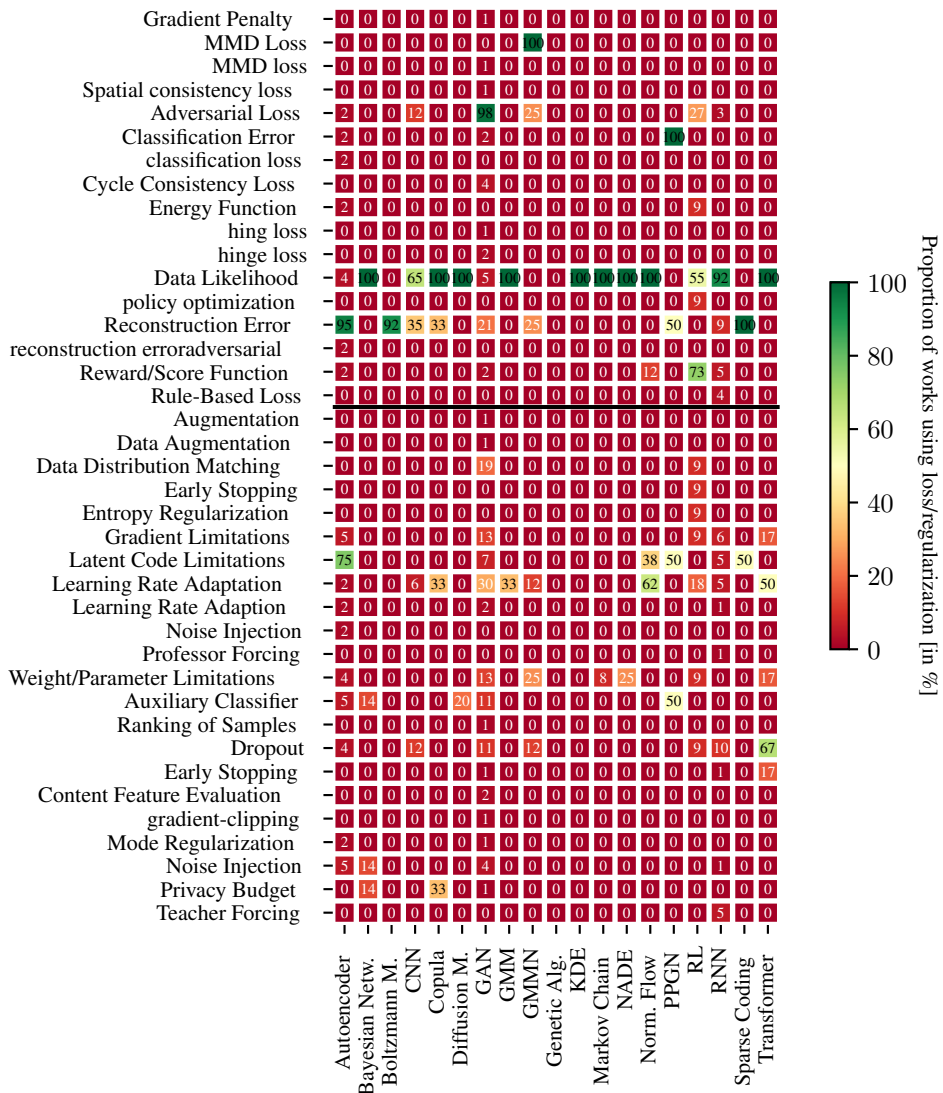
Figure 54: Loss (above) and additional optimization (below black line) heatmap for different model types.

**Privacy Budget**  Used to prevent differentially-private models from disclosing too much information from the original data.

**Teacher Forcing**  Applied to RNNs to avoid error propagation by feeding the correct token of the real data instead of the faulty predicted token to the next recurrent step.

## 3.9  Data Sets

Common data sets are essential for model training and evaluation. They allow researchers to compare their models against others in a meaningful manner. These data sets must be commonly available to the research community for that purpose. In Figure 55, we investigate if the generative models presented in this survey use public or private data sets. We find that 355 out of 388 models (excluding virtual environments) use at least one publicly available data set for their evaluation, while 33 do not disclose their data, of which only five models are used for privacy preservation, where restrictions often apply (e.g., healthcare). Additionally, 47 models utilize both private and public data sets for their evaluations. Each paper presented uses 1.74 data sets on average.

We further present the most often used public data sets used in our survey: MNIST [504], CIFAR-10 [505], celebA [506], ImageNet [507], LSUN (Bedrooms) [508], MS COCO [509] and SVHN [510] are collections of labeled or
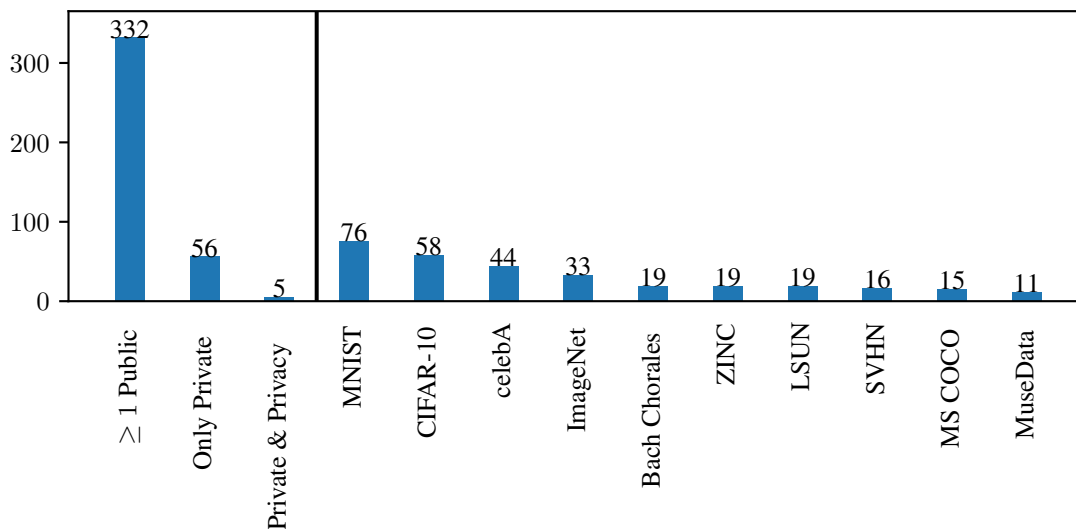
Figure 55: Amount of models (excluding virtual environments) using at least one public, only private, and only private data sets while needing to enforce privacy constraints. Additionally, we cover the most popular data sets with more than ten occurrences.

unlabeled images, which also are the most popular application area of SDG as we identified in Section 3.5. Other domains where commonly used data sets exist are music with the Bach Chorales [511] and MuseData[3] data sets, and graphs/molecules, where ZINC [512] is popular.

## 3.10   Model Performance

In Figure 56, we visualize the models and their relationships to other approaches that they claim to outperform in their respective evaluations. We can see that almost all arrows run from the top to the bottom, meaning that newer models tend to outperform older ones. Most models only evaluate against a small number of other works (on average 0.65 of the presented approaches), leading to a small in-degree and usually also out-degree. Some approaches like DCGAN [267] and [111] are popular to compare against, as shown by their large out-degree. GANs are overall most often compared against (83 times), followed by autoencoders (62 times), RNNs (44 times), CNNs (25 times), normalizing flow models (12 times), Boltzmann machines (7 times), RL approaches (7 times), diffusion models (5 times), NADEs (3 times), transformers (3 times), GMMNs (2 times), and PPGNs (1 time).

GANs and transformers tend to outperform other model types like CNNs and autoencoders. In graph/molecule generation domain, RNNs also hold up well. Since the resurgence of diffusion models in 2020 [461], they also outperform GANs in unconditional and text-conditional image generation.

We also investigate the evaluation metrics used by the presented models: The overall most commonly used metric is the Negative Log-Likelihood (NLL), which describes the probabilities assigned to the observed ground truth by the model. As Borji [19] points out, a low NLL score does not necessarily result in high data quality, and the metric is difficult to compute for high-dimensional data. Other common metrics for image data and GANs in particular are Inception Score (IS) and Fréchet Inception Distance (FID), which use a pre-trained image classifier to compare the real and fake data distributions [19]. We also often encounter evaluations by humans using a **MOS!** (**MOS!**) on a scale from one to five or one to ten to rate several features of the produced data. A comprehensive evaluation based on these metrics is not possible, due to most works using different metrics specialized for their task, and even if the same metric is used, the data sets (see Section 3.9) are often different.

## 3.11   Privacy

Many modern applications of ML are in areas such as the health care sector, where sensitive data of real persons has to be processed. This leads to the problem that large amounts of data for training and evaluation of models are required

---

[3]www.musedata.org

Figure 56: The relationships between models and their performance predecessors. The fill color indicates the model category as in Figure 44 while the border color shows how often other models have outperformed a model.
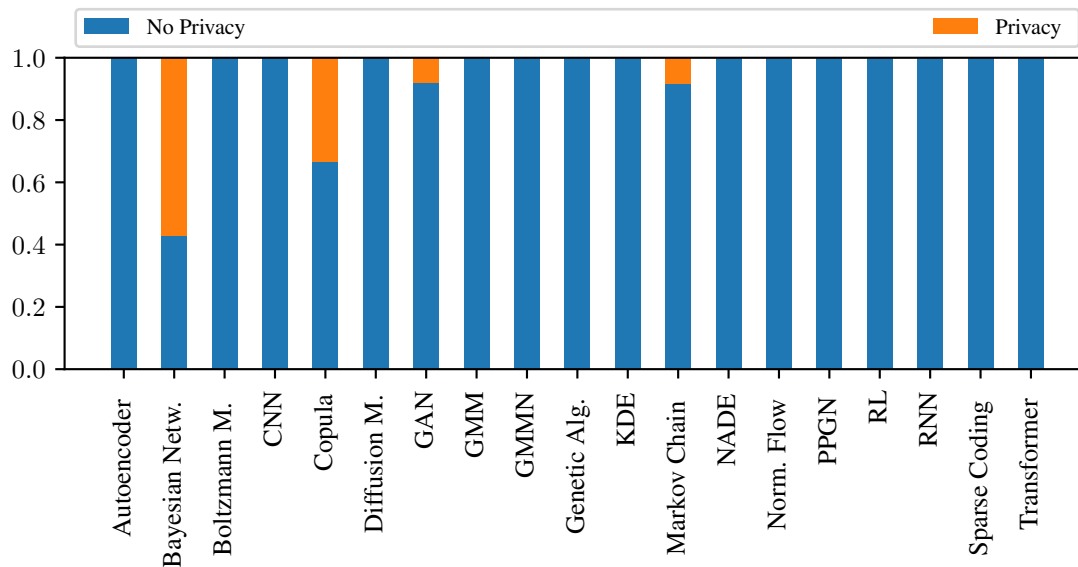
Figure 57: Amount of models with privacy considerations per model category absolute (top) and normalized (bottom).
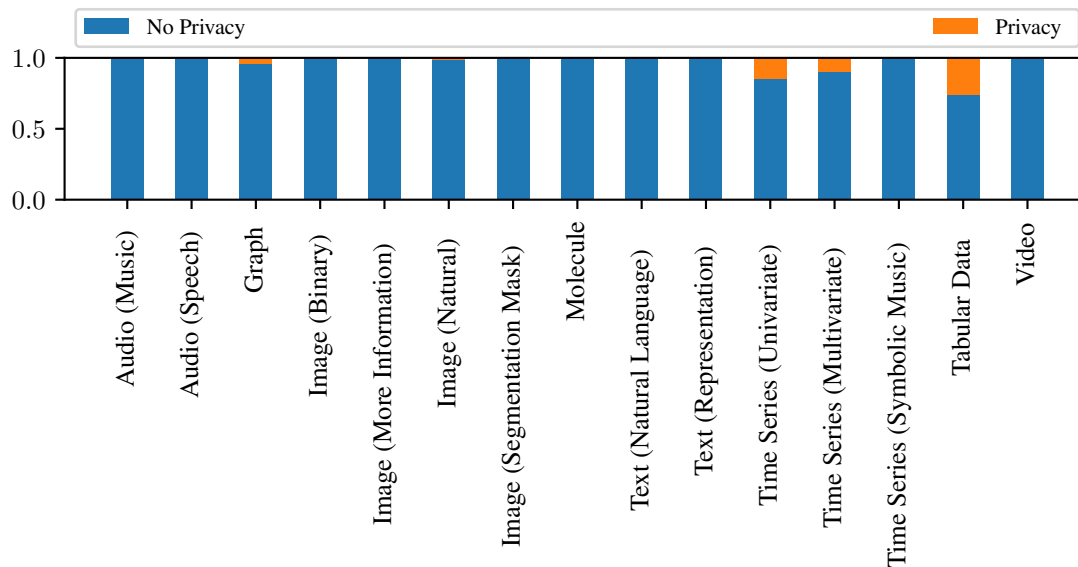


Figure 58: Privacy considerations per output data type absolute (top) and normalized (bottom).

but cannot be disclosed. SDG can be a sensible solution to this problem [95], but proving that no sensitive information is leaked by the SDG model requires special techniques such as differential privacy [52].

In Figure 57, we look at our surveyed models grouped by category and evaluate whether they provide a privacy guarantee. Usually, simplistic models with limited learning capabilities, such as BNs, Markov chains, and copulas, are used to generate private data. These models have the advantage that they can be inspected and modified by humans and with simple distribution modifications or noise injection. Thus, it becomes difficult to determine and extract information from real or personal data. More complex neural network models are seldom used, which is likely because they lack the advantages above and learn lots of details about their training data. That is, to make models more privacy-preserving, the realism and, ultimately, the utility of the data is reduced [317]. GANs, however, are an exception: They have the

advantage that their generator never sees the real training data, and it is demonstrated that with thoughtful design of the loss function and the model architecture, privacy guarantees for these models can be given [317].

In Figure 58, we evaluate the types of private data provided by our presented models. The most used data type is tabular data, which encapsulates EHRs and most other personal information often encountered in the real world. The only other data types we encountered were graphs and time series, which can be used to store mobility trajectories of persons or other medical data (e.g., an electrocardiogram). We did not encounter private SDG of higher-dimensional data like audio, images, video, or text.

## 3.12 Summary

In the earlier sections, we took a look at various aspects of SDG models that we now summarize: The most popular model type is the GAN, which is flexible and, by design, can create large amounts of data because it does not directly train on the training data. RNNs and CNNs are used as standalone models but also serve as building blocks for GANs, RL approaches, diffusion models, and autoencoders. They are suitable for generating sequences of samples and high-dimensional data (e.g., audio, images) respectively.

We find visual data generation, that is, images and videos, to be by far the most essential use case for SDG. Virtual environments are exclusively applied to this domain. At the same time, GANs and autoencoders are flexibly employed, and RNNs are preferably used for sequential data such as time series and text. Most models sample the data "in one go". In contrast, iterative sample refinement was less popular despite achieving similar results until recently, when diffusion models quickly became competitive against GANs for image generation.

Our performance evaluation found that newer models usually outperform older ones. Especially GANs, transformers, diffusion models, and RNNs, sometimes combined with RL, often come out on top. But we also encountered two significant problems: First, no common standardized evaluation metric for SDG models exists. FID, IS, NLL, and human evaluation are frequently used but are not suited for all tasks. The second problem is the evaluation data, which is also not standardized. We identified several common data sets for different domains, but direct comparison of models was often impossible due to different metric and data set combinations. Another aspect of model performance that was seldom mentioned is the computational complexity of models, meaning the training and sampling resources and time and the amount of training data required to achieve the quality of the presented results. A more systematic evaluation approach could solve the comparability problem: To measure the quality of generated images, for example, a model author could compare against a fixed set of other popular models (e.g., DCGAN [267]), on a larger set of common data sets (e.g., celebA, MNIST, CIFAR-10) using predetermined metrics (e.g., FID, IS). This common foundation would allow for a more precise comparison of approaches.

According to our findings and also the research of others [513], privacy-preserving data generation is still in the early stages of development. It is limited to low-dimensional data like table entries or time series and is usually performed with simple models observable by humans like Markov chains or Bayesian networks. The only more complex model that seems to be suitable for this task is the GAN, whose generator never sees the actual data. The main challenge also identified by others [513] is the trade-off between data utility and privacy, which means that the modifications required to make the data more private reduce the representativeness. Another problem is that the more complex and powerful neural-network-based approaches are known [514] to covertly encode individual samples from the training data in their parameters that can be reconstructed.

## 3.13 Guideline for Synthetic Data Generation

After classifying a large amount of SDG models and presenting our findings in written and visual form, we finally provide a guideline for model selection for various use cases. We explain our suggestions in written form and finally illustrate them as a decision tree in Figure 59:

- Use recent and up-to-date models. More powerful models are usually more expensive to train.
- For image generation and related purposes, diffusion models and GANs are the best options quality-wise. GANs are better for situations where less training data is available. Autoencoders are a less powerful but also a less resource-demanding alternative.
- Generation of sequential data like symbolic music, time series, most graph/molecule representations, and text is a suitable task for transformers and RNNs. Markov chains are a less powerful but simple and efficient alternative that humans can inspect.
- Autoencoders, especially VAEs, are good unsupervised feature learners that can disentangle and recombine features from training data to produce new data with desired properties.

- CNNs are a common building block of models to make them suitable to process high-dimensional data like large-resolution images and audio waveforms, whose sizes otherwise significantly slow down the training process or increase model size beyond a processable point.

- For the guidance of SDG models towards producing data with specific properties (e.g., molecular properties/validity), two approaches have proven to be simple yet powerful: RL, especially in combination with RNNs, allows users to define rewards for SDG models to produce the desired output. Models like PPGNs pair a trained GAN generator with a classifier to fine-tune for a specific task. Such approaches can greatly reduce the required training data.

- For private data generation, Markov chains or Bayesian networks with modified probabilities or noise injection in the sampling process can be used. GANs are a modern and more complex solution to private data generation when applied correctly [317].

- When visual data with highly accurate labels and high configurability is required, virtual environments are a good option, but they often require a considerable amount of human interaction to build and configure.
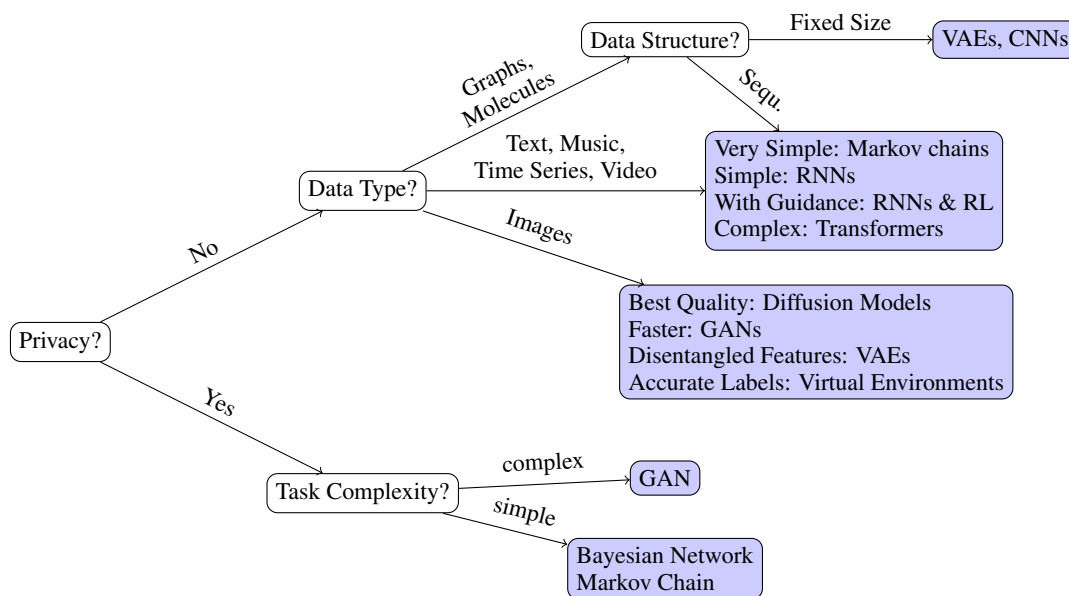
Figure 59: A simple decision tree for model selection. (Sequ. = Sequential)

# 4 Related Work

This section presents prior papers that compare or classify generative models for synthetic data. In Section 4.1, we provide an overview of previous works that selectively comprise and compare models in specific domains like healthcare privacy or graph generation. In Section 4.2, we focus on literature that aims to organize approaches for SDG comprehensively by specific aspects to provide an overview or guidance to novice users.

## 4.1 Domain- or Model-Specific Overviews and Comparisons

Fernández and Vico [515] summarize the research done in the field of algorithmic music composition to provide a comprehensive survey of various generation approaches: Grammars, knowledge-based systems, Markov chains, artificial neural networks, evolutionary (genetic) algorithms, and cellular automata.

Goodfellow et al. [3] create a tutorial on GANs. They show how these models work and how they can be improved for specific tasks. Further, different specialized applications of GANs in literature are shown and explained. GANs are also compared to other approaches, such as belief networks, autoencoders, and Boltzmann machines, regarding how the likelihood of generated samples can be computed.

Briot et al. [11] issue a large-scale survey on music generation via deep learning. They cover different types of musical content (melody, polyphony, performed by humans/machines), representations (Formats: MIDI, piano roll, text.

Encodings: Scalar, one-hot, many-hot), strategies (single-step/iterative feed-forward, sampling, etc.), challenges (e.g., variability, interactivity, originality) and deep neural network architectures. In this work, standard feed-forward and recurrent neural networks, autoencoders, and RBM architectures are covered and compared with the five criteria above.

He et al. [516] provide an overview of different models and benchmarks for image captioning. They capture end-to-end encoder-decoder frameworks like CNN-RNN combinations and also cover an attention mechanism applied to subregions of the image to improve the decoding. Other approaches are compositional frameworks that generate and arrange tags to generate captions, GANs, autoencoders, and RL. The authors aim to highlight the importance of image-to-text generation and encourage newcomers to contribute to this topic.

Jørgensen et al. [13] review the VAE as an alternative to quantum-mechanical computations with lower computational cost to generate new molecular structures and predict their properties. They also discuss approaches to improve the realism of generated data by using grammar-based instead of character-based encoders and decoders and propose other models like GANs and RL agents.

Korakakis et al. [5] provide an overview of virtual environments and their usage in literature. They illustrate how synthetic data obtained from CAD renders and video games can be used to improve object detection or classification models or train RL agents. The authors identify a trend of steadily increasing usage of virtual environments for cheap yet effective model learning, especially for computer vision tasks.

Gaidon et al. [8] present nine papers exploring novel ways of generating and using synthetic data for computer vision tasks. They summarize these nine works and also cover some of the problems encountered by the literature, such as generation challenges or the "sim2real" domain gap, which describes the problem of fitting models trained on synthetic data to real applications despite the generated data often being different in some way (e.g., lack of photo-realism). Nevertheless, more and higher-quality synthetic data might help overcome the limitations of current computer vision algorithms, according to the authors.

Kulkarni et al. [437] evaluate the performance of RNNs, GANs, and copulas on the generation of synthetic human mobility trajectories in terms of privacy preservation, long-range dependencies, the statistical similarity of the distributions, training and generation time, circadian rhythms, and semantic and geographic similarity. They conclude that copulas have preferable statistical and semantic properties over the neural network models, which also consume more time and are less computationally efficient. Further, they assess that a utility metric to measure and maximize privacy and statistical similarity jointly could improve the usability of synthetic trajectories, but is not yet available.

Hong et al. [4] give a detailed introduction to GANs and various recently proposed objective functions for them. Further, the combination of a GAN with an autoencoder is discussed. Finally, multiple applications of GANs in different tasks and fields are covered, and the pros and cons of this model type, such as convergence towards an optimal solution, are highlighted.

Yi et al. [517] review GANs and their applications in medical imaging. They collect literature using GANs for medical purposes such as image synthesis, segmentation, and reconstruction/repair and classify them according to GAN method (e.g., pix2pix [335]), adversarial loss type and quantitative measures (e.g., Wasserstein distance [285]) used.

Iqbal et al. [10] survey deep learning text generation models and the progress made from 2015 onwards. They focus on RNNs such as LSTMs, GRUs and bidirectional RNNs, CNNs, VAEs, and GANs in combination with RL. Various representations (Word2Vec, Glove, FastText), optimization techniques (stochastic gradient descent, RMSProp, AdaGrad, Adam), activation functions (Sigmoid, ReLu), and evaluation methods (Rouge, BLEU) for text generation models are also introduced.

Tsirikoglou et al. [9] collect and compare different image synthesis and augmentation methods. They identify that the visual data generation pipeline consists of two parts: *Content/scene generation*, which means generating the features of the virtual environment, and *rendering*, which simulates the light transport and perception of sensors. Further, synthetic visual training data has four requirements to be useful: Feature variation and coverage, domain realism, automatic generation of annotations and meta-data, and scalability to large numbers of data points. Over 40 generative models from recent literature are categorized by their modeling and rendering approach and compared regarding image quality and what tasks they are applied to.

Guo et al. [7] extensively cover and analyze the recent literature of deep generative models for graph generation, including BNs, VAEs, GANs, RNNs, flow-based learning and RL. They provide taxonomies of models for conditional and unconditional graph generation and describe the evaluation metrics applicable in this domain. Finally, the application fields of deep graph generation, such as the analysis of interaction dynamics in social networks, the creation of molecules, or anomaly detection, are discussed.

Seib et al. [6] discuss different techniques to improve neural network training results on computer vision tasks in urban and traffic environments without acquiring additional real-world data. The topics explored are data augmentation, transfer learning, which describes the fine-tuning of pre-trained models for another task with few data samples, and approaches to generating synthetic data. Different 3D engines (Unreal Engine and Unity) and video games (GTA V) and their usage in literature are covered. A future outlook towards GANs and their image-to-image translation capabilities is also provided.

Abufadda et al. [518] examine and summarize related works about SDG in the healthcare domain, especially presenting many GAN approaches. For each paper, they highlight the research field, used methods and results, and indicate, whether the models are suitable for their task. The authors conclude that generalized solutions for utility and efficiency evaluations would improve the model selection process.

Dankar et al. [12] propose an overall utility score for masked synthetic data sets, where features of the original data have to be kept secret. For that, available metrics are categorized by the measure they aim to preserve, and from each of the four categories (attribute, bivariate, population, and application fidelity), one suitable metric is chosen to determine the final utility of a generative model. The utility measurement approach is evaluated on four recent models and 19 data sets with different sizes and features. They conclude that each privacy-enhancing technology decreases data utility, and the acceptable or necessary decrease to achieve privacy is unknown.

All the literature above provides useful insight into specific application fields of synthetic data and model types. However, the scope of the individual works is usually limited to one domain and a small selection of models and architectures, making them rather unsuitable as a comprehensive introduction to SDG.

## 4.2 Comprehensive Reviews

Turhan et al. [14] conduct a comprehensive review of generative models, especially focussing on deep learning models like GANs and autoencoders for image generation. They highlight use cases for these models and classify them into five categories: Unsupervised fundamental models (RBM, DBN and DBM), autoencoder-based models, autoregressive models (CNNs and RNNs), GAN-based models and autoencoder-GAN hybrid models. Further, a relation diagram of all presented models is compiled that is particularly useful to beginners in this topic.

Oussidi et al. [15] consolidate promising types of generative models like RBMs, DBNs, DBMs, VAEs and GANs and describe the three models PixelRNN [115], DRAW [97] and NADE [146] in detail. They also put generative models into two categories:

**Cost function-based models**  Models that optimize parameters based on cost/loss, like autoencoders and GANs.

**Energy-based models**  The joint probability is defined by an energy function, which measures the compatibility of variable configurations [279]. This approach is used by Boltzmann machines and their derivatives.

They study their advantages, limitations, and potential for the future. They find that energy-based models are more complex to combine than directed graphical models like feed-forward neural networks, and deep networks often suffer from vanishing or exploding gradients during training. Hence, the bottom layers barely learn anything, while the top layers quickly reach an optimal state.

Wang et al. [503] propose an architecture- and loss-based taxonomy for GANs and highlight the significant advances made with them in recent years in the field of computer vision. They further discuss the challenges in these tasks, such as image quality, diversity, and training stability, but also the risks involved in being able to generate high-quality fake data, such as fake evidence of crimes or events.

Harshvardhan et al. [16] compile a comprehensive survey of generative models, highlight some noteworthy contributions from literature, and implement and evaluate each presented model to help the readers pick the best solution for their use case. Their high-level review incorporates GMMs, HMMs, Latent Dirichlet Allocation (LDA), Boltzmann machines, VAEs and GANs. Further, the models are also classified by their learning type (un-/semi- and supervised learning) and their model architecture (machine learning or subset deep learning).

Eigenschink et al. [17] present a data-driven framework that evaluates synthetic data generation models independently of their internal workings. The authors exclusively cover models for synthetic sequential data to complement previous reviews that only focused on particular data (e.g., time-series, videos, text) or model (e.g., GANs) types. The models compared in this work are RNNs, CNNs, transformers, autoencoders, autoregressive neural networks, and GANs. They are compared in terms of

**Representativeness**  How well the synthetic data captures distributions and dependencies between the distributions, for instance, hair colors and eye distances in a face image data set.

**Novelty** Do new samples resemble samples from the original data set, or are they new observations from the latent distribution? This is especially important for use cases concerned with privacy, like healthcare, because original samples must not be leaked.

**Realism** Similar to Representativeness, a statistical measure, but on a per-sample level: Often, human judgment is used to determine whether synthetic samples are realistic.

**Diversity** Measure of similarity between individual synthetic data points, for example, the average Euclidean distance to their nearest neighbors [271].

**Coherence** Often implicitly evaluated with realism, coherence describes whether the internal structure of single synthetic data points is consistent. Varying Notes in music, for example, are natural and necessary, but random genre changes are not and indicate bad coherence.

Further, the importance of individual criteria is evaluated in different domains such as Natural Language Processing (NLP), audio processing, or anonymization of healthcare data, concluding that representativeness and realism are most important for NLP. Speech, music, and video tasks mostly rely on realism and coherence, private synthetic EHRs in healthcare applications need representativeness, novelty, and realism, and mobility trajectories additionally require coherence.

Nikolenko [2] released a book about synthetic data and how it is currently used for ML. The book mainly covers the computer vision topic, especially the collaboration of deep learning models for classification and GANs as synthetic data generators, and also provides a historical perspective on it, but other application fields like computer security, bioinformatics and NLP also make an appearance. Common problems of generative models, such as privacy concerns and the challenge of domain shift, are also treated in separate sections, and potential solutions are discussed. Finally, an outlook on potential future improvements for SDG is given in the form of RL or the incorporation of domain knowledge in generative models.

The works presented here have a more comprehensive view on SDG methods, some even providing a historical perspective [2]. Even though these surveys and reviews are all recent, released in 2018 or later, new approaches like transformers [18] are often missed entirely, and many of them focus primarily on GANs or autoencoders. Also, essential domains like music, which could play an important future role in the film and video game industry, are often forgotten. Another important use of a comprehensive overview is the comparison of approaches and guidance of users to select an appropriate model for their use case, which is only provided extensively by [16] and coarsely by [17]. The main difference between our work and the presented works is listed in Table 1.

| Approach | Data Type(s) | Model(s) | #Works | Investigated Aspects |
|---|---|---|---|---|
| [515] | Music | 6 | 267 | Creativity |
| [11] | Music | 7 | 12 | Creativity, Model capabilities |
| [10] | Text | 5 | 30 | Similarity to real data |
| [14] | Image | 2 | 45 | Model relationships |
| [15] | Image | 5 | 13 | Limitations/Advantages |
| [503] | Image | 1 | 36 | Performance, Architecture |
| [16] | Comprehensive | 6 | 20-30 | Lim./Adv., Perf., Implementation |
| [17] | Sequential Data | 6 | 17 | Creativity, Data Quality |
| **This article** | **8** types (15 sub-types) | **20** (42 sub-types) | **417** | **9** (see Section 3.1) |

Table 1: Comparison of this survey to others.

## 5 Conclusion

The surge in applications in ML has transformed various fields, but limited training data, expensive acquisition, and privacy laws hinder ML model efficacy. SDG emerges as a solution, yet the diverse and rapidly evolving landscape of SDG models, spanning decades of development, poses a challenge for decision-makers. To this end, we conducted a comprehensive survey of 417 SDG model papers, resulting in the categorization of these models into 20 distinct types (42 sub-types). The classification was performed based on criteria extracted from related work and identified during our survey. The key findings from our comprehensive classification are as follows:

- Computer vision is the most popular application field, and GANs are the most popular SDG models.
- Different data types require different model types. RNNs and transformers are more suitable for sequential data, while CNNs, GANs or autoencoders are mostly used for data with static size.

- We also observed a trend of combining different model types. GANs and autoencoders often act as frameworks with RNNs or CNNs as building blocks. RNNs are often combined with RL components to guide the generation process.

- Our analysis revealed challenges in performance evaluation, citing the absence of standardized metrics and datasets. That is, a common set of metrics, datasets, and reference models is needed to enhance comparability among SDG models.

- We identified a nascent stage in the development of privacy-preserving data generation, where simplistic models like Markov chains, BN, and genetic algorithms are prevalent, with GAN being the only more complex neural network-based model.

We are convinced that our work provides a valuable resource for researchers entering the field, aiding them in selecting suitable approaches for their specific purposes. Furthermore, it is crucial to underscore the necessity for future research to (i) delve into the training and sampling costs of SDG models for a more comprehensive classification and (ii) establish a systematic evaluation approach that enhances the overall understanding, the usability, and the comparability among SDG models.

# References

[1] Konstantinos Stathoulopoulos, Joel Klinger, and Juan Mateos-Garcia. Is ai eating software? an analysis of ai/ml research trends using scientific pre-prints. Apr 2018. Accessed: 2022-05-24.

[2] Sergey I Nikolenko et al. *Synthetic data for deep learning*. Springer, 2021.

[3] Ian Goodfellow. Nips 2016 tutorial: Generative adversarial networks. *arXiv preprint arXiv:1701.00160*, 2016.

[4] Yongjun Hong, Uiwon Hwang, Jaeyoon Yoo, and Sungroh Yoon. How generative adversarial networks and their variants work: An overview. *ACM Computing Surveys (CSUR)*, 52(1):1–43, 2019.

[5] Michalis Korakakis, Phivos Mylonas, and Evaggelos Spyrou. A short survey on modern virtual environments that utilize ai and synthetic data. In *MCIS*, page 34, 2018.

[6] Viktor Seib, Benjamin Lange, and Stefan Wirtz. Mixing real and synthetic data to enhance neural network training–a review of current approaches. *arXiv preprint arXiv:2007.08781*, 2020.

[7] Xiaojie Guo and Liang Zhao. A systematic survey on deep generative models for graph generation. *arXiv preprint arXiv:2007.06686*, 2020.

[8] Adrien Gaidon, Antonio Lopez, and Florent Perronnin. The reasonable effectiveness of synthetic visual data. *International Journal of Computer Vision*, 126(9):899–901, Sep 2018.

[9] Apostolia Tsirikoglou, Gabriel Eilertsen, and Jonas Unger. A survey of image synthesis methods for visual machine learning. In *Computer Graphics Forum*, volume 39, pages 426–451. Wiley Online Library, 2020.

[10] Touseef Iqbal and Shaima Qureshi. The survey: Text generation models in deep learning. *Journal of King Saud University - Computer and Information Sciences*, 2020.

[11] Jean-Pierre Briot, Gaëtan Hadjeres, and François-David Pachet. Deep learning techniques for music generation–a survey. *arXiv preprint arXiv:1709.01620*, 2017.

[12] Fida K Dankar, Mahmoud K Ibrahim, and Leila Ismail. A multi-dimensional evaluation of synthetic data generators. *IEEE Access*, 10:11147–11158, 2022.

[13] Peter B Jørgensen, Mikkel N Schmidt, and Ole Winther. Deep generative models for molecular science. *Molecular informatics*, 37(1-2):1700133, 2018.

[14] Ceren Güzel Turhan and Hasan Sakir Bilge. Recent trends in deep generative models: a review. In *2018 3rd International Conference on Computer Science and Engineering (UBMK)*, pages 574–579, 2018.

[15] Achraf Oussidi and Azeddine Elhassouny. Deep generative models: Survey. In *2018 International Conference on Intelligent Systems and Computer Vision (ISCV)*, pages 1–8. IEEE, 2018.

[16] GM Harshvardhan, Mahendra Kumar Gourisaria, Manjusha Pandey, and Siddharth Swarup Rautaray. A comprehensive survey and analysis of generative models in machine learning. *Computer Science Review*, 38:100285, 2020.

[17] Peter Eigenschink, Stefan Vamosi, Ralf Vamosi, Chang Sun, Thomas Reutterer, and Klaudius Kalcher. Deep generative models for synthetic data. *ACM Computing Surveys*, 2021.

[18] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.

[19] Ali Borji. Pros and cons of gan evaluation measures: New developments. *Computer Vision and Image Understanding*, 215:103329, 2022.

[20] Jake VanderPlas. *Python data science handbook: Essential tools for working with data*. O'Reilly Media, Inc., 2016.

[21] Aaron Van den Oord and Benjamin Schrauwen. Factoring variations in natural images with deep gaussian mixture models. *Advances in neural information processing systems*, 27, 2014.

[22] Heiga Zen and Andrew Senior. Deep mixture density networks for acoustic modeling in statistical parametric speech synthesis. In *2014 IEEE international conference on acoustics, speech and signal processing (ICASSP)*, pages 3844–3848. IEEE, 2014.

[23] Sakyajit Bhattacharya, Oishee Mazumder, Dibyendu Roy, Aniruddha Sinha, and Avik Ghose. Synthetic data generation through statistical explosion: Improving classification accuracy of coronary artery disease using ppg. In *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1165–1169. IEEE, 2020.

[24] Richard Durbin, Sean R Eddy, Anders Krogh, and Graeme Mitchison. *Biological sequence analysis: probabilistic models of proteins and nucleic acids*. Cambridge university press, 1998.

[25] José Luis Triviño-Rodriguez and Rafael Morales-Bueno. Using multiattribute prediction suffix graphs to predict and generate music. *Computer Music Journal*, 25(3):62–79, 2001.

[26] François Pachet. Beyond the cybernetic jam fantasy: The continuator. *IEEE Computer Graphics and Applications*, 24(1):31–35, 2004.

[27] Stanisław A Raczyński, Satoru Fukayama, and Emmanuel Vincent. Melody harmonization with interpolated probabilistic models. *Journal of New Music Research*, 42(3):223–235, 2013.

[28] Maximos Kaliakatsos-Papakostas and Emilios Cambouropoulos. Probabilistic harmonization with fixed intermediate chord constraints. In *ICMC*, 2014.

[29] Vincent Bindschaedler and Reza Shokri. Synthesizing plausible privacy-preserving location traces. In *2016 IEEE Symposium on Security and Privacy (SP)*, pages 546–563. IEEE, 2016.

[30] MarcAurelio Ranzato, Arthur Szlam, Joan Bruna, Michael Mathieu, Ronan Collobert, and Sumit Chopra. Video (language) modeling: a baseline for generative models of natural videos. *arXiv preprint arXiv:1412.6604*, 2014.

[31] Yoshua Bengio, Réjean Ducharme, and Pascal Vincent. A neural probabilistic language model. *Advances in Neural Information Processing Systems*, 13, 2000.

[32] Gabriele Barbieri, François Pachet, Pierre Roy, and Mirko Degli Esposti. Markov constraints for generating lyrics with style. In *Ecai*, volume 242, pages 115–120, 2012.

[33] Pierre Roy and François Pachet. Enforcing meter in finite-length markov sequences. In *Twenty-Seventh AAAI Conference on Artificial Intelligence*, 2013.

[34] Jonathan P Forsyth and Juan P Bello. Generating musical accompaniment using finite state transducers. In *16th International Conference on Digital Audio Effects (DAFx-13)*, pages 1–7, 2013.

[35] Alexandre Papadopoulos, Pierre Roy, and François Pachet. Avoiding plagiarism in markov sequence generation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 28, 2014.

[36] Alexandre Papadopoulos, Pierre Roy, and François Pachet. Assisted lead sheet composition using flowcomposer. In *International Conference on Principles and Practice of Constraint Programming*, pages 769–785. Springer, 2016.

[37] Raymond P Whorley and Darrell Conklin. Music generation from statistical models of harmony. *Journal of New Music Research*, 45(2):160–183, 2016.

[38] Barbara Draghi, Zhenchen Wang, Puja Myles, and Allan Tucker. Bayesboost: Identifying and handling bias using synthetic data generators. In *Third International Workshop on Learning with Imbalanced Domains: Theory and Applications*, pages 49–62. PMLR, 2021.

[39] Todd Andrew Stephenson. An introduction to bayesian network theory and usage. Technical report, Idiap, 2000.

[40] Haipeng Guo and William Hsu. A survey of algorithms for real-time bayesian network inference. In *Join Workshop on Real Time Decision Support and Diagnosis Systems*, 2002.

[41] Diederik P Kingma. Fast gradient-based inference with continuous latent variable models in auxiliary form. *arXiv preprint arXiv:1306.0733*, 2013.

[42] Gideon Schwarz. Estimating the dimension of a model. *The annals of statistics*, pages 461–464, 1978.

[43] Mikko Koivisto and Kismat Sood. Exact bayesian structure discovery in bayesian networks. *The Journal of Machine Learning Research*, 5:549–573, 2004.

[44] Tomi Silander and Petri Myllymaki. A simple approach for finding the globally optimal bayesian network structure. *arXiv preprint arXiv:1206.6875*, 2012.

[45] Ajit P Singh and Andrew W Moore. *Finding optimal Bayesian networks by dynamic programming*. Citeseer, 2005.

[46] Changhe Yuan, Brandon Malone, and Xiaojian Wu. Learning optimal bayesian networks using a* search. In *Twenty-Second International Joint Conference on Artificial Intelligence*, 2011.

[47] Jim Young, Patrick Graham, and Richard Penny. Using bayesian networks to create synthetic data. *Journal of Official Statistics*, 25(4):549, 2009.

[48] Susanne G Bøttcher and Claus Dethlefsen. deal: A package for learning bayesian networks. *Journal of statistical software*, 8:1–40, 2003.

[49] Patrick Graham, Jim Young, and Richard Penny. Multiply imputed synthetic data: Evaluation of hierarchical bayesian imputation models. *Journal of Official Statistics*, 25(2):245, 2009.

[50] Syunpei Suzuki and Tetsuro Kitahara. Four-part harmonization using bayesian networks: pros and cons of introducing chord nodes. *Journal of New Music Research*, 43(3):331–353, 2014.

[51] Jun Zhang, Graham Cormode, Cecilia M Procopiuc, Divesh Srivastava, and Xiaokui Xiao. Privbayes: Private data release via bayesian networks. *ACM Transactions on Database Systems (TODS)*, 42(4):1–41, 2017.

[52] Cynthia Dwork, Aaron Roth, et al. The algorithmic foundations of differential privacy. *Foundations and Trends® in Theoretical Computer Science*, 9(3–4):211–407, 2014.

[53] Wenzheng Chen, Huan Wang, Yangyan Li, Hao Su, Zhenhua Wang, Changhe Tu, Dani Lischinski, Daniel Cohen-Or, and Baoquan Chen. Synthesizing training images for boosting human 3d pose estimation. In *2016 Fourth International Conference on 3D Vision (3DV)*, pages 479–488. IEEE, 2016.

[54] Haoyue Ping, Julia Stoyanovich, and Bill Howe. Datasynthesizer: Privacy-preserving synthetic datasets. In *Proceedings of the 29th International Conference on Scientific and Statistical Database Management*, pages 1–5, 2017.

[55] Allan Tucker, Zhenchen Wang, Ylenia Rotalinti, and Puja Myles. Generating high-fidelity synthetic patient data for assessing machine learning healthcare software. *NPJ digital medicine*, 3(1):1–13, 2020.

[56] Peter Spirtes, Clark N Glymour, Richard Scheines, and David Heckerman. *Causation, prediction, and search*. MIT press, 2000.

[57] Yingrui Chen, Mark Elliot, and Joseph Sakshaug. Genetic algorithms in matrix representation and its application in synthetic data. In *UNECE Worksession on Statistical Confidentiality 2017*. 2017.

[58] Shingchern D You and Po-Sheng Liu. Automatic chord generation system using basic music theory and genetic algorithm. In *2016 IEEE International Conference on Consumer Electronics-Taiwan (ICCE-TW)*, pages 1–2. IEEE, 2016.

[59] Chien-Hung Liu and Chuan-Kang Ting. Polyphonic accompaniment using genetic algorithm with music theory. In *2012 IEEE Congress on Evolutionary Computation*, pages 1–7. IEEE, 2012.

[60] Ruslan Salakhutdinov and Geoffrey Hinton. Deep boltzmann machines. In David van Dyk and Max Welling, editors, *Proceedings of the Twelth International Conference on Artificial Intelligence and Statistics*, volume 5 of *Proceedings of Machine Learning Research*, pages 448–455, Hilton Clearwater Beach Resort, Clearwater Beach, Florida USA, 16–18 Apr 2009. PMLR.

[61] Paul Smolensky. Information processing in dynamical systems: Foundations of harmony theory. Technical report, Colorado Univ at Boulder Dept of Computer Science, 1986.

[62] Yuming Hua, Junhai Guo, and Hua Zhao. Deep belief networks and deep learning. In *Proceedings of 2015 International Conference on Intelligent Computing and Internet of Things*, pages 1–4. IEEE, 2015.

[63] Honglak Lee, Roger Grosse, Rajesh Ranganath, and Andrew Y Ng. Convolutional deep belief networks for scalable unsupervised learning of hierarchical representations. In *Proceedings of the 26th annual international conference on machine learning*, pages 609–616, 2009.

[64] Stefan Lattner, Maarten Grachten, and Gerhard Widmer. Imposing higher-level structure in polyphonic music generation using convolutional restricted boltzmann machines and constraints. *Journal of Creative Music Systems*, 2:1–31, 2018.

[65] Geoffrey E Hinton, Simon Osindero, and Yee-Whye Teh. A fast learning algorithm for deep belief nets. *Neural computation*, 18(7):1527–1554, 2006.

[66] Christopher Tosh. Mixing rates for the alternating gibbs sampler over restricted boltzmann machines and friends. In *International Conference on Machine Learning*, pages 840–849. PMLR, 2016.

[67] Greg Bickerman, Sam Bosley, Peter Swire, and Robert M Keller. Learning to create jazz melodies using deep belief nets. In *ICCC*, pages 228–237, 2010.

[68] Felix Sun. Deephear – composing and harmonizing music with neural networks. `https://fephsun.github.io/2015/09/01/neural-music.html`, 2015. Accessed: 2022-07-29.

[69] Ilya Sutskever and Geoffrey Hinton. Learning multilevel distributed representations for high-dimensional sequences. In *Artificial intelligence and statistics*, pages 548–555. PMLR, 2007.

[70] Ilya Sutskever, Geoffrey E Hinton, and Graham W Taylor. The recurrent temporal restricted boltzmann machine. *Advances in neural information processing systems*, 21, 2008.

[71] Roni Mittelman, Benjamin Kuipers, Silvio Savarese, and Honglak Lee. Structured recurrent temporal restricted boltzmann machines. In *International Conference on Machine Learning*, pages 1647–1655. PMLR, 2014.

[72] Yann LeCun, Fu Jie Huang, and Leon Bottou. Learning methods for generic object recognition with invariance to pose and lighting. In *Proceedings of the 2004 IEEE Computer Society Conference on Computer Vision and Pattern Recognition, 2004. CVPR 2004.*, volume 2, pages II–104. IEEE, 2004.

[73] Graham W Taylor, Geoffrey E Hinton, and Sam T Roweis. Two distributed-state models for generating high-dimensional time series. *Journal of Machine Learning Research*, 12(3), 2011.

[74] Roland Memisevic and Geoffrey Hinton. Unsupervised learning of image transformations. In *2007 IEEE Conference on Computer Vision and Pattern Recognition*, pages 1–8. IEEE, 2007.

[75] Roland Memisevic and Geoffrey E Hinton. Learning to represent spatial transformations with factored higher-order boltzmann machines. *Neural computation*, 22(6):1473–1492, 2010.

[76] Pascal Vincent, Hugo Larochelle, Yoshua Bengio, and Pierre-Antoine Manzagol. Extracting and composing robust features with denoising autoencoders. In *Proceedings of the 25th international conference on Machine learning*, pages 1096–1103, 2008.

[77] Pascal Vincent, Hugo Larochelle, Isabelle Lajoie, Yoshua Bengio, Pierre-Antoine Manzagol, and Léon Bottou. Stacked denoising autoencoders: Learning useful representations in a deep network with a local denoising criterion. *Journal of machine learning research*, 11(12), 2010.

[78] Andy Sarroff and Michael A Casey. Musical audio synthesis using autoencoding neural nets. In *Proceedings of the International Society for Music Information Retrieval Conference (ISMIR2014)*. International Society for Music Information Retrieval, 2014.

[79] Karol Gregor, Ivo Danihelka, Andriy Mnih, Charles Blundell, and Daan Wierstra. Deep autoregressive networks. In *International Conference on Machine Learning*, pages 1242–1250. PMLR, 2014.

[80] Jiwei Li, Minh-Thang Luong, and Dan Jurafsky. A hierarchical neural autoencoder for paragraphs and documents. *arXiv preprint arXiv:1506.01057*, 2015.

[81] Antti Rasmus, Mathias Berglund, Mikko Honkala, Harri Valpola, and Tapani Raiko. Semi-supervised learning with ladder networks. In C. Cortes, N. Lawrence, D. Lee, M. Sugiyama, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 28. Curran Associates, Inc., 2015.

[82] Lars Maaløe, Casper Kaae Sønderby, Søren Kaae Sønderby, and Ole Winther. Auxiliary deep generative models. In *International conference on machine learning*, pages 1445–1453. PMLR, 2016.

[83] Rim Assouel, Mohamed Ahmed, Marwin H Segler, Amir Saffari, and Yoshua Bengio. Defactor: Differentiable edge factorization-based probabilistic graph generation. *arXiv preprint arXiv:1811.09766*, 2018.

[84] Thomas N Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. *arXiv preprint arXiv:1609.02907*, 2016.

[85] Dmitry Ulyanov, Andrea Vedaldi, and Victor Lempitsky. It takes (only) two: Adversarial generator-encoder networks. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 32, 2018.

[86] Shain Shahid Chowdhury, Soukaïna Filali Boubrahimi, and Shah Muhammad Hamdi. Time series data augmentation using time-warped auto-encoders. In *2021 20th IEEE International Conference on Machine Learning and Applications (ICMLA)*, pages 467–470, 2021.

[87] Peter Dayan, Geoffrey E Hinton, Radford M Neal, and Richard S Zemel. The helmholtz machine. *Neural computation*, 7(5):889–904, 1995.

[88] Yifeng Li and Xiaodan Zhu. Exploring helmholtz machine and deep belief net in the exponential family perspective. In *ICML 2018 Workshop on Theoretical Foundations and Applications of Deep Generative Models*, 2018.

[89] Geoffrey E Hinton, Peter Dayan, Brendan J Frey, and Radford M Neal. The" wake-sleep" algorithm for unsupervised neural networks. *Science*, 268(5214):1158–1161, 1995.

[90] Yoshua Bengio, Li Yao, Guillaume Alain, and Pascal Vincent. Generalized denoising auto-encoders as generative models, 2013.

[91] Salah Rifai, Pascal Vincent, Xavier Muller, Xavier Glorot, and Yoshua Bengio. Contractive auto-encoders: Explicit invariance during feature extraction. In *Icml*, 2011.

[92] Salah Rifai, Yoshua Bengio, Yann Dauphin, and Pascal Vincent. A generative process for sampling contractive auto-encoders. *arXiv preprint arXiv:1206.6434*, 2012.

[93] Yoshua Bengio, Grégoire Mesnil, Yann Dauphin, and Salah Rifai. Better mixing via deep representations. In *International conference on machine learning*, pages 552–560. PMLR, 2013.

[94] Yoshua Bengio, Eric Laufer, Guillaume Alain, and Jason Yosinski. Deep generative stochastic networks trainable by backprop. In *International Conference on Machine Learning*, pages 226–234. PMLR, 2014.

[95] Lei Xu et al. *Synthesizing tabular data using conditional GAN*. PhD thesis, Massachusetts Institute of Technology, 2020.

[96] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.

[97] Karol Gregor, Ivo Danihelka, Alex Graves, Danilo Rezende, and Daan Wierstra. Draw: A recurrent neural network for image generation. In *International Conference on Machine Learning*, pages 1462–1471. PMLR, 2015.

[98] Danilo Rezende, Ivo Danihelka, Karol Gregor, Daan Wierstra, et al. One-shot generalization in deep generative models. In *International conference on machine learning*, pages 1521–1529. PMLR, 2016.

[99] Max Jaderberg, Karen Simonyan, Andrew Zisserman, et al. Spatial transformer networks. *Advances in neural information processing systems*, 28, 2015.

[100] Irina Higgins, Loic Matthey, Arka Pal, Christopher Burgess, Xavier Glorot, Matthew Botvinick, Shakir Mohamed, and Alexander Lerchner. beta-vae: Learning basic visual concepts with a constrained variational framework. 2016.

[101] Xi Chen, Yan Duan, Rein Houthooft, John Schulman, Ilya Sutskever, and Pieter Abbeel. Infogan: Interpretable representation learning by information maximizing generative adversarial nets. *Advances in neural information processing systems*, 29, 2016.

[102] Tejas D Kulkarni, William F Whitney, Pushmeet Kohli, and Josh Tenenbaum. Deep convolutional inverse graphics network. *Advances in neural information processing systems*, 28, 2015.

[103] Jakub Tomczak and Max Welling. Vae with a vampprior. In *International Conference on Artificial Intelligence and Statistics*, pages 1214–1223. PMLR, 2018.

[104] Otto Fabius and Joost R Van Amersfoort. Variational recurrent auto-encoders. *arXiv preprint arXiv:1412.6581*, 2014.

[105] Daniel Jiwoong Im, Sungjin Ahn, Roland Memisevic, and Yoshua Bengio. Denoising criterion for variational auto-encoding framework. *arXiv preprint arXiv:1511.06406*, 2015.

[106] Yuri Burda, Roger Grosse, and Ruslan Salakhutdinov. Importance weighted autoencoders. *arXiv preprint arXiv:1509.00519*, 2015.

[107] Samuel R Bowman, Luke Vilnis, Oriol Vinyals, Andrew M Dai, Rafal Jozefowicz, and Samy Bengio. Generating sentences from a continuous space. *arXiv preprint arXiv:1511.06349*, 2015.

[108] Jonathan Huang and Kevin Murphy. Efficient inference in occlusion-aware generative models of images. *arXiv preprint arXiv:1511.06362*, 2015.

[109] Anders Boesen Lindbo Larsen, Søren Kaae Sønderby, Hugo Larochelle, and Ole Winther. Autoencoding beyond pixels using a learned similarity metric. In *International conference on machine learning*, pages 1558–1566. PMLR, 2016.

[110] Xinchen Yan, Jimei Yang, Kihyuk Sohn, and Honglak Lee. Attribute2image: Conditional image generation from visual attributes. In *European conference on computer vision*, pages 776–791. Springer, 2016.

[111] 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. *arXiv preprint arXiv:1610.02415*, 2016.

[112] David Weininger. Smiles, a chemical language and information system. 1. introduction to methodology and encoding rules. *Journal of chemical information and computer sciences*, 28(1):31–36, 1988.

[113] Xi Chen, Diederik P Kingma, Tim Salimans, Yan Duan, Prafulla Dhariwal, John Schulman, Ilya Sutskever, and Pieter Abbeel. Variational lossy autoencoder. *arXiv preprint arXiv:1611.02731*, 2016.

[114] Ishaan Gulrajani, Kundan Kumar, Faruk Ahmed, Adrien Ali Taiga, Francesco Visin, David Vazquez, and Aaron Courville. Pixelvae: A latent variable model for natural images. *arXiv preprint arXiv:1611.05013*, 2016.

[115] Aaron Van Oord, Nal Kalchbrenner, and Koray Kavukcuoglu. Pixel recurrent neural networks. In *International conference on machine learning*, pages 1747–1756. PMLR, 2016.

[116] Matt J Kusner, Brooks Paige, and José Miguel Hernández-Lobato. Grammar variational autoencoder. In *International conference on machine learning*, pages 1945–1954. PMLR, 2017.

[117] Adam Roberts, Jesse Engel, and Douglas Eck. Hierarchical variational autoencoders for music. In *NIPS Workshop on Machine Learning for Creativity and Design*, volume 3, 2017.

[118] Zichao Yang, Zhiting Hu, Ruslan Salakhutdinov, and Taylor Berg-Kirkpatrick. Improved variational autoencoders for text modeling using dilated convolutions. In *International conference on machine learning*, pages 3881–3890. PMLR, 2017.

[119] Aaron Van Den Oord, Oriol Vinyals, et al. Neural discrete representation learning. *Advances in neural information processing systems*, 30, 2017.

[120] Stanislau Semeniuta, Aliaksei Severyn, and Erhardt Barth. A hybrid convolutional variational autoencoder for text generation. *arXiv preprint arXiv:1702.02390*, 2017.

[121] David Ha and Douglas Eck. A neural representation of sketch drawings. *arXiv preprint arXiv:1704.03477*, 2017.

[122] Alexey Tikhonov, Ivan P Yamshchikov, et al. Music generation with variational recurrent autoencoder supported by history. *arXiv preprint arXiv:1705.05458*, 2017.

[123] Jesse Engel, Matthew Hoffman, and Adam Roberts. Latent constraints: Learning to generate conditionally from unconditional generative models. *arXiv preprint arXiv:1711.05772*, 2017.

[124] Peter Bjørn Jørgensen, Murat Mesta, Suranjan Shil, Juan Maria García Lastra, Karsten Wedel Jacobsen, Kristian Sommer Thygesen, and Mikkel N Schmidt. Machine learning-based screening of complex molecules for polymer solar cells. *The Journal of chemical physics*, 148(24):241735, 2018.

[125] Martin Simonovsky and Nikos Komodakis. Graphvae: Towards generation of small graphs using variational autoencoders. In *International conference on artificial neural networks*, pages 412–422. Springer, 2018.

[126] Adam Roberts, Jesse Engel, Colin Raffel, Curtis Hawthorne, and Douglas Eck. A hierarchical latent vector model for learning long-term structure in music. In *International conference on machine learning*, pages 4364–4373. PMLR, 2018.

[127] Wengong Jin, Regina Barzilay, and Tommi Jaakkola. Junction tree variational autoencoder for molecular graph generation. In *International conference on machine learning*, pages 2323–2332. PMLR, 2018.

[128] Justin Gilmer, Samuel S Schoenholz, Patrick F Riley, Oriol Vinyals, and George E Dahl. Neural message passing for quantum chemistry. In *International conference on machine learning*, pages 1263–1272. PMLR, 2017.

[129] Qi Liu, Miltiadis Allamanis, Marc Brockschmidt, and Alexander Gaunt. Constrained graph variational autoencoders for molecule design. *Advances in neural information processing systems*, 31, 2018.

[130] Hanjun Dai, Yingtao Tian, Bo Dai, Steven Skiena, and Le Song. Syntax-directed variational autoencoder for structured data. *arXiv preprint arXiv:1802.08786*, 2018.

[131] Qingrong Chen, Chong Xiang, Minhui Xue, Bo Li, Nikita Borisov, Dali Kaarfar, and Haojin Zhu. Differentially private data generative models. *arXiv preprint arXiv:1812.02274*, 2018.

[132] Martin Abadi, Andy Chu, Ian Goodfellow, H Brendan McMahan, Ilya Mironov, Kunal Talwar, and Li Zhang. Deep learning with differential privacy. In *Proceedings of the 2016 ACM SIGSAC conference on computer and communications security*, pages 308–318, 2016.

[133] Huaibo Huang, zhihang li, Ran He, Zhenan Sun, and Tieniu Tan. Introvae: Introspective variational autoencoders for photographic image synthesis. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc., 2018.

[134] Dinghan Shen, Asli Celikyilmaz, Yizhe Zhang, Liqun Chen, Xin Wang, Jianfeng Gao, and Lawrence Carin. Towards generating long and coherent text with multi-level latent variable models. *arXiv preprint arXiv:1902.00154*, 2019.

[135] Wenlin Wang, Zhe Gan, Hongteng Xu, Ruiyi Zhang, Guoyin Wang, Dinghan Shen, Changyou Chen, and Lawrence Carin. Topic-guided variational autoencoders for text generation. *arXiv preprint arXiv:1903.07137*, 2019.

[136] Xavier Bresson and Thomas Laurent. A two-step graph convolutional decoder for molecule generation. *arXiv preprint arXiv:1906.03412*, 2019.

[137] Bidisha Samanta, Abir De, Gourhari Jana, Vicenç Gómez, Pratim Kumar Chattaraj, Niloy Ganguly, and Manuel Gomez-Rodriguez. Nevae: A deep generative model for molecular graphs. *Journal of machine learning research. 2020 Apr; 21 (114): 1-33*, 2020.

[138] Xiaojie Guo, Liang Zhao, Zhao Qin, Lingfei Wu, Amarda Shehu, and Yanfang Ye. Interpretable deep graph generation with node-edge co-disentanglement. In *Proceedings of the 26th ACM SIGKDD international conference on knowledge discovery & data mining*, pages 1697–1707, 2020.

[139] Alfredo Nazabal, Pablo M Olmos, Zoubin Ghahramani, and Isabel Valera. Handling incomplete heterogeneous data using vaes. *Pattern Recognition*, 107:107501, 2020.

[140] Daniel Flam-Shepherd, Tony Wu, and Alan Aspuru-Guzik. Graph deconvolutional generation. *arXiv preprint arXiv:2002.07087*, 2020.

[141] Marco Podda, Davide Bacciu, and Alessio Micheli. A deep generative model for fragment-based molecule generation. In *International Conference on Artificial Intelligence and Statistics*, pages 2240–2250. PMLR, 2020.

[142] Danilo Jimenez Rezende, Shakir Mohamed, and Daan Wierstra. Stochastic backpropagation and approximate inference in deep generative models. In *International conference on machine learning*, pages 1278–1286. PMLR, 2014.

[143] Vincent Michalski, Roland Memisevic, and Kishore Konda. Modeling deep temporal dependencies with recurrent grammar cells"". *Advances in neural information processing systems*, 27, 2014.

[144] Roland Memisevic. Gradient-based learning of higher-order image features. In *2011 International Conference on Computer Vision*, pages 1591–1598. IEEE, 2011.

[145] Mathieu Germain, Karol Gregor, Iain Murray, and Hugo Larochelle. Made: Masked autoencoder for distribution estimation. In *International conference on machine learning*, pages 881–889. PMLR, 2015.

[146] Hugo Larochelle and Iain Murray. The neural autoregressive distribution estimator. In *Proceedings of the fourteenth international conference on artificial intelligence and statistics*, pages 29–37. JMLR Workshop and Conference Proceedings, 2011.

[147] Benigno Uria, Iain Murray, and Hugo Larochelle. Rnade: The real-valued neural autoregressive density-estimator. *Advances in Neural Information Processing Systems*, 26, 2013.

[148] Benigno Uria, Iain Murray, and Hugo Larochelle. A deep and tractable density estimator. In *International Conference on Machine Learning*, pages 467–475. PMLR, 2014.

[149] Tapani Raiko, Yao Li, Kyunghyun Cho, and Yoshua Bengio. Iterative neural autoregressive distribution estimator nade-k. *Advances in neural information processing systems*, 27, 2014.

[150] Benigno Uria, Marc-Alexandre Côté, Karol Gregor, Iain Murray, and Hugo Larochelle. Neural autoregressive distribution estimation. *The Journal of Machine Learning Research*, 17(1):7184–7220, 2016.

[151] Francesco Tonolini, Bjørn Sand Jensen, and Roderick Murray-Smith. Variational sparse coding. In *Uncertainty in Artificial Intelligence*, pages 690–700. PMLR, 2020.

[152] Bruno A Olshausen and David J Field. Sparse coding of sensory inputs. *Current Opinion in Neurobiology*, 14(4):481–487, 2004.

[153] Zhaowen Wang, Ding Liu, Jianchao Yang, Wei Han, and Thomas Huang. Deep networks for image super-resolution with sparse prior. In *Proceedings of the IEEE international conference on computer vision*, pages 370–378, 2015.

[154] Zachary Lipton. A critical review of recurrent neural networks for sequence learning. 05 2015.

[155] Haonan Yu, Jiang Wang, Zhiheng Huang, Yi Yang, and Wei Xu. Video paragraph captioning using hierarchical recurrent neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4584–4593, 2016.

[156] Nicolas Boulanger-Lewandowski, Yoshua Bengio, and Pascal Vincent. Modeling temporal dependencies in high-dimensional sequences: Application to polyphonic music generation and transcription. *arXiv preprint arXiv:1206.6392*, 2012.

[157] Alex Graves. Generating sequences with recurrent neural networks. *arXiv preprint arXiv:1308.0850*, 2013.

[158] Oriol Vinyals, Alexander Toshev, Samy Bengio, and Dumitru Erhan. Show and tell: A neural image caption generator. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3156–3164, 2015.

[159] Jeffrey Donahue, Lisa Anne Hendricks, Sergio Guadarrama, Marcus Rohrbach, Subhashini Venugopalan, Kate Saenko, and Trevor Darrell. Long-term recurrent convolutional networks for visual recognition and description. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2625–2634, 2015.

[160] Nitish Srivastava, Elman Mansimov, and Ruslan Salakhudinov. Unsupervised learning of video representations using lstms. In *International conference on machine learning*, pages 843–852. PMLR, 2015.

[161] Elman Mansimov, Emilio Parisotto, Jimmy Lei Ba, and Ruslan Salakhutdinov. Generating images from captions with attention. *arXiv preprint arXiv:1511.02793*, 2015.

[162] Natasha Jaques, Shixiang Gu, Richard E Turner, and Douglas Eck. Tuning recurrent neural networks with re-inforcement learning. *arXiv preprint arXiv:1611.02796*, 2016.

[163] Elliot Waite et al. Generating long-term structure in songs and stories. *Web blog post. Magenta*, 15(4), 2016. Accessed: 2022-05-13.

[164] Gaëtan Hadjeres and Frank Nielsen. Interactive music generation with positional constraints using anticipation-rnns. *arXiv preprint arXiv:1709.06404*, 2017.

[165] Sageev Oore, Ian Simon, Sander Dieleman, Douglas Eck, and Karen Simonyan. This time with feeling: Learning expressive musical performance. *Neural Computing and Applications*, 32(4):955–967, 2020.

[166] Yoshua Bengio, Nicolas Boulanger-Lewandowski, and Razvan Pascanu. Advances in optimizing recurrent networks. *arXiv preprint arXiv:1212.0901*, 2012.

[167] Tomáš Mikolov, Ilya Sutskever, Anoop Deoras, Hai-Son Le, Stefan Kombrink, and Jan Cernocky. Subword language modeling with neural networks. *preprint (http://www. fit. vutbr. cz/imikolov/rnnlm/char. pdf)*, 8(67), 2012.

[168] Andrés E Coca, Débora C Corrêa, and Liang Zhao. Computer-aided music composition with lstm neural network and chaotic inspiration. In *IJCNN*, pages 1–7, 2013.

[169] Justin Bayer, Christian Osendorfer, Daniela Korhammer, Nutan Chen, Sebastian Urban, and Patrick van der Smagt. On fast dropout and its applicability to recurrent networks. *arXiv preprint arXiv:1311.0701*, 2013.

[170] Razvan Pascanu, Caglar Gulcehre, Kyunghyun Cho, and Yoshua Bengio. How to construct deep recurrent neural networks. *arXiv preprint arXiv:1312.6026*, 2013.

[171] Subhashini Venugopalan, Huijuan Xu, Jeff Donahue, Marcus Rohrbach, Raymond Mooney, and Kate Saenko. Translating videos to natural language using deep recurrent neural networks. *arXiv preprint arXiv:1412.4729*, 2014.

[172] Xingxing Zhang and Mirella Lapata. Chinese poetry generation with recurrent neural networks. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 670–680, 2014.

[173] Justin Bayer and Christian Osendorfer. Learning stochastic recurrent networks. *arXiv preprint arXiv:1411.7610*, 2014.

[174] Jan Koutnik, Klaus Greff, Faustino Gomez, and Juergen Schmidhuber. A clockwork rnn. In *International Conference on Machine Learning*, pages 1863–1871. PMLR, 2014.

[175] I Liu, Bhiksha Ramakrishnan, et al. Bach in 2014: Music composition with recurrent neural network. *arXiv preprint arXiv:1412.3191*, 2014.

[176] Martin Riedmiller and I Rprop. Rprop-description and implementation details. 1994.

[177] Kratarth Goel, Raunaq Vohra, and Jajati Keshari Sahoo. Polyphonic music generation by modeling temporal dependencies using a rnn-dbn. In *International Conference on Artificial Neural Networks*, pages 217–224. Springer, 2014.

[178] Junhua Mao, Wei Xu, Yi Yang, Jiang Wang, Zhiheng Huang, and Alan Yuille. Deep captioning with multimodal recurrent neural networks (m-rnn). *arXiv preprint arXiv:1412.6632*, 2014.

[179] Lucas Theis and Matthias Bethge. Generative image modeling using spatial lstms. *Advances in neural information processing systems*, 28, 2015.

[180] Alex Graves and Jürgen Schmidhuber. Offline handwriting recognition with multidimensional recurrent neural networks. *Advances in neural information processing systems*, 21, 2008.

[181] Lucas Theis, Reshad Hosseini, and Matthias Bethge. Mixtures of conditional gaussian scale mixtures applied to multiscale image representations. 2012.

[182] Mathias Berglund, Tapani Raiko, Mikko Honkala, Leo Kärkkäinen, Akos Vetek, and Juha T Karhunen. Bidirectional recurrent neural networks as generative models. *Advances in neural information processing systems*, 28, 2015.

[183] Junyoung Chung, Kyle Kastner, Laurent Dinh, Kratarth Goel, Aaron C Courville, and Yoshua Bengio. A recurrent latent variable model for sequential data. *Advances in neural information processing systems*, 28, 2015.

[184] Samy Bengio, Oriol Vinyals, Navdeep Jaitly, and Noam Shazeer. Scheduled sampling for sequence prediction with recurrent neural networks. *Advances in neural information processing systems*, 28, 2015.

[185] Xinlei Chen and C Lawrence Zitnick. Mind's eye: A recurrent visual representation for image caption generation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2422–2431, 2015.

[186] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.

[187] Andrej Karpathy and Li Fei-Fei. Deep visual-semantic alignments for generating image descriptions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3128–3137, 2015.

[188] Raunaq Vohra, Kratarth Goel, and Jajati Keshari Sahoo. Modeling temporal dependencies in data using a dbn-lstm. In *2015 IEEE International Conference on Data Science and Advanced Analytics (DSAA)*, pages 1–4. IEEE, 2015.

[189] Qi Lyu, Zhiyong Wu, Jun Zhu, and Helen Meng. Modelling high-dimensional sequences with lstm-rtrbm: Application to polyphonic music generation. In *Twenty-Fourth International Joint Conference on Artificial Intelligence*, 2015.

[190] Kelvin Xu, Jimmy Ba, Ryan Kiros, Kyunghyun Cho, Aaron Courville, Ruslan Salakhudinov, Rich Zemel, and Yoshua Bengio. Show, attend and tell: Neural image caption generation with visual attention. In *International conference on machine learning*, pages 2048–2057. PMLR, 2015.

[191] Subhashini Venugopalan, Marcus Rohrbach, Jeffrey Donahue, Raymond Mooney, Trevor Darrell, and Kate Saenko. Sequence to sequence-video to text. In *Proceedings of the IEEE international conference on computer vision*, pages 4534–4542, 2015.

[192] Li Yao, Atousa Torabi, Kyunghyun Cho, Nicolas Ballas, Christopher Pal, Hugo Larochelle, and Aaron Courville. Describing videos by exploiting temporal structure. In *Proceedings of the IEEE international conference on computer vision*, pages 4507–4515, 2015.

[193] Marco Fraccaro, Søren Kaae Sønderby, Ulrich Paquet, and Ole Winther. Sequential neural models with stochastic layers. *Advances in neural information processing systems*, 29, 2016.

[194] Florian Colombo, Samuel P Muscinelli, Alexander Seeholzer, Johanni Brea, and Wulfram Gerstner. Algorithmic composition of melodies with deep recurrent neural networks. *arXiv preprint arXiv:1606.07251*, 2016.

[195] Zheng Sun, Jiaqi Liu, Zewang Zhang, Jingwen Chen, Zhao Huo, Ching Hua Lee, and Xiao Zhang. Composing music with grammar argumented neural networks and note-level encoding. *arXiv preprint arXiv:1611.05416*, 2016.

[196] Dzmitry Bahdanau, Philemon Brakel, Kelvin Xu, Anirudh Goyal, Ryan Lowe, Joelle Pineau, Aaron Courville, and Yoshua Bengio. An actor-critic algorithm for sequence prediction. *arXiv preprint arXiv:1607.07086*, 2016.

[197] Zhilin Yang, Ye Yuan, Yuexin Wu, William W Cohen, and Russ R Salakhutdinov. Review networks for caption generation. *Advances in neural information processing systems*, 29, 2016.

[198] Quanzeng You, Hailin Jin, Zhaowen Wang, Chen Fang, and Jiebo Luo. Image captioning with semantic attention. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4651–4659, 2016.

[199] Keunwoo Choi, George Fazekas, and Mark Sandler. Text-based lstm networks for automatic music composition. *arXiv preprint arXiv:1604.05358*, 2016.

[200] Hang Chu, Raquel Urtasun, and Sanja Fidler. Song from pi: A musically plausible network for pop music generation. *arXiv preprint arXiv:1611.03477*, 2016.

[201] Bob L Sturm, Joao Felipe Santos, Oded Ben-Tal, and Iryna Korshunova. Music transcription modelling and composition using deep learning. *arXiv preprint arXiv:1604.08723*, 2016.

[202] Feynman Liang. Bachbot: Automatic composition in the style of bach chorales. *University of Cambridge*, 8:19–48, 2016.

[203] Allen Huang and Raymond Wu. Deep learning for music. *arXiv preprint arXiv:1606.04930*, 2016.

[204] Soroush Mehri, Kundan Kumar, Ishaan Gulrajani, Rithesh Kumar, Shubham Jain, Jose Sotelo, Aaron Courville, and Yoshua Bengio. Samplernn: An unconditional end-to-end neural audio generation model. *arXiv preprint arXiv:1612.07837*, 2016.

[205] Iulian Serban, Alessandro Sordoni, Yoshua Bengio, Aaron Courville, and Joelle Pineau. Building end-to-end dialogue systems using generative hierarchical neural network models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 30, 2016.

[206] Alex M Lamb, Anirudh Goyal ALIAS PARTH GOYAL, Ying Zhang, Saizheng Zhang, Aaron C Courville, and Yoshua Bengio. Professor forcing: A new algorithm for training recurrent networks. *Advances in neural information processing systems*, 29, 2016.

[207] Jose Sotelo, Soroush Mehri, Kundan Kumar, Joao Felipe Santos, Kyle Kastner, Aaron Courville, and Yoshua Bengio. Char2wav: End-to-end speech synthesis. 2017.

[208] Julian Georg Zilly, Rupesh Kumar Srivastava, Jan Koutník, and Jürgen Schmidhuber. Recurrent highway networks. In *International conference on machine learning*, pages 4189–4198. PMLR, 2017.

[209] Rupesh K Srivastava, Klaus Greff, and Jürgen Schmidhuber. Training very deep networks. *Advances in neural information processing systems*, 28, 2015.

[210] Florian Colombo, Alexander Seeholzer, and Wulfram Gerstner. Deep artificial composer: A creative neural network model for automated melody generation. In *International Conference on Evolutionary and Biologically Inspired Music and Art*, pages 81–96. Springer, 2017.

[211] Chenxi Liu, Junhua Mao, Fei Sha, and Alan Yuille. Attention correctness in neural image captioning. In *Thirty-first AAAI conference on artificial intelligence*, 2017.

[212] Zhe Gan, Chuang Gan, Xiaodong He, Yunchen Pu, Kenneth Tran, Jianfeng Gao, Lawrence Carin, and Li Deng. Semantic compositional networks for visual captioning. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5630–5639, 2017.

[213] Gaëtan Hadjeres, François Pachet, and Frank Nielsen. Deepbach: a steerable model for bach chorales generation. In *International Conference on Machine Learning*, pages 1362–1371. PMLR, 2017.

[214] Dimos Makris, Maximos Kaliakatsos-Papakostas, Ioannis Karydis, and Katia Lida Kermanidis. Combining lstm and feed forward neural networks for conditional rhythm composition. In *International conference on engineering applications of neural networks*, pages 570–582. Springer, 2017.

[215] Natasha Jaques, Shixiang Gu, Dzmitry Bahdanau, José Miguel Hernández-Lobato, Richard E Turner, and Douglas Eck. Sequence tutor: Conservative fine-tuning of sequence generation models with kl-control. In *International Conference on Machine Learning*, pages 1645–1654. PMLR, 2017.

[216] Hyungui Lim, Seungyeon Rhyu, and Kyogu Lee. Chord generation from symbolic melody using blstm networks. *arXiv preprint arXiv:1712.01011*, 2017.

[217] Iulian Serban, Alessandro Sordoni, Ryan Lowe, Laurent Charlin, Joelle Pineau, Aaron Courville, and Yoshua Bengio. A hierarchical latent variable encoder-decoder model for generating dialogues. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 31, 2017.

[218] Patrick Hutchings and Jon McCormack. Using autonomous agents to improvise music compositions in real-time. In *International conference on evolutionary and biologically inspired music and art*, pages 114–127. Springer, 2017.

[219] Ian Simon and Sageev Oore. Performance rnn: Generating music with expressive timing and dynamics. https://magenta.tensorflow.org/performance-rnn, 2017.

[220] Daniel D Johnson. Generating polyphonic music using tied parallel networks. In *International conference on evolutionary and biologically inspired music and art*, pages 128–143. Springer, 2017.

[221] Rachel Manzelli, Vijay Thakkar, Ali Siahkamari, and Brian Kulis. Conditioning deep generative raw audio models for structured automatic music. *arXiv preprint arXiv:1806.09905*, 2018.

[222] Yujia Li, Oriol Vinyals, Chris Dyer, Razvan Pascanu, and Peter Battaglia. Learning deep generative models of graphs. *arXiv preprint arXiv:1803.03324*, 2018.

[223] Nal Kalchbrenner, Erich Elsen, Karen Simonyan, Seb Noury, Norman Casagrande, Edward Lockhart, Florian Stimberg, Aaron Oord, Sander Dieleman, and Koray Kavukcuoglu. Efficient neural audio synthesis. In *International Conference on Machine Learning*, pages 2410–2419. PMLR, 2018.

[224] Adam Santoro, Ryan Faulkner, David Raposo, Jack Rae, Mike Chrzanowski, Theophane Weber, Daan Wierstra, Oriol Vinyals, Razvan Pascanu, and Timothy Lillicrap. Relational recurrent neural networks. *Advances in neural information processing systems*, 31, 2018.

[225] Huanru Henry Mao, Taylor Shin, and Garrison Cottrell. Deepj: Style-specific music generation. In *2018 IEEE 12th International Conference on Semantic Computing (ICSC)*, pages 377–382. IEEE, 2018.

[226] Jiaxuan You, Rex Ying, Xiang Ren, William Hamilton, and Jure Leskovec. Graphrnn: Generating realistic graphs with deep auto-regressive models. In *International conference on machine learning*, pages 5708–5717. PMLR, 2018.

[227] Jonathan Shen, Ruoming Pang, Ron J Weiss, Mike Schuster, Navdeep Jaitly, Zongheng Yang, Zhifeng Chen, Yu Zhang, Yuxuan Wang, Rj Skerrv-Ryan, et al. Natural tts synthesis by conditioning wavenet on mel spectrogram predictions. In *2018 IEEE international conference on acoustics, speech and signal processing (ICASSP)*, pages 4779–4783. IEEE, 2018.

[228] Aaron van den Oord, Sander Dieleman, Heiga Zen, Karen Simonyan, Oriol Vinyals, Alex Graves, Nal Kalchbrenner, Andrew Senior, and Koray Kavukcuoglu. Wavenet: A generative model for raw audio. *arXiv preprint arXiv:1609.03499*, 2016.

[229] Renjie Liao, Yujia Li, Yang Song, Shenlong Wang, Will Hamilton, David K Duvenaud, Raquel Urtasun, and Richard Zemel. Efficient graph generation with graph recurrent attention networks. *Advances in neural information processing systems*, 32, 2019.

[230] Davide Bacciu, Alessio Micheli, and Marco Podda. Graph generation by sequential edge prediction. In *27th European Symposium on Artificial Neural Networks, Computational Intelligence and Machine Learning, ESANN 2019*, pages 95–100. ESANN (i6doc. com), 2019.

[231] Mariya Popova, Mykhailo Shvets, Junier Oliva, and Olexandr Isayev. Molecularrnn: Generating realistic molecular graphs with optimized properties. *arXiv preprint arXiv:1905.13372*, 2019.

[232] Mahdi Khodayar, Jianhui Wang, and Zhaoyu Wang. Deep generative graph distribution learning for synthetic power grids. *arXiv preprint arXiv:1901.09674*, 2019.

[233] Andres Hernandez-Matamoros, Hamido Fujita, and Hector Perez-Meana. A novel approach to create synthetic biomedical signals using birnn. *Information Sciences*, 541:218–241, 2020.

[234] Nikhil Goyal, Harsh Vardhan Jain, and Sayan Ranu. Graphgen: a scalable approach to domain-agnostic labeled graph generation. In *Proceedings of The Web Conference 2020*, pages 1253–1263, 2020.

[235] Nicola Privato, Omar Rampado, and Alberto Novello. A creative tool for the musician combining lstm and markov chains in max/msp. In *International Conference on Computational Intelligence in Music, Sound, Art and Design (Part of EvoStar)*, pages 228–242. Springer, 2022.

[236] Saad Albawi, Tareq Abed Mohammed, and Saad Al-Zawi. Understanding of a convolutional neural network. In *2017 international conference on engineering and technology (ICET)*, pages 1–6. Ieee, 2017.

[237] William Lotter, Gabriel Kreiman, and David Cox. Unsupervised learning of visual structure using predictive generative networks. *arXiv preprint arXiv:1511.06380*, 2015.

[238] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. *Advances in neural information processing systems*, 27, 2014.

[239] Joan Bruna, Pablo Sprechmann, and Yann LeCun. Super-resolution with deep convolutional sufficient statistics. *arXiv preprint arXiv:1511.05666*, 2015.

[240] Nal Kalchbrenner, Aäron Oord, Karen Simonyan, Ivo Danihelka, Oriol Vinyals, Alex Graves, and Koray Kavukcuoglu. Video pixel networks. In *International Conference on Machine Learning*, pages 1771–1779. PMLR, 2017.

[241] Scott Reed, Aäron van den Oord, Nal Kalchbrenner, Victor Bapst, Matt Botvinick, and Nando De Freitas. Generating interpretable images with controllable structure. 2016.

[242] Leon A Gatys, Alexander S Ecker, and Matthias Bethge. Image style transfer using convolutional neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2414–2423, 2016.

[243] Jiwon Kim, Jung Kwon Lee, and Kyoung Mu Lee. Deeply-recursive convolutional network for image super-resolution. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1637–1645, 2016.

[244] Chao Dong, Chen Change Loy, Kaiming He, and Xiaoou Tang. Learning a deep convolutional network for image super-resolution. In *European conference on computer vision*, pages 184–199. Springer, 2014.

[245] Chao Dong, Chen Change Loy, Kaiming He, and Xiaoou Tang. Image super-resolution using deep convolutional networks. *IEEE transactions on pattern analysis and machine intelligence*, 38(2):295–307, 2015.

[246] Tim Salimans, Andrej Karpathy, Xi Chen, and Diederik P Kingma. Pixelcnn++: Improving the pixelcnn with discretized logistic mixture likelihood and other modifications. *arXiv preprint arXiv:1701.05517*, 2017.

[247] Alexander Mordvintsev, Christopher Olah, and Mike Tyka. Deepdream-a code example for visualizing neural networks. *Google Research*, 2(5), 2015.

[248] Alexander Mordvintsev, Christopher Olah, and Mike Tyka. Inceptionism: Going deeper into neural networks. 2015.

[249] Aaron Van den Oord, Nal Kalchbrenner, Lasse Espeholt, Oriol Vinyals, Alex Graves, et al. Conditional image generation with pixelcnn decoders. *Advances in neural information processing systems*, 29, 2016.

[250] Justin Johnson, Alexandre Alahi, and Li Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In *European conference on computer vision*, pages 694–711. Springer, 2016.

[251] Wenzhe Shi, Jose Caballero, Ferenc Huszár, Johannes Totz, Andrew P Aitken, Rob Bishop, Daniel Rueckert, and Zehan Wang. Real-time single image and video super-resolution using an efficient sub-pixel convolutional neural network. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1874–1883, 2016.

[252] Wei Ping, Kainan Peng, Andrew Gibiansky, Sercan Ömer Arik, Ajay Kannan, Sharan Narang, Jonathan Raiman, and John Miller. Deep voice 3: 2000-speaker neural text-to-speech. 2017.

[253] Xi Chen, Nikhil Mishra, Mostafa Rohaninejad, and Pieter Abbeel. Pixelsnail: An improved autoregressive generative model. *arXiv preprint arXiv:1712.09763*, 2017.

[254] Nikhil Mishra, Mostafa Rohaninejad, Xi Chen, and Pieter Abbeel. A simple neural attentive meta-learner. *arXiv preprint arXiv:1707.03141*, 2017.

[255] Jacob Menick and Nal Kalchbrenner. Generating high fidelity images with subscale pixel networks and multidimensional upscaling. *arXiv preprint arXiv:1812.01608*, 2018.

[256] Yifan Wang, Lijun Wang, Hongyu Wang, and Peihua Li. End-to-end image super-resolution via deep and shallow convolutional networks. *IEEE Access*, 7:31959–31970, 2019.

[257] Peter J Liu, Mohammad Saleh, Etienne Pot, Ben Goodrich, Ryan Sepassi, Lukasz Kaiser, and Noam Shazeer. Generating wikipedia by summarizing long sequences. *arXiv preprint arXiv:1801.10198*, 2018.

[258] Niki Parmar, Ashish Vaswani, Jakob Uszkoreit, Lukasz Kaiser, Noam Shazeer, Alexander Ku, and Dustin Tran. Image transformer. In *International Conference on Machine Learning*, pages 4055–4064. PMLR, 2018.

[259] Cheng-Zhi Anna Huang, Ashish Vaswani, Jakob Uszkoreit, Noam Shazeer, Ian Simon, Curtis Hawthorne, Andrew M Dai, Matthew D Hoffman, Monica Dinculescu, and Douglas Eck. Music transformer. *arXiv preprint arXiv:1809.04281*, 2018.

[260] Peter Shaw, Jakob Uszkoreit, and Ashish Vaswani. Self-attention with relative position representations. *arXiv preprint arXiv:1803.02155*, 2018.

[261] Rewon Child, Scott Gray, Alec Radford, and Ilya Sutskever. Generating long sequences with sparse transformers. *arXiv preprint arXiv:1904.10509*, 2019.

[262] Chen Sun, Austin Myers, Carl Vondrick, Kevin Murphy, and Cordelia Schmid. Videobert: A joint model for video and language representation learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7464–7473, 2019.

[263] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.

[264] Chia-Cheng Liu, Harris Chan, Kevin Luk, and AI Borealis. Auto-regressive graph generation modeling with improved evaluation methods. In *33rd Conference on Neural Information Processing Systems. Vancouver, Canada*, 2019.

[265] Jie Feng, Xueliang Feng, Jiantong Chen, Xianghai Cao, Xiangrong Zhang, Licheng Jiao, and Tao Yu. Generative adversarial networks based on collaborative learning and attention mechanism for hyperspectral image classification. *Remote Sensing*, 12:1149, 04 2020.

[266] Olivier Breuleux, Yoshua Bengio, and Pascal Vincent. Quickly generating representative samples from an rbm-derived process. *Neural computation*, 23(8):2058–2073, 2011.

[267] Alec Radford, Luke Metz, and Soumith Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. *arXiv preprint arXiv:1511.06434*, 2015.

[268] Maayan Frid-Adar, Eyal Klang, Michal Amitai, Jacob Goldberger, and Hayit Greenspan. Synthetic data augmentation using gan for improved liver lesion classification. In *2018 IEEE 15th international symposium on biomedical imaging (ISBI 2018)*, pages 289–293. IEEE, 2018.

[269] Luke Metz, Ben Poole, David Pfau, and Jascha Sohl-Dickstein. Unrolled generative adversarial networks. *arXiv preprint arXiv:1611.02163*, 2016.

[270] Ahmed Elgammal, Bingchen Liu, Mohamed Elhoseiny, and Marian Mazzone. Can: Creative adversarial networks, generating" art" by learning about styles and deviating from style norms. *arXiv preprint arXiv:1706.07068*, 2017.

[271] Chris Donahue, Julian McAuley, and Miller Puckette. Adversarial audio synthesis. *arXiv preprint arXiv:1802.04208*, 2018.

[272] Ishan Durugkar, Ian Gemp, and Sridhar Mahadevan. Generative multi-adversarial networks. *arXiv preprint arXiv:1611.01673*, 2016.

[273] Tong Che, Yanran Li, Athul Paul Jacob, Yoshua Bengio, and Wenjie Li. Mode regularized generative adversarial networks. *arXiv preprint arXiv:1612.02136*, 2016.

[274] Jiajun Wu, Chengkai Zhang, Tianfan Xue, Bill Freeman, and Josh Tenenbaum. Learning a probabilistic latent space of object shapes via 3d generative-adversarial modeling. *Advances in neural information processing systems*, 29, 2016.

[275] Yizhe Zhang, Zhe Gan, and Lawrence Carin. Generating text via adversarial training. In *NIPS workshop on Adversarial Training*, volume 21, pages 21–32, 2016.

[276] Yizhe Zhang, Zhe Gan, Kai Fan, Zhi Chen, Ricardo Henao, Dinghan Shen, and Lawrence Carin. Adversarial feature matching for text generation. In *International Conference on Machine Learning*, pages 4006–4015. PMLR, 2017.

[277] Ming-Yu Liu and Oncel Tuzel. Coupled generative adversarial networks. *Advances in neural information processing systems*, 29, 2016.

[278] Junbo Zhao, Michael Mathieu, and Yann LeCun. Energy-based generative adversarial network. *arXiv preprint arXiv:1609.03126*, 2016.

[279] Yann LeCun, Sumit Chopra, Raia Hadsell, M Ranzato, and F Huang. A tutorial on energy-based learning. *Predicting structured data*, 1(0), 2006.

[280] Olof Mogren. C-rnn-gan: Continuous recurrent neural networks with adversarial training. *arXiv preprint arXiv:1611.09904*, 2016.

[281] Andrew Brock, Theodore Lim, James M Ritchie, and Nick Weston. Neural photo editing with introspective adversarial networks. *arXiv preprint arXiv:1609.07093*, 2016.

[282] Sebastian Nowozin, Botond Cseke, and Ryota Tomioka. f-gan: Training generative neural samplers using variational divergence minimization. *Advances in neural information processing systems*, 29, 2016.

[283] Daniel Jiwoong Im, Chris Dongjoo Kim, Hui Jiang, and Roland Memisevic. Generating images with recurrent adversarial networks. *arXiv preprint arXiv:1602.05110*, 2016.

[284] Tim Salimans, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, Xi Chen, and Xi Chen. Improved techniques for training gans. In D. Lee, M. Sugiyama, U. Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 29. Curran Associates, Inc., 2016.

[285] Martin Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein generative adversarial networks. In *International conference on machine learning*, pages 214–223. PMLR, 2017.

[286] Ishaan Gulrajani, Faruk Ahmed, Martin Arjovsky, Vincent Dumoulin, and Aaron C Courville. Improved training of wasserstein gans. *Advances in neural information processing systems*, 30, 2017.

[287] David Berthelot, Thomas Schumm, and Luke Metz. Began: Boundary equilibrium generative adversarial networks. *arXiv preprint arXiv:1703.10717*, 2017.

[288] 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.

[289] Kevin Lin, Dianqi Li, Xiaodong He, Zhengyou Zhang, and Ming-Ting Sun. Adversarial ranking for language generation. *Advances in neural information processing systems*, 30, 2017.

[290] Paulina Grnarova, Kfir Y Levy, Aurelien Lucchi, Thomas Hofmann, and Andreas Krause. An online learning approach to generative adversarial networks. *arXiv preprint arXiv:1706.03269*, 2017.

[291] Yoon Kim, Kelly Zhang, Alexander M Rush, Yann LeCun, et al. Adversarially regularized autoencoders for generating discrete structures. *arXiv preprint arXiv:1706.04223*, 2:12, 2017.

[292] Youssef Mroueh, Tom Sercu, and Vaibhava Goel. Mcgan: Mean and covariance feature matching gan. In *International conference on machine learning*, pages 2527–2535. PMLR, 2017.

[293] Youssef Mroueh and Tom Sercu. Fisher gan. *Advances in Neural Information Processing Systems*, 30, 2017.

[294] Mihaela Rosca, Balaji Lakshminarayanan, David Warde-Farley, and Shakir Mohamed. Variational approaches for auto-encoding generative adversarial networks. *arXiv preprint arXiv:1706.04987*, 2017.

[295] Lantao Yu, Weinan Zhang, Jun Wang, and Yong Yu. Seqgan: Sequence generative adversarial nets with policy gradient. In *Proceedings of the AAAI conference on artificial intelligence*, volume 31, 2017.

[296] Sang-gil Lee, Uiwon Hwang, Seonwoo Min, and Sungroh Yoon. A seqgan for polyphonic music generation. 2017.

[297] Masaki Saito, Eiichi Matsumoto, and Shunta Saito. Temporal generative adversarial nets with singular value clipping. In *Proceedings of the IEEE international conference on computer vision*, pages 2830–2839, 2017.

[298] Akash Srivastava, Lazar Valkov, Chris Russell, Michael U Gutmann, and Charles Sutton. Veegan: Reducing mode collapse in gans using implicit variational learning. *Advances in neural information processing systems*, 30, 2017.

[299] R Devon Hjelm, Athul Paul Jacob, Tong Che, Adam Trischler, Kyunghyun Cho, and Yoshua Bengio. Boundary-seeking generative adversarial networks. *arXiv preprint arXiv:1702.08431*, 2017.

[300] Edward Choi, Siddharth Biswal, Bradley Malin, Jon Duke, Walter F Stewart, and Jimeng Sun. Generating multi-label discrete patient records using generative adversarial networks. In *Machine learning for healthcare conference*, pages 286–305. PMLR, 2017.

[301] Li-Chia Yang, Szu-Yu Chou, and Yi-Hsuan Yang. Midinet: A convolutional generative adversarial network for symbolic-domain music generation. *arXiv preprint arXiv:1703.10847*, 2017.

[302] Vincent Dumoulin and Francesco Visin. A guide to convolution arithmetic for deep learning. *arXiv preprint arXiv:1603.07285*, 2016.

[303] Tero Karras, Timo Aila, Samuli Laine, and Jaakko Lehtinen. Progressive growing of gans for improved quality, stability, and variation. *arXiv preprint arXiv:1710.10196*, 2017.

[304] Naveen Kodali, Jacob Abernethy, James Hays, and Zsolt Kira. How to train your dragan. *arXiv preprint arXiv:1705.07215*, 2(4), 2017.

[305] Felix Juefei-Xu, Vishnu Naresh Boddeti, and Marios Savvides. Gang of gans: Generative adversarial networks with maximum margin ranking. *arXiv preprint arXiv:1704.04865*, 2017.

[306] Jiaxian Guo, Sidi Lu, Han Cai, Weinan Zhang, Yong Yu, and Jun Wang. Long text generation via adversarial training with leaked information. *arXiv preprint arXiv:1709.08624*, 2017.

[307] Vincent Dumoulin, Ishmael Belghazi, Ben Poole, Olivier Mastropietro, Alex Lamb, Martin Arjovsky, and Aaron Courville. Adversarially learned inference. *arXiv preprint arXiv:1606.00704*, 2016.

[308] Jianwei Yang, Anitha Kannan, Dhruv Batra, and Devi Parikh. Lr-gan: Layered recursive generative adversarial networks for image generation. *arXiv preprint arXiv:1703.01560*, 2017.

[309] David Warde-Farley and Yoshua Bengio. Improving generative adversarial networks with denoising feature matching. In *International Conference on Learning Representations*, 2017.

[310] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017.

[311] Sergey Tulyakov, Ming-Yu Liu, Xiaodong Yang, and Jan Kautz. Mocogan: Decomposing motion and content for video generation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1526–1535, 2018.

[312] 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 AAAI Conference on Artificial Intelligence*, volume 32, 2018.

[313] Ayush Jaiswal, Wael AbdAlmageed, Yue Wu, and Premkumar Natarajan. Capsulegan: Generative adversarial capsule network. In *Proceedings of the European conference on computer vision (ECCV) workshops*, pages 0–0, 2018.

[314] Sara Sabour, Nicholas Frosst, and Geoffrey E Hinton. Dynamic routing between capsules. *Advances in neural information processing systems*, 30, 2017.

[315] Nicola De Cao and Thomas Kipf. Molgan: An implicit generative model for small molecular graphs. *arXiv preprint arXiv:1805.11973*, 2018.

[316] Aleksandar Bojchevski, Oleksandr Shchur, Daniel Zügner, and Stephan Günnemann. Netgan: Generating graphs via random walks. In *International conference on machine learning*, pages 610–619. PMLR, 2018.

[317] Noseong Park, Mahmoud Mohammadi, Kshitij Gorde, Sushil Jajodia, Hongkyu Park, and Youngmin Kim. Data synthesis based on generative adversarial networks. *arXiv preprint arXiv:1806.03384*, 2018.

[318] Shreyas Patel, Ashutosh Kakadiya, Maitrey Mehta, Raj Derasari, Rahul Patel, and Ratnik Gandhi. Correlated discrete data generation using adversarial training. *arXiv preprint arXiv:1804.00925*, 2018.

[319] Lei Xu and Kalyan Veeramachaneni. Synthesizing tabular data using generative adversarial networks. *arXiv preprint arXiv:1811.11264*, 2018.

[320] Xiang Wei, Boqing Gong, Zixia Liu, Wei Lu, and Liqiang Wang. Improving the improved training of wasserstein gans: A consistency term and its dual effect. *arXiv preprint arXiv:1803.01541*, 2018.

[321] Ashish Bora, Eric Price, and Alexandros G Dimakis. Ambientgan: Generative models from lossy measurements. In *International conference on learning representations*, 2018.

[322] Namrata Anand and Possu Huang. Generative modeling for protein structures. *Advances in neural information processing systems*, 31, 2018.

[323] Kun Ouyang, Reza Shokri, David S Rosenblum, and Wenzhuo Yang. A non-parametric generative model for human trajectories. In *IJCAI*, volume 18, pages 3812–3817, 2018.

[324] Quan Hoang, Tu Dinh Nguyen, Trung Le, and Dinh Phung. MGAN: Training generative adversarial nets with multiple generators. In *International Conference on Learning Representations*, 2018.

[325] Weili Nie, Nina Narodytska, and Ankit Patel. Relgan: Relational generative adversarial networks for text generation. In *International conference on learning representations*, 2018.

[326] Ngoc-Trung Tran, Tuan-Anh Bui, and Ngai-Man Cheung. Dist-gan: An improved gan using distance constraints. In *Proceedings of the European Conference on Computer Vision (ECCV)*, September 2018.

[327] Takeru Miyato, Toshiki Kataoka, Masanori Koyama, and Yuichi Yoshida. Spectral normalization for generative adversarial networks. *arXiv preprint arXiv:1802.05957*, 2018.

[328] Hao He, Hao Wang, Guang-He Lee, and Yonglong Tian. Bayesian modelling and monte carlo inference for GAN. In *International Conference on Learning Representations*, 2019.

[329] Aude Genevay, Gabriel Peyre, and Marco Cuturi. Learning generative models with sinkhorn divergences. In Amos Storkey and Fernando Perez-Cruz, editors, *Proceedings of the Twenty-First International Conference on Artificial Intelligence and Statistics*, volume 84 of *Proceedings of Machine Learning Research*, pages 1608–1617. PMLR, 09–11 Apr 2018.

[330] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4401–4410, 2019.

[331] James Jordon, Jinsung Yoon, and Mihaela Van Der Schaar. Pate-gan: Generating synthetic data with differential privacy guarantees. In *International conference on learning representations*, 2018.

[332] Nicolas Papernot, Shuang Song, Ilya Mironov, Ananth Raghunathan, Kunal Talwar, and Úlfar Erlingsson. Scalable private learning with pate. *arXiv preprint arXiv:1802.08908*, 2018.

[333] Xingyuan Chen, Yanzhe Li, Peng Jin, Jiuhua Zhang, Xinyu Dai, Jiajun Chen, and Gang Song. Adversarial sub-sequence for text generation. *arXiv preprint arXiv:1905.12835*, 2019.

[334] Emanuele Ghelfi, Paolo Galeone, Michele De Simoni, and Federico Di Mattia. Adversarial pixel-level generation of semantic images. *arXiv preprint arXiv:1906.12195*, 2019.

[335] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Efros. Image-to-image translation with conditional adversarial networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1125–1134, 2017.

[336] Mrinal Kanti Baowaly, Chia-Ching Lin, Chao-Lin Liu, and Kuan-Ta Chen. Synthesizing electronic health records using improved generative adversarial networks. *Journal of the American Medical Informatics Association*, 26(3):228–241, 2019.

[337] Pooyan Sedigh, Rasoul Sadeghian, and Mehdi Tale Masouleh. Generating synthetic medical images by using gan to improve cnn performance in skin cancer classification. In *2019 7th International Conference on Robotics and Mechatronics (ICRoM)*, pages 497–502. IEEE, 2019.

[338] Xinyu Gong, Shiyu Chang, Yifan Jiang, and Zhangyang Wang. Autogan: Neural architecture search for generative adversarial networks. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2019.

[339] Zinan Lin, Alankar Jain, Chen Wang, Giulia C. Fanti, and Vyas Sekar. Generating high-fidelity, synthetic time series datasets with doppelganger. *arXiv preprint*, September 2019.

[340] Jinsung Yoon, Daniel Jarrett, and Mihaela van der Schaar. Time-series generative adversarial networks. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.

[341] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of stylegan. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8110–8119, 2020.

[342] Andrew Yale, Saloni Dash, Ritik Dutta, Isabelle Guyon, Adrien Pavao, and Kristin P Bennett. Generation and evaluation of privacy preserving synthetic health data. *Neurocomputing*, 416:244–255, 2020.

[343] Saadia Binte Alam, Moazzem Hossain, and Syoji Kobashi. Synthetic brain image generation for adhd prediction based on progressive growing generative adversarial network. In *International Symposium on Affective Science and Engineering ISASE2020*, pages 1–5. Japan Society of Kansei Engineering, 2020.

[344] Jyoti Islam and Yanqing Zhang. Gan-based synthetic brain pet image generation. *Brain informatics*, 7(1):1–12, 2020.

[345] Debapriya Hazra and Yung-Cheol Byun. Synsiggan: Generative adversarial networks for synthetic biomedical signal generation. *Biology*, 9(12):441, 2020.

[346] Chen Gao, Yunpeng Chen, Si Liu, Zhenxiong Tan, and Shuicheng Yan. Adversarialnas: Adversarial neural architecture search for gans. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020.

[347] Tianlin Xu, Li Kevin Wenliang, Michael Munn, and Beatrice Acciaio. Cot-gan: Generating sequential data via causal optimal transport. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 8798–8809. Curran Associates, Inc., 2020.

[348] Shengyu Zhao, Zhijian Liu, Ji Lin, Jun-Yan Zhu, and Song Han. Differentiable augmentation for data-efficient gan training. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 7559–7570. Curran Associates, Inc., 2020.

[349] Christine Dewi, Rung-Ching Chen, Yan-Ting Liu, and Shao-Kuo Tai. Synthetic data generation using dcgan for improved traffic sign recognition. *Neural Computing and Applications*, pages 1–16, 2021.

[350] Sana Imtiaz, Muhammad Arsalan, Vladimir Vlassov, and Ramin Sadre. Synthetic and private smart health care data generation using gans. In *2021 International Conference on Computer Communications and Networks (ICCCN)*, pages 1–7. IEEE, 2021.

[351] Vajira Thambawita, Jonas L Isaksen, Steven A Hicks, Jonas Ghouse, Gustav Ahlberg, Allan Linneberg, Niels Grarup, Christina Ellervik, Morten Salling Olesen, Torben Hansen, et al. Deepfake electrocardiograms using generative adversarial networks are the beginning of the end for privacy issues in medicine. *Scientific reports*, 11(1):1–8, 2021.

[352] Xiaomin Li, Vangelis Metsis, Huangyingrui Wang, and Anne Hee Hiong Ngu. Tts-gan: A transformer-based time-series generative adversarial network. In Martin Michalowski, Syed Sibte Raza Abidi, and Samina Abidi, editors, *Artificial Intelligence in Medicine*, pages 133–143, Cham, 2022. Springer International Publishing.

[353] Mehdi Mirza and Simon Osindero. Conditional generative adversarial nets. *arXiv preprint arXiv:1411.1784*, 2014.

[354] Michael Mathieu, Camille Couprie, and Yann LeCun. Deep multi-scale video prediction beyond mean square error. *arXiv preprint arXiv:1511.05440*, 2015.

[355] Jun-Yan Zhu, Philipp Krähenbühl, Eli Shechtman, and Alexei A Efros. Generative visual manipulation on the natural image manifold. In *European conference on computer vision*, pages 597–613. Springer, 2016.

[356] Scott Reed, Zeynep Akata, Xinchen Yan, Lajanugen Logeswaran, Bernt Schiele, and Honglak Lee. Generative adversarial text to image synthesis. In *International conference on machine learning*, pages 1060–1069. PMLR, 2016.

[357] Scott E Reed, Zeynep Akata, Santosh Mohan, Samuel Tenka, Bernt Schiele, and Honglak Lee. Learning what and where to draw. *Advances in neural information processing systems*, 29, 2016.

[358] Yunjey Choi, Minje Choi, Munyoung Kim, Jung-Woo Ha, Sunghun Kim, and Jaegul Choo. Stargan: Unified generative adversarial networks for multi-domain image-to-image translation. *arXiv preprint arXiv:1711.09020*, 2017.

[359] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *Proceedings of the IEEE international conference on computer vision*, pages 2223–2232, 2017.

[360] Murat Kocaoglu, Christopher Snyder, Alexandros G Dimakis, and Sriram Vishwanath. Causalgan: Learning causal implicit generative models with adversarial training. *arXiv preprint arXiv:1709.02023*, 2017.

[361] Han Zhang, Ian Goodfellow, Dimitris Metaxas, and Augustus Odena. Self-attention generative adversarial networks. In *International conference on machine learning*, pages 7354–7363. PMLR, 2019.

[362] Augustus Odena, Christopher Olah, and Jonathon Shlens. Conditional image synthesis with auxiliary classifier gans. In *International conference on machine learning*, pages 2642–2651. PMLR, 2017.

[363] Arjun Krishna and Klaus Mueller. Medical (ct) image generation with style. In *15th International Meeting on Fully Three-Dimensional Image Reconstruction in Radiology and Nuclear Medicine*, volume 11072, page 1107234. International Society for Optics and Photonics, 2019.

[364] Eloi Alonso, Bastien Moysset, and Ronaldo Messina. Adversarial generation of handwritten text images conditioned on sequences. In *2019 international conference on document analysis and recognition (ICDAR)*, pages 481–486. IEEE, 2019.

[365] Andrew Brock, Jeff Donahue, and Karen Simonyan. Large scale gan training for high fidelity natural image synthesis. *arXiv preprint arXiv:1809.11096*, 2019.

[366] Mario Lučić, Michael Tschannen, Marvin Ritter, Xiaohua Zhai, Olivier Bachem, and Sylvain Gelly. High-fidelity image generation with fewer labels. In *International conference on machine learning*, pages 4183–4192. PMLR, 2019.

[367] Jon Gauthier. Conditional generative adversarial nets for convolutional face generation. *Class project for Stanford CS231N: convolutional neural networks for visual recognition, Winter semester*, 2014(5):2, 2014.

[368] Jost Tobias Springenberg. Unsupervised and semi-supervised learning with categorical generative adversarial networks. *arXiv preprint arXiv:1511.06390*, 2015.

[369] Guim Perarnau, Joost Van De Weijer, Bogdan Raducanu, and Jose M Álvarez. Invertible conditional gans for image editing. *arXiv preprint arXiv:1611.06355*, 2016.

[370] Levent Karacan, Zeynep Akata, Aykut Erdem, and Erkut Erdem. Learning to generate images of outdoor scenes from attributes and semantic layouts. *arXiv preprint arXiv:1612.00215*, 2016.

[371] Carl Vondrick, Hamed Pirsiavash, and Antonio Torralba. Generating videos with scene dynamics. *Advances in neural information processing systems*, 29, 2016.

[372] Ashish Shrivastava, Tomas Pfister, Oncel Tuzel, Joshua Susskind, Wenda Wang, and Russell Webb. Learning from simulated and unsupervised images through adversarial training. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2107–2116, 2017.

[373] Kiana Ehsani, Roozbeh Mottaghi, and Ali Farhadi. Segan: Segmenting and generating the invisible. *arXiv preprint arXiv:1703.10239*, 2017.

[374] 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.

[375] Konstantinos Bousmalis, Nathan Silberman, David Dohan, Dumitru Erhan, and Dilip Krishnan. Unsupervised pixel-level domain adaptation with generative adversarial networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3722–3731, 2017.

[376] Yuan Xue, Tao Xu, Han Zhang, Rodney Long, and Xiaolei Huang. Segan: Adversarial network with multi-scale $l\_1$ loss for medical image segmentation. *arXiv preprint arXiv:1706.01805*, 2017.

[377] Zili Yi, Hao Zhang, Ping Tan, and Minglun Gong. Dualgan: Unsupervised dual learning for image-to-image translation. In *Proceedings of the IEEE international conference on computer vision*, pages 2849–2857, 2017.

[378] Denis Volkhonskiy, Ivan Nazarov, and Evgeny Burnaev. Steganographic generative adversarial networks. *arXiv preprint arXiv:1703.05502*, 2017.

[379] Christian Ledig, Lucas Theis, Ferenc Huszár, Jose Caballero, Andrew Cunningham, Alejandro Acosta, Andrew Aitken, Alykhan Tejani, Johannes Totz, Zehan Wang, et al. Photo-realistic single image super-resolution using a generative adversarial network. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4681–4690, 2017.

[380] Chin-Cheng Hsu, Hsin-Te Hwang, Yi-Chiao Wu, Yu Tsao, and Hsin-Min Wang. Voice conversion from unaligned corpora using variational autoencoding wasserstein generative adversarial networks. *arXiv preprint arXiv:1704.00849*, 2017.

[381] Chongxuan Li, Taufik Xu, Jun Zhu, and Bo Zhang. Triple generative adversarial nets. *Advances in neural information processing systems*, 30, 2017.

[382] Cristóbal Esteban, Stephanie L Hyland, and Gunnar Rätsch. Real-valued (medical) time series generation with recurrent conditional gans. *arXiv preprint arXiv:1706.02633*, 2017.

[383] Arnab Ghosh, Viveka Kulharia, Vinay Namboodiri, Philip HS Torr, and Puneet K Dokania. Multi-agent diverse generative adversarial networks. *arXiv preprint arXiv:1704.02906*, 2017.

[384] Mahesh Gorijala and Ambedkar Dukkipati. Image generation and editing with variational info generative adversarialnetworks. *arXiv preprint arXiv:1701.04568*, 2017.

[385] W Dai, J Doyle, X Liang, H Zhang, N Dong, Y Li, and EP Xing. Scan: Structure correcting adversarial network for chest x-rays organ segmentation. *arXiv preprint arXiv:1703.08770*, 2017.

[386] Hao Dong, Simiao Yu, Chao Wu, and Yike Guo. Semantic image synthesis via adversarial learning. In *Proceedings of the IEEE international conference on computer vision*, pages 5706–5714, 2017.

[387] Jun-Yan Zhu, Richard Zhang, Deepak Pathak, Trevor Darrell, Alexei A Efros, Oliver Wang, and Eli Shechtman. Toward multimodal image-to-image translation. *Advances in neural information processing systems*, 30, 2017.

[388] Ming-Yu Liu, Thomas Breuel, and Jan Kautz. Unsupervised image-to-image translation networks. *Advances in neural information processing systems*, 30, 2017.

[389] Yongyi Lu, Yu-Wing Tai, and Chi-Keung Tang. Conditional cyclegan for attribute guided face image generation. *arXiv preprint arXiv:1705.09966*, 2, 2017.

[390] Guillermo L Grinblat, Lucas C Uzal, and Pablo M Granitto. Class-splitting generative adversarial networks. *arXiv preprint arXiv:1709.07359*, 2017.

[391] Lichao Zhang, Abel Gonzalez-Garcia, Joost Van De Weijer, Martin Danelljan, and Fahad Shahbaz Khan. Synthetic data generation for end-to-end thermal infrared tracking. *IEEE Transactions on Image Processing*, 28(4):1837–1850, 2018.

[392] Chaoyue Wang, Chang Xu, Chaohui Wang, and Dacheng Tao. Perceptual adversarial networks for image-to-image transformation. *IEEE Transactions on Image Processing*, 27(8):4066–4079, 2018.

[393] Seonghyeon Nam, Yunji Kim, and Seon Joo Kim. Text-adaptive generative adversarial networks: manipulating images with natural language. *Advances in neural information processing systems*, 31, 2018.

[394] Piotr Bojanowski, Edouard Grave, Armand Joulin, and Tomas Mikolov. Enriching word vectors with subword information. *Transactions of the association for computational linguistics*, 5:135–146, 2017.

[395] Ting-Chun Wang, Ming-Yu Liu, Jun-Yan Zhu, Guilin Liu, Andrew Tao, Jan Kautz, and Bryan Catanzaro. Video-to-video synthesis. *arXiv preprint arXiv:1808.06601*, 2018.

[396] Daniel Sáez Trigueros, Li Meng, and Margaret Hartnett. Generating photo-realistic training data to improve face recognition accuracy. *arXiv preprint arXiv:1811.00112*, 2018.

[397] Baris Gecer, Binod Bhattarai, Josef Kittler, and Tae-Kyun Kim. Semi-supervised adversarial learning to generate photorealistic face images of new identities from 3d morphable model. In *Proceedings of the European conference on computer vision (ECCV)*, pages 217–234, 2018.

[398] Xiaojie Guo, Lingfei Wu, and Liang Zhao. Deep graph translation. *arXiv preprint arXiv:1805.09980*, 2018.

[399] William Fedus, Ian Goodfellow, and Andrew M Dai. Maskgan: better text generation via filling in the_. *arXiv preprint arXiv:1801.07736*, 2018.

[400] Hoo-Chang Shin, Neil A Tenenholtz, Jameson K Rogers, Christopher G Schwarz, Matthew L Senjem, Jeffrey L Gunter, Katherine P Andriole, and Mark Michalski. Medical image synthesis for data augmentation and anonymization using generative adversarial networks. In *International workshop on simulation and synthesis in medical imaging*, pages 1–11. Springer, 2018.

[401] Aleksei Triastcyn and Boi Faltings. Generating artificial data for private deep learning. *arXiv preprint arXiv:1803.03148*, 2018.

[402] Jesse Engel, Kumar Krishna Agrawal, Shuo Chen, Ishaan Gulrajani, Chris Donahue, and Adam Roberts. Gansynth: Adversarial neural audio synthesis. *arXiv preprint arXiv:1902.08710*, 2019.

[403] Shuangfei Fan and Bert Huang. Labeled graph generative adversarial networks. *arXiv preprint arXiv:1906.03220*, 2019.

[404] Reihaneh Torkzadehmahani, Peter Kairouz, and Benedict Paten. Dp-cgan: Differentially private synthetic data and label generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, pages 0–0, 2019.

[405] Ilya Mironov. Rényi differential privacy. In *2017 IEEE 30th computer security foundations symposium (CSF)*, pages 263–275. IEEE, 2017.

[406] Dawei Zhou, Lecheng Zheng, Jiejun Xu, and Jingrui He. Misc-gan: A multi-scale generative model for graphs. *Frontiers in big Data*, 2:3, 2019.

[407] Chieh Hubert Lin, Chia-Che Chang, Yu-Sheng Chen, Da-Cheng Juan, Wei Wei, and Hwann-Tzong Chen. Coco-gan: Generation by parts via conditional coordinating. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2019.

[408] Karim Armanious, Chenming Jiang, Marc Fischer, Thomas Küstner, Tobias Hepp, Konstantin Nikolaou, Sergios Gatidis, and Bin Yang. Medgan: Medical image translation using gans. *Computerized medical imaging and graphics*, 79, 2020.

[409] Łukasz Maziarka, Agnieszka Pocha, Jan Kaczmarczyk, Krzysztof Rataj, Tomasz Danel, and Michał Warchoł. Mol-cyclegan: a generative model for molecular optimization. *Journal of Cheminformatics*, 12(1):1–18, 2020.

[410] Sina Rashidian, Fusheng Wang, Richard Moffitt, Victor Garcia, Anurag Dutt, Wei Chang, Vishwam Pandya, Janos Hajagos, Mary Saltz, and Joel Saltz. Smooth-gan: towards sharp and smooth synthetic ehr data generation. In *International Conference on Artificial Intelligence in Medicine*, pages 37–48. Springer, 2020.

[411] Salih Sarp, Murat Kuzlu, Emmanuel Wilson, and Ozgur Guler. Wg2an: Synthetic wound image generation using generative adversarial network. *The Journal of Engineering*, 2021(5):286–294, 2021.

[412] Javaria Amin, Muhammad Sharif, Nadia Gul, Seifedine Kadry, and Chinmay Chakraborty. Quantum machine learning architecture for covid-19 classification based on synthetic data generation using conditional adversarial neural network. *Cognitive Computation*, pages 1–12, 2021.

[413] Emily L Denton, Soumith Chintala, Rob Fergus, et al. Deep generative image models using a laplacian pyramid of adversarial networks. *Advances in neural information processing systems*, 28, 2015.

[414] Han Zhang, Tao Xu, Hongsheng Li, Shaoting Zhang, Xiaogang Wang, Xiaolei Huang, and Dimitris N Metaxas. Stackgan: Text to photo-realistic image synthesis with stacked generative adversarial networks. In *Proceedings of the IEEE international conference on computer vision*, pages 5907–5915, 2017.

[415] Xun Huang, Yixuan Li, Omid Poursaeed, John Hopcroft, and Serge Belongie. Stacked generative adversarial networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5077–5086, 2017.

[416] John T Guibas, Tejpal S Virdi, and Peter S Li. Synthetic medical images from dual generative adversarial networks. *arXiv preprint arXiv:1709.01872*, 2017.

[417] Yifan Jiang, Shiyu Chang, and Zhangyang Wang. Transgan: Two pure transformers can make one strong gan, and that can scale up. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems*, volume 34, pages 14745–14758. Curran Associates, Inc., 2021.

[418] Jeff Donahue, Philipp Krähenbühl, and Trevor Darrell. Adversarial feature learning. *arXiv preprint arXiv:1605.09782*, 2016.

[419] Vincent Dumoulin, Ishmael Belghazi, Ben Poole, Olivier Mastropietro, Alex Lamb, Martin Arjovsky, and Aaron Courville. Adversarially learned inference. *arXiv preprint arXiv:1606.00704*, 2016.

[420] Jeff Donahue and Karen Simonyan. Large scale adversarial representation learning. *Advances in neural information processing systems*, 32, 2019.

[421] Alireza Makhzani, Jonathon Shlens, Navdeep Jaitly, Ian Goodfellow, and Brendan Frey. Adversarial autoencoders. *arXiv preprint arXiv:1511.05644*, 2015.

[422] Yujia Li, Kevin Swersky, and Rich Zemel. Generative moment matching networks. In *International conference on machine learning*, pages 1718–1727. PMLR, 2015.

[423] Ilya Tolstikhin, Olivier Bousquet, Sylvain Gelly, and Bernhard Schoelkopf. Wasserstein auto-encoders. *arXiv preprint arXiv:1711.01558*, 2017.

[424] Artur Kadurin, Alexander Aliper, Andrey Kazennov, Polina Mamoshina, Quentin Vanhaelen, Kuzma Khrabrov, and Alex Zhavoronkov. The cornucopia of meaningful leads: Applying deep adversarial autoencoders for new molecule development in oncology. *Oncotarget*, 8(7):10883, 2017.

[425] Artur Kadurin, Sergey Nikolenko, Kuzma Khrabrov, Alex Aliper, and Alex Zhavoronkov. drugan: an advanced generative adversarial autoencoder model for de novo generation of new molecules with desired molecular properties in silico. *Molecular pharmaceutics*, 14(9):3098–3104, 2017.

[426] Pedro Costa, Adrian Galdran, Maria Ines Meyer, Meindert Niemeijer, Michael Abràmoff, Ana Maria Mendonça, and Aurélio Campilho. End-to-end adversarial retinal image synthesis. *IEEE transactions on medical imaging*, 37(3):781–791, 2017.

[427] Oscar Pastor-Serrano, Danny Lathouwers, and Zoltán Perkó. A semi-supervised autoencoder framework for joint generation and classification of breathing. *Computer Methods and Programs in Biomedicine*, 209:106312, 2021.

[428] Chun-Liang Li, Wei-Cheng Chang, Yu Cheng, Yiming Yang, and Barnabás Póczos. Mmd gan: Towards deeper understanding of moment matching network. *Advances in neural information processing systems*, 30, 2017.

[429] Gintare Karolina Dziugaite, Daniel M Roy, and Zoubin Ghahramani. Training generative neural networks via maximum mean discrepancy optimization. *arXiv preprint arXiv:1505.03906*, 2015.

[430] Yong Ren, Jun Zhu, Jialian Li, and Yucen Luo. Conditional generative moment-matching networks. *Advances in Neural Information Processing Systems*, 29, 2016.

[431] Shinnosuke Takamichi, Tomoki Koriyama, and Hiroshi Saruwatari. Sampling-based speech parameter generation using moment-matching networks. *arXiv preprint arXiv:1704.03626*, 2017.

[432] Wei Wang, Yuan Sun, and Saman Halgamuge. Improving mmd-gan training with repulsive loss function. *arXiv preprint arXiv:1812.09916*, 2018.

[433] Wenlong Liao, Yusen Wang, Yuelong Wang, Kody Powell, Qi Liu, and Zhe Yang. Scenario generation for cooling, heating, and power loads using generative moment matching networks. *CSEE Journal of Power and Energy Systems*, 2022.

[434] Anh Nguyen, Alexey Dosovitskiy, Jason Yosinski, Thomas Brox, and Jeff Clune. Synthesizing the preferred inputs for neurons in neural networks via deep generator networks. *Advances in neural information processing systems*, 29, 2016.

[435] Anh Nguyen, Jeff Clune, Yoshua Bengio, Alexey Dosovitskiy, and Jason Yosinski. Plug & play generative networks: Conditional iterative generation of images in latent space. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4467–4477, 2017.

[436] Natasa Tagasovska, Damien Ackerer, and Thibault Vatter. Copulas as high-dimensional generative models: Vine copula autoencoders. *Advances in neural information processing systems*, 32, 2019.

[437] Vaibhav Kulkarni, Natasa Tagasovska, Thibault Vatter, and Benoit Garbinato. Generative models for simulating mobility trajectories. *arXiv preprint arXiv:1811.12801*, 2018.

[438] Haoran Li, Li Xiong, and Xiaoqian Jiang. Differentially private synthesization of multi-dimensional data using copula functions. In *Advances in database technology: proceedings. International conference on extending database technology*, volume 2014, page 475. NIH Public Access, 2014.

[439] Chence Shi, Minkai Xu, Zhaocheng Zhu, Weinan Zhang, Ming Zhang, and Jian Tang. Graphaf: a flow-based autoregressive model for molecular graph generation. *arXiv preprint arXiv:2001.09382*, 2020.

[440] Michael Schlichtkrull, Thomas N Kipf, Peter Bloem, Rianne van den Berg, Ivan Titov, and Max Welling. Modeling relational data with graph convolutional networks. In *European semantic web conference*, pages 593–607. Springer, 2018.

[441] Laurent Dinh, David Krueger, and Yoshua Bengio. Nice: Non-linear independent components estimation. *arXiv preprint arXiv:1410.8516*, 2014.

[442] Laurent Dinh, Jascha Sohl-Dickstein, and Samy Bengio. Density estimation using real nvp. *arXiv preprint arXiv:1605.08803*, 2016.

[443] Durk P Kingma and Prafulla Dhariwal. Glow: Generative flow with invertible 1x1 convolutions. *Advances in neural information processing systems*, 31, 2018.

[444] Ricky TQ Chen, Yulia Rubanova, Jesse Bettencourt, and David K Duvenaud. Neural ordinary differential equations. *Advances in neural information processing systems*, 31, 2018.

[445] Will Grathwohl, Ricky TQ Chen, Jesse Bettencourt, Ilya Sutskever, and David Duvenaud. Ffjord: Free-form continuous dynamics for scalable reversible generative models. *arXiv preprint arXiv:1810.01367*, 2018.

[446] Kaushalya Madhawa, Katushiko Ishiguro, Kosuke Nakago, and Motoki Abe. Graphnvp: An invertible flow model for generating molecular graphs. *arXiv preprint arXiv:1905.11600*, 2019.

[447] Ahmed Alaa, Alex James Chan, and Mihaela van der Schaar. Generative time-series modeling with fourier flows. In *International Conference on Learning Representations*, 2021.

[448] Vincent François-Lavet, Peter Henderson, Riashat Islam, Marc G Bellemare, and Joelle Pineau. An introduction to deep reinforcement learning. *arXiv preprint arXiv:1811.12560*, 2018.

[449] Junhyuk Oh, Xiaoxiao Guo, Honglak Lee, Richard L Lewis, and Satinder Singh. Action-conditional video prediction using deep networks in atari games. *Advances in neural information processing systems*, 28, 2015.

[450] Biao Jia, Chen Fang, Jonathan Brandt, Byungmoon Kim, and Dinesh Manocha. Paintbot: A reinforcement learning approach for natural media painting. *arXiv preprint arXiv:1904.02201*, 2019.

[451] Arjun Krishna, Kedar Bartake, Chuang Niu, Ge Wang, Youfang Lai, Xun Jia, and Klaus Mueller. Image synthesis for data augmentation in medical ct using deep reinforcement learning. *arXiv preprint arXiv:2103.10493*, 2021.

[452] Jiwei Li, Will Monroe, Alan Ritter, Michel Galley, Jianfeng Gao, and Dan Jurafsky. Deep reinforcement learning for dialogue generation. *arXiv preprint arXiv:1606.01541*, 2016.

[453] Jiwei Li, Will Monroe, Tianlin Shi, Sébastien Jean, Alan Ritter, and Dan Jurafsky. Adversarial learning for neural dialogue generation. *arXiv preprint arXiv:1701.06547*, 2017.

[454] Gabriel Lima Guimaraes, Benjamin Sanchez-Lengeling, Carlos Outeiral, Pedro Luis Cunha Farias, and Alán Aspuru-Guzik. Objective-reinforced generative adversarial networks (organ) for sequence generation models. *arXiv preprint arXiv:1705.10843*, 2017.

[455] Marcus Olivecrona, Thomas Blaschke, Ola Engkvist, and Hongming Chen. Molecular de-novo design through deep reinforcement learning. *Journal of cheminformatics*, 9(1):1–14, 2017.

[456] Zhan Shi, Xinchi Chen, Xipeng Qiu, and Xuanjing Huang. Toward diverse text generation with inverse reinforcement learning. *arXiv preprint arXiv:1804.11258*, 2018.

[457] Ronald J Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, 8(3):229–256, 1992.

[458] Jiaxuan You, Bowen Liu, Zhitao Ying, Vijay Pande, and Jure Leskovec. Graph convolutional policy network for goal-directed molecular graph generation. *Advances in neural information processing systems*, 31, 2018.

[459] Harish Kumar and Balaraman Ravindran. Polyphonic music composition with lstm neural networks and reinforcement learning. *arXiv preprint arXiv:1902.01973*, 2019.

[460] Daniel Jarrett, Ioana Bica, and Mihaela van der Schaar. Time-series generation by contrastive imitation. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems*, volume 34, pages 28968–28982. Curran Associates, Inc., 2021.

[461] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*, 33:6840–6851, 2020.

[462] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *International Conference on Machine Learning*, pages 2256–2265. PMLR, 2015.

[463] Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models. In *International Conference on Machine Learning*, pages 8162–8171. PMLR, 2021.

[464] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in Neural Information Processing Systems*, 34:8780–8794, 2021.

[465] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 2022.

[466] Nikolaus Mayer, Eddy Ilg, Philipp Fischer, Caner Hazirbas, Daniel Cremers, Alexey Dosovitskiy, and Thomas Brox. What makes good synthetic training data for learning disparity and optical flow estimation? *International Journal of Computer Vision*, 126(9):942–960, 2018.

[467] Daniel J Butler, Jonas Wulff, Garrett B Stanley, and Michael J Black. A naturalistic open source movie for optical flow evaluation. In *European conference on computer vision*, pages 611–625. Springer, 2012.

[468] Ton Roosendaal. Sintel. 2010. Accessed: 2022-07-18.

[469] Blender - free and open 3d creation software. `https://www.blender.org/`. Accessed: 2022-07-18.

[470] Ankur Handa, Thomas Whelan, John McDonald, and Andrew J Davison. A benchmark for rgb-d visual odometry, 3d reconstruction and slam. In *2014 IEEE international conference on Robotics and automation (ICRA)*, pages 1524–1531. IEEE, 2014.

[471] The persistence of vision raytracer. `https://www.povray.org/`. Accessed: 2022-07-18.

[472] Hao Su, Charles R Qi, Yangyan Li, and Leonidas J Guibas. Render for cnn: Viewpoint estimation in images using cnns trained with rendered 3d model views. In *Proceedings of the IEEE international conference on computer vision*, pages 2686–2694, 2015.

[473] Xingchao Peng, Baochen Sun, Karim Ali, and Kate Saenko. Learning deep object detectors from 3d models. In *Proceedings of the IEEE international conference on computer vision*, pages 1278–1286, 2015.

[474] Ankur Handa, Viorica Pătrăucean, Simon Stent, and Roberto Cipolla. Scenenet: An annotated model generator for indoor scene understanding. In *2016 IEEE International Conference on Robotics and Automation (ICRA)*, pages 5737–5743. IEEE, 2016.

[475] Ankur Handa, Viorica Patraucean, Vijay Badrinarayanan, Simon Stent, and Roberto Cipolla. Understanding real world indoor scenes with synthetic data. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4077–4085, 2016.

[476] Martin Peris, Sara Martull, Atsuto Maki, Yasuhiro Ohkawa, and Kazuhiro Fukui. Towards a simulation driven stereo vision system. In *Proceedings of the 21st International Conference on Pattern Recognition (ICPR2012)*, pages 1038–1042. IEEE, 2012.

[477] Javier Molina, José A Pajuelo, Marcos Escudero-Viñolo, Jesús Bescós, and José M Martínez. A natural and synthetic corpus for benchmarking of hand gesture recognition systems. *Machine Vision and Applications*, 25(4):943–954, 2014.

[478] Baochen Sun and Kate Saenko. From virtual to reality: Fast adaptation of virtual object detectors to real domains. In *BMVC*, volume 1, page 3, 2014.

[479] Konstantinos Rematas, Tobias Ritschel, Mario Fritz, and Tinne Tuytelaars. Image-based synthesis and re-synthesis of viewpoints guided by 3d models. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3898–3905, 2014.

[480] Jeremie Papon and Markus Schoeler. Semantic pose using deep networks trained on synthetic rgb-d. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 774–782, 2015.

[481] Philipp Fischer, Alexey Dosovitskiy, Eddy Ilg, Philip Häusser, Caner Hazırbaş, Vladimir Golkov, Patrick Van der Smagt, Daniel Cremers, and Thomas Brox. Flownet: Learning optical flow with convolutional networks. *arXiv preprint arXiv:1504.06852*, 2015.

[482] Adam Kortylewski, Andreas Schneider, Thomas Gerig, Bernhard Egger, Andreas Morel-Forster, and Thomas Vetter. Training deep face recognition systems with synthetic data. *arXiv preprint arXiv:1802.05891*, 2018.

[483] Markus Philipp, Neal Bacher, Jonas Nienhaus, Lars Hauptmann, Laura Lang, Anna Alperovich, Marielena Gutt-Will, Andrea Mathis, Stefan Saur, Andreas Raabe, et al. Synthetic data generation for optical flow evaluation in the neurosurgical domain. *Current directions in biomedical engineering*, 7(1):67–71, 2021.

[484] Aleksei Boikov, Vladimir Payor, Roman Savelev, and Alexandr Kolesnikov. Synthetic data generation for steel defect detection and classification using deep learning. *Symmetry*, 13(7):1176, 2021.

[485] German Ros, Laura Sellart, Joanna Materzynska, David Vazquez, and Antonio M Lopez. The synthia dataset: A large collection of synthetic images for semantic segmentation of urban scenes. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3234–3243, 2016.

[486] Vladimir Haltakov, Christian Unger, and Slobodan Ilic. Framework for generation of synthetic ground truth data for driver assistance applications. In *German conference on pattern recognition*, pages 323–332. Springer, 2013.

[487] Vdrift - open-source driving simulation. `https://vdrift.net/`. Accessed: 2022-07-18.

[488] Stephan R Richter, Vibhav Vineet, Stefan Roth, and Vladlen Koltun. Playing for data: Ground truth from computer games. In *European conference on computer vision*, pages 102–118. Springer, 2016.

[489] Matthew Johnson-Roberson, Charles Barto, Rounak Mehta, Sharath Nittur Sridhar, Karl Rosaen, and Ram Vasudevan. Driving in the matrix: Can virtual worlds replace human-generated annotations for real world tasks? *arXiv preprint arXiv:1610.01983*, 2016.

[490] Unity real-time development platform. `https://unity.com/`. Accessed: 2022-07-18.

[491] Alireza Shafaei, James J Little, and Mark Schmidt. Play and learn: Using video games to train computer vision models. *arXiv preprint arXiv:1608.01745*, 2016.

[492] Adrien Gaidon, Qiao Wang, Yohann Cabon, and Eleonora Vig. Virtual worlds as proxy for multi-object tracking analysis. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4340–4349, 2016.

[493] Andreas Geiger, Philip Lenz, and Raquel Urtasun. Are we ready for autonomous driving? the kitti vision benchmark suite. In *2012 IEEE conference on computer vision and pattern recognition*, pages 3354–3361. IEEE, 2012.

[494] Stephan R Richter, Zeeshan Hayder, and Vladlen Koltun. Playing for benchmarks. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2213–2222, 2017.

[495] Amlan Kar, Aayush Prakash, Ming-Yu Liu, Eric Cameracci, Justin Yuan, Matt Rusiniak, David Acuna, Antonio Torralba, and Sanja Fidler. Meta-sim: Learning to generate synthetic datasets. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4551–4560, 2019.

[496] Marc G Bellemare, Yavar Naddaf, Joel Veness, and Michael Bowling. The arcade learning environment: An evaluation platform for general agents. *Journal of Artificial Intelligence Research*, 47:253–279, 2012.

[497] Michał Kempka, Marek Wydmuch, Grzegorz Runc, Jakub Toczek, and Wojciech Jaśkowski. Vizdoom: A doom-based ai research platform for visual reinforcement learning. In *2016 IEEE conference on computational intelligence and games (CIG)*, pages 1–8. IEEE, 2016.

[498] Fereshteh Sadeghi and Sergey Levine. Cad2rl: Real single-image flight without a single real image. *arXiv preprint arXiv:1611.04201*, 2016.

[499] Oriol Vinyals, Timo Ewalds, Sergey Bartunov, Petko Georgiev, Alexander Sasha Vezhnevets, Michelle Yeo, Alireza Makhzani, Heinrich Küttler, John Agapiou, Julian Schrittwieser, et al. Starcraft ii: A new challenge for reinforcement learning. *arXiv preprint arXiv:1708.04782*, 2017.

[500] Charles Beattie, Joel Z Leibo, Denis Teplyashin, Tom Ward, Marcus Wainwright, Heinrich Küttler, Andrew Lefrancq, Simon Green, Víctor Valdés, Amir Sadik, et al. Deepmind lab. *arXiv preprint arXiv:1612.03801*, 2016.

[501] Gabriel Synnaeve, Nantas Nardelli, Alex Auvolat, Soumith Chintala, Timothée Lacroix, Zeming Lin, Florian Richoux, and Nicolas Usunier. Torchcraft: a library for machine learning research on real-time strategy games. *arXiv preprint arXiv:1611.00625*, 2016.

[502] Matthew Johnson, Katja Hofmann, Tim Hutton, and David Bignell. The malmo platform for artificial intelligence experimentation. In *Ijcai*, pages 4246–4247. Citeseer, 2016.

[503] Zhengwei Wang, Qi She, and Tomas E Ward. Generative adversarial networks in computer vision: A survey and taxonomy. *arXiv preprint arXiv:1906.01529*, 2019.

[504] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.

[505] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.

[506] Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild. In *Proceedings of International Conference on Computer Vision (ICCV)*, December 2015.

[507] Jia Deng, Wei Dong, Richard Socher, Li-Jia 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. Ieee, 2009.

[508] Fisher Yu, Ari Seff, Yinda Zhang, Shuran Song, Thomas Funkhouser, and Jianxiong Xiao. Lsun: Construction of a large-scale image dataset using deep learning with humans in the loop. *arXiv preprint arXiv:1506.03365*, 2015.

[509] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European conference on computer vision*, pages 740–755. Springer, 2014.

[510] Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, Bo Wu, and Andrew Y Ng. Reading digits in natural images with unsupervised feature learning. 2011.

[511] Darrell Conklin. Bach Chorales. UCI Machine Learning Repository.

[512] Teague Sterling and John J. Irwin. Zinc 15 – ligand discovery for everyone. *Journal of Chemical Information and Modeling*, 55(11):2324–2337, 2015. PMID: 26479676.

[513] Fida K Dankar and Mahmoud Ibrahim. Fake it till you make it: Guidelines for effective synthetic data generation. *Applied Sciences*, 11(5):2158, 2021.

[514] Shan Chang and Chao Li. Privacy in neural network learning: Threats and countermeasures. *IEEE Network*, 32(4):61–67, 2018.

[515] Jose D Fernández and Francisco Vico. Ai methods in algorithmic composition: A comprehensive survey. *Journal of Artificial Intelligence Research*, 48:513–582, 2013.

[516] Xiaodong He and Li Deng. Deep learning for image-to-text generation: A technical overview. *IEEE Signal Processing Magazine*, 34(6):109–116, 2017.

[517] Xin Yi, Ekta Walia, and Paul Babyn. Generative adversarial network in medical imaging: A review. *Medical image analysis*, 58:101552, 2019.

[518] Mohammad Abufadda and Khalid Mansour. A survey of synthetic data generation for machine learning. In *2021 22nd International Arab Conference on Information Technology (ACIT)*, pages 1–7. IEEE, 2021.

# A  Acronyms

**AAE**  Adversarial Autoencoder

**AF**  Autoregressive Flow

**AI**  Artificial Intelligence

**AIC**  Akaike information criterion

**ALI**  Adversarially Learned Inference

**BIC**  Bayesian information criterion

**BN**  Bayesian Network

**BPTT**  Backpropagation Through Time

**BiGAN**  Bidirectional GAN

**C2ST**  Classifier Two Sample Test

**CAD**  Computer Aided Design

**CAE**  Contractive Autoencoder

**CAN**  Creative Adversarial Network

**CDBN**  Conditional Deep Belief Network

**CNN**  Convolutional Neural Network

**CRBM**  Conditional Restricted Boltzmann Machine

**CT**  Computed Tomography

**DAE**  Denoising Autoencoder

**DAG**  Directed Acyclic Graph

**DL**  Deep Learning

**DBM**  Deep Boltzmann Machine

**DBN**  Deep Belief Network

**DCGAN**  Deep Convolutional GAN

**DDPM**  Denoising Diffusion Probabilistic Model

**DLGM**  Deep Latent Gaussian Model

**ECG**  Electrocardiogram

**EHR**  Electronic Health Record

**ELBO**  Evidence Lower-Bound

**FID**  Fréchet Inception Distance

**GA**  Genetic Algorithm

**GAE**  Gated Autoencoder

**GAN**  Generative Adversarial Net

**GMM**  Gaussian Mixture Model

**GMMN**  Generative Moment Matching Network

**GNN**  Graph Neural Network

**GPU**  Graphics Processing Unit

**GRU**  Gated Recurrent Unit

**GSN**  Generative Stochastic Network

**HMM**  Hidden Markov Model

**HSI**  Hyperspectral Image

**IS**  Inception Score

**KDE**  Kernel Density Estimators

**LDA** Latent Dirichlet Allocation

**LRCN** Long-term Recurrent Convolutional Network

**LSTM** Long Short-Term Memory

**MCMC** Markov chain Monte Carlo

**ML** Machine Learning

**MLP** Multilayer Perceptron

**MMD** Maximum Mean Discrepancy

**NADE** Neural Autoregressive Distribution Estimator

**NLL** Negative Log-Likelihood

**NLP** Natural Language Processing

**ODE** Ordinary Differential Equation

**PPG** Photoplethysmogram

**PPGN** Plug & Play Generative Network

**PSD** Private Spatial Decomposition

**RBM** Restricted Boltzmann Machine

**RL** Reinforcement Learning

**RNN** Recurrent Neural Network

**RTRBM** Recurrent Temporal Boltzmann Machine

**ReLU** rectified linear unit

**SDG** Synthetic Data Generation

**SLAM** Simultaneous Localization and Mapping

**SRTRBM** Structured Recurrent Temporal Boltzmann Machine

**SVHN** Street View House Numbers

**TFD** Toronto Face Database

**TRBM** Temporal Restricted Boltzmann Machine

**VAE** Variational Autoencoder

**VCAE** Vine Copula Autoencoder

**WAE** Wasserstein Autoencoder