

Artificial Creativity

Federico Matteoni
Intelligent Systems for Pattern Recognition
A.Y. 2022/2023
Mat. 530257

Abstract—This work is a survey on some common machine learning systems aimed at image generation that surfaced in the last years. These systems have recently been used by the general public in viral fashion to generate artistic content, ranging from artistic repropositions of submitted photos to completely new concepts produced by text-based prompts. Although exceptional, these systems also raised some concerns regarding copyright and intellectual property.

I. INTRODUCTION

In January 2021, OpenAI labs presented DALL-E [18], a text-to-image generative autoregressive transformer, followed by an unofficial demo. In that demo, users could simply type a text prompt and after few moments a small selection of artistic interpretations of that prompt would be shown. While not being the first text-to-image machine learning model to be presented, it can be argued that DALL-E is the first that had a big impact on the public: in the following months many examples of images generated with that tool would be posted on social media, and many more sophisticated models would surface, taking advantage of the increasing interest in the image generation and *artificial creativity* field.

Many of the most recent and successful of these systems stem from the "Diffusion Probabilistic Models" [21] family of machine learning approaches, while some of the earliest examples are Generative Adversarial Networks [13] with alternative techniques for the sample generation [15] compared to the classical GAN approaches.

II. MODELS

Image generation, and broadly speaking every kind of data generated through machine learning, is achieved with a family of models called "Generative Models". These models, in contrast with the family of "Discriminative Models", aim to generate new samples of data after learning its distribution, while the goal of discriminative models is to distinguish between different samples of data. Mathematically speaking, given a set of samples X and the corresponding labels Y , a generative model learns $p(X, Y)$ while a discriminative model learns $p(X | Y)$. Generative models can be explicit, when they directly learn $p(X)$, or implicit, if they learn how to sample from $p(X)$ without directly learning it.

Text-to-image generation can be achieved with many different kinds of generative models, but we will focus on three main family of models that recently achieved the best results: Generative Adversarial Networks, Transformers and Diffusion Models.

A. Generative Adversarial Networks

A Generative Adversarial Network (GAN) [13] is composed of two main parts:

- A Generator \mathcal{G} that learns to generate plausible data with the goal of confusing the discriminator
- A Discriminator \mathcal{D} that learns to distinguish between the real data and the fake data, penalizing the generator when it produces non-plausible samples

The trick that GANs employ to be able to sample data from complex, high-dimensional distributions, is to sample from a very simple distribution (usually a Gaussian, $\mathcal{N}(0, 1)$) and to train a differentiable distribution (usually a Neural Network, \mathcal{G}) to transform the sampled random noise into data from the target distribution:

- Sample $z \sim \mathcal{N}(0, 1)$
- Generate $x = \mathcal{G}(z)$
- Discriminate $\mathcal{D}(x) = \mathcal{D}(\mathcal{G}(z))$

Given real data x and random noise $z \sim \mathcal{N}(0, 1)$, the objective function of a GAN is C defined in (1)

$$C = \min_{\mathcal{G}} \max_{\mathcal{D}} (\mathbb{E}_x[\log \mathcal{D}(x)] - \mathbb{E}_z[\log(1 - \mathcal{D}(\mathcal{G}(z)))] \quad (1)$$

This represents a hard two-player game, where the optimal solution resides in a saddle point that corresponds to the concurrent min max.

Most of the works that employ GANs for data generation focus on developing different strategies that adapt the general architecture of these models to the interested data domain. Among the first examples of image generations that reached the general public is thispersondoesnotexist.com, a website published by nVidia which employs a particular GAN architecture, called StyleGAN, to generate highly realistic human portraits. Released in December 2018 with StyleGAN [15] and updated in December 2019 with StyleGAN2 [16], this project may represent the first impact that the general public had with state-of-the-art automatic image generation via machine learning models.

The original motivation behind StyleGAN was to unbox and better understand the image synthesis process and the properties of the latent spaces employed by the GANs, to both better control the generative process and to better compare different generators between each other. By taking inspiration from the style-transfer literature [14], which studies techniques of merging the content of an image with the artistic style of another, the StyleGAN proposed a redesigned generative process: the input latent code (random noise) z is mapped

into an intermediate latent space \mathcal{W} which is used to control the generator \mathcal{G} at each layer together with added Gaussian noise. This is in contrast to traditional generators, which feed the input directly to the generating network thus using the latent data just in the input layer. StyleGAN uses latent data to control the generative process at each layer, with \mathcal{W} fed via adaptive instance normalization (AdaIN, (2)): it normalizes each feature map x_i separately given the style $y = (y_s, y_b)$ which is synthesized from \mathcal{W} via learned transformations.

$$\text{AdaIN}(x_i, y) = y_{si} \frac{x_i - \mu(x_i)}{\sigma(x_i)} + y_{bi} \quad (2)$$

A main difference from classical style transfer [14] is that the style y of the target image is computed from \mathcal{W} , the latent data, instead of an example image.

The noise which is fed independently to each layer of the generative process is used to generate stochastic detail in the final images: these details range from the placement of the skin pores to other details like freckles and hair. By using random noise at each layer the network does not need to use activations from the previous layer, thus reusing data from the only source of input available, to generate stochastic effects: this reduces the occurrence of repetitive patterns, which can be particularly perceivable in the case of human faces. According to style transfer literature, the style of an image is reliably encoded by spatially invariant statistics like channel-wise mean or variance etc., and in StyleGAN the globally defined style affects global effects like lightning, pose and the like, while locally defined noise added independently to each pixel is only used to generate stochastic variation: if the generator tried to control a global effect, e.g. pose, with the noise, this would lead to spatially inconsistent decisions and then penalization by the discriminator. So the network learns to use global and local channels appropriately without explicit guidance, fully employing the power of the adversarial process.

Besides thispersondoesnotexists, StyleGANs have been used in other online demos like thisartworkdoesnotexists and thiscatdoesnotexists: an example of the latter is Fig. 1.



Fig. 1. A cat portrait generated by a StyleGAN.

B. Transformers

Transformers [23] are models that make extensive use of the Attention mechanism [22]. Attention is a technique that enhances small portions of the input data that are deemed important depending on the context. Given h_i as input, e.g. tokens of latent information, and context information S , in its simplest form an attention module computes C in three passes:

- Relevance, where each encoding h_i is fused with current context S yielding the relevance η_i between the two

$$\eta_i = \tanh(S, h_i) \quad (3)$$

- Normalization via a softmax operator

$$\alpha_i = \frac{e^{\eta_i}}{\sum_j e^{\eta_j}} \quad (4)$$

- Aggregation by soft attention voting

$$C = \sum_i \alpha_i h_i \quad (5)$$

Transformers are pure attention models with no recurrence that employ a variant of the attention mechanism called self-attention, used to draw global dependencies of the data. Self-attention maps each input to a set of three vectors, called Query (Q), Key (K) and Value (V): these vectors are used to compute the output vector, which is a weighted sum of the Vs where the weights are computed through a compatibility function of the Q and the corresponding K. The version used traditionally in the Transformers is called Scaled Dot-Product Attention, in Eq. 6 where d_k is the dimension of Q and K.

$$\text{Attention}(Q, K, V) = \text{Softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V \quad (6)$$

A Transformer is made of an encoder and a decoder, both composed by stack of self-attention and fully connected layers. In particular, Autoregressive Transformers use masked attention in order to hide future positions, thus depending only on previous outputs and inputs.

DALL-E [18] is an autoregressive transformer that models the prompt text and image as a single stream of tokens. In order to not use pixels directly as image tokens, and to prevent that much of the modeling capacity is spent capturing high-frequency details (like the stochastic variations of the displacement of hairs and freckles) rather than low-frequency details (like pose and lightning), the training is split into two stages [18]:

- A Discrete Variational Autoencoder (dVAE) is used to compress the images into a grid of tokens, each with 8192 possible values
- Up to 256 text tokens are concatenated with the 1024 image tokens to train an autoregressive transformer, that models the joint distribution over the text and image tokens.

The training is performed by maximizing the ELBO of the joint likelihood of the model distribution over images x , captions y and tokens z of the encoded image. By factorizing

the distribution with $p_{\theta\psi}(x, y, z) = p_\theta(x | y, z)p_\psi(y, z)$ we have the following lower bound:

$$\ln p_{\theta\psi}(x, y) \geq \mathbb{E}_{z \sim q_\phi(z | x)} [\ln p_\theta(x | y, z) - \beta KL(q_\phi(y, z | x) || p_\psi(y, z))] \quad (7)$$

Where q_ϕ is the distribution over the images tokens generated by the dVAE given the image, p_θ is the distribution over the images generated by the decoder given the tokens, p_ψ is the joint distribution over the text and image tokens modeled by the transformer.

The resulting model, called DALL-E and shown in Fig. 2, has been presented in January 2021 [1] and gained online relevance as soon as it came out thanks to a unofficial open-source version formerly called DALL-E mini [11] of which an example of the latest version, called Craiyon, is Fig. 3. Much of the success that DALL-E faced since its presentation came from the many social media communities that used the tool to quickly generate and share memes, creating a viral phenomenon that contributed to an increased interest in text-to-image models research and in general in the artistic applications of machine learning models.



Fig. 2. "A low-angle view of a capybara sitting on a mountain" by DALL-E



Fig. 3. "A low-angle view of a capybara sitting on a mountain" by Craiyon

C. Diffusion Probabilistic Models

One of the main limitations of machine learning, regarding generative models in particular, is the tractability of the probability distributions involved. Learning, sampling and evaluation are often still analytically or computationally intractable when using highly flexible families of models, thus requiring a tradeoff between these two conflicting goals: flexibility and tractability. Diffusion Probabilistic Models (DM) [21] aim to achieve both, by taking inspiration from statistical physics. DMs are built by two main processes:

- A Diffusion Process which converts any complex data distribution into a simple one (e.g. a Gaussian)
- A Reverse Process which defines the generative model

The diffusion process gradually converts an input data distribution $q(x^{(0)})$ into a tractable distribution $p(y)$ by repeated application of a Markov Diffusion Kernel (8) where β is the diffusion rate

$$T_p(y | y'; \beta) \quad (8)$$

$$p(y) = \int T_p(y | y'; \beta) p(y') dy' \quad (9)$$

$$q(x^{(t)} | x^{(t-1)}) = T_p(x^{(t)} | x^{(t-1)}; \beta) \quad (10)$$

The forward trajectory corresponds to performing T steps of diffusion starting from the data distribution and is defined in (11), which exhibits clear Markovian dynamics showing a prior and a first-order transition

$$q(x^{(0...T)}) = q(x^{(0)}) \prod_{t=1}^T q(x^{(t)} | x^{(t-1)}) \quad (11)$$

The forward diffusion kernel $q(x^{(t)} | x^{(t-1)})$ for example may correspond to a Gaussian Diffusion into a Gaussian distribution with identity covariance, defined in (12)

$$q(x^{(t)} | x^{(t-1)}) = \mathcal{N}(x^{(t-1)} \sqrt{1 - \beta_t}, \mathbf{I}\beta_t) \quad (12)$$

The generative distribution is trained to describe the same trajectory of the forward process but in reverse, thus yielding the reverse process described in (13)

$$p(x^{(0...T)}) = p(x^{(T)}) \prod_{t=1}^T p(x^{(t-1)} | x^{(t)}) \quad (13)$$

During learning, and in the Gaussian diffusion kernel scenario taken as example, only the mean and covariance of the kernel need to be estimated:

- $f_\mu(x^{(t)}, t)$ defines the mean of the reverse Markov transitions
- $f_\sigma(x^{(t)}, t)$, defines the covariance of the reverse Markov transitions

Both can be defined via neural networks or other function fitting techniques and the cost of these functions defines the computational cost of the overall algorithm. Those are used in the reverse diffusion kernel defined in (14)

$$p(x^{(t-1)} | x^{(t)}) = \mathcal{N}(f_\mu(x^{(t)}, t), f_\sigma(x^{(t)}, t)) \quad (14)$$

The main limitation of classical DMs is that they operate directly in pixel space and require the sequential run of a large number of steps of the Markov chain, thus resulting in expensive and long optimization (e.g. 150-1000 GPU days [12]) and inference processes (circa 5 days to produce 50k images [12]). This has multiple consequences, beginning with the carbon footprint, the unavailability of such computational resources to many researchers and the unapplicability of these models to real-word scenarios.

To increase the accessibility of these models, the Latent Diffusion Model (LDM) [19] family of approaches have been proposed. These models split training into two distinct phases:

- Train an Autoencoder [10] that provides a lower-dimensional representational space, more efficient than pixel space and perceptually equivalent
- Train the DM in the learned latent, smaller and more efficient, space

The resulting model, a LDM, exhibits better scaling properties and more efficient image generation than traditional DMs. A notable advantage is that the autoencoder may need to be trained just once, and then be used for multiple different LDMs or for completely different models.

A successful application of LDMs is the Stable Diffusion model [2] released in 2022, an open source implementation of an LDMs trained on the LAION-5B dataset [3]. This model has been used in several projects [4] and in a online demo [5]: an example of an image generated by this demo is Fig. 4. Another notable example is Midjourney [6], released in 2022 as well, a commercial application of a DM with the goal of being a tool for artists. An example of an image produced by Midjourney is Fig. 5.

LDMs recently have been implemented in mobile apps, highlighting their generation efficiency and general-purpose. This is one of the main differences between GANs (e.g. StyleGAN) and DMs (e.g. Stable Diffusion), with the first achieving great result only on data with limited variability with an overall small flexibility, while the latter showing more flexibility.



Fig. 4. "A teacher explaining artificial intelligence with holograms" by Stable Diffusion 2.1



Fig. 5. "Ukiyo-e, robot, detailed" by Midjourney

III. DATASETS

In just few years, from models capable of generating high quality samples from a specific domain like human portraits or cats, we now have models that are capable of generating arbitrary images starting from a simple text prompt and in just few moments. Aside from the novel ideas and solutions regarding training and generation, very important contributors of this fast improvement are the datasets on which the models have been trained.

A. LAION

LAION-400M [20], released in 2021, and its newest version LAION-5B [3], released in 2022, are openly available datasets composed by image-text pairs. Massive datasets, like the one used to train GPT-3 or to scale up various models, are rarely released publicly maybe due to industrial secret or legal reasons. LAION-400M was initially released in order to address this issue, by providing a curated and filtered massive dataset in order to push forward the research involving these data. LAION-5B contains more than 5 billion image-text pairs, for a total of 240 terabytes of data. It also provides kNN indices to speed up search, tools for building the dataset, subset of the datasets useful for testing or more focused tasks and other metadata. LAION datasets are built through a post-processing of data contained in the Common Crawl dataset.

B. Common Crawl

Common Crawl [7] is a corpus containing data collected since 2008 and composed by raw web pages, metadata and text for a total of circa 420 terabytes. Its access is completely free. It works via a bot, now called CCBot/2.0, that scrapes the web, accesses each webpage identifying itself and gathering the data it needs. The data is then stored in servers where it is processed and archived in the WARC format, a standard format for libraries which enables storing multiple resources and metadata into a single archive file.

IV. OTHER ARTISTIC FIELDS

Up until now the focus of this review has been text-to-image generation. The main reason for this is due to the recent improvements in this fields, and the social media resonance that models like Midjourney have obtained thanks to the high-quality results. But text-to-image generation isn't the only artistic field that has been targeted by AI researchers, and some notable examples are the following.

A. Text Generation: ChatGPT

Text generation is among the first fields that researchers studied as possible applications of machine learning models. Text generation requires models that are able of Natural Language Processing (NLP) and this field alone has sparked many developments, of which all other fields benefited from. One example are transformers, and one of the most successful application of transformers is ChatGPT [8], a model derived from the general-purpose GPT-3 (Generative Pretrained Transformer). ChatGPT can describe itself: "I am a large language model trained by OpenAI. Language models are trained to predict the next word in a sequence of words, based on the context of the words that come before it. This allows me to understand and generate human-like text. I was trained on a dataset of more than 8 million web pages, including articles, books, and websites.".

B. 3D Objects Generation: Point-E

The generation of 3D objects may be seen as an extension of text-to-image generation in which errors of the model as well as generation efficiency are even more crucial aspects. Point-E [17] is a combination of a text-to-image diffusion model that generates a view of the object and a diffusion model which produces the 3D object (specifically, a cloud of points in the three-dimensional space) by being condition on the generated image. A key aspect of Point-E is the generation efficiency: while only producing a 3D point cloud, thus generating objects that are not on-par with state-of-the-art quality, it achieves so in just a couple of minutes with a single GPU compared to the multiple GPU-hours required by current methods.

C. Music Generation: Museformer

Music has always been regarded as a crucial field for developing machine learning models that are able to convincingly imitate human intelligence, as how humans perceive music is still under study. Thanks to competitions like AI Song Contest, music generation is one of the fields where novel training methods and other techniques can be put to test and compared with other researchers. A notable example is Museformer [24], a transformer-based model trained on the Lakh MIDI dataset, a collection of MIDI-based versions of a million popular contemporary music tracks. Museformer foregoes the full attention mechanism of classical transformers for a combination of fine and coarse attention engineered specifically for music generation, by attending to specific tokens that are more relevant in music structures. The result is an efficient model that can produce longer tracks than competitors.

V. CONCLUSION

In just few years, the field of text-to-image generations has seen a great deal of improvement both in terms of sample quality as well as generation speed. Current models like Stable Diffusion and Midjourney are able to generate high-quality images starting from user-provided text prompts in just a couple of minutes, and often with an artistic quality that is comparable to highly-regarded artists. This is a crucial observation, because it has sparked a lot of attention and dissatisfaction with the AI research community regarding the training methods and the dataset that have been used, particularly in terms of copyright. Another key aspect is the rising need of the development of systems like invisible watermarks [9] that enables the undisputable marking of AI-generated images.

REFERENCES

- [1] <https://openai.com/blog/dall-e/>. [Online; Accessed 2022-12-19].
- [2] <https://github.com/CompVis/stable-diffusion>. [Online; Accessed 2022-12-18].
- [3] <https://laion.ai/blog/laion-5b/>. [Online; Accessed 2022-12-18].
- [4] <https://www.banana.dev/blog/project-ideas-built-with-stable-diffusion>. [Online; Accessed 2022-12-18].
- [5] <https://huggingface.co/spaces/stabilityai/stable-diffusion>. [Online; Accessed 2022-12-18].
- [6] <https://www.midjourney.com/home/>. [Online; Accessed: 2022-12-16].
- [7] <https://commoncrawl.org/>. [Online; Accessed 2022-12-20].
- [8] <https://openai.com/blog/chatgpt/>. [Online; Accessed 2022-12-21].
- [9] <https://github.com/ShieldMnt/invisible-watermark/>. [Online; Accessed 2022-12-21].
- [10] D. Bank, N. Koenigstein, and R. Giryes. Autoencoders. *CoRR*, abs/2003.05991, 2020.
- [11] B. Dayma. <https://github.com/borisdayma/dalle-mini/>. [Online; Accessed 2022-12-19].
- [12] P. Dhariwal and A. Nichol. Diffusion models beat gans on image synthesis. *CoRR*, abs/2105.05233, 2021.
- [13] I. J. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative adversarial networks. 2014.
- [14] X. Huang and S. J. Belongie. Arbitrary style transfer in real-time with adaptive instance normalization. *CoRR*, abs/1703.06868, 2017.
- [15] T. Karras, S. Laine, and T. Aila. A style-based generator architecture for generative adversarial networks. *CoRR*, abs/1812.04948, 2018.
- [16] T. Karras, S. Laine, M. Aittala, J. Hellsten, J. Lehtinen, and T. Aila. Analyzing and improving the image quality of stylegan. *CoRR*, abs/1912.04958, 2019.
- [17] A. Nichol, H. Jun, P. Dhariwal, P. Mishkin, and M. Chen. Point-e: A system for generating 3d point clouds from complex prompts. 2022.
- [18] A. Ramesh, M. Pavlov, G. Goh, S. Gray, C. Voss, A. Radford, M. Chen, and I. Sutskever. Zero-shot text-to-image generation, 2021.
- [19] R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer. High-resolution image synthesis with latent diffusion models, 2022.
- [20] C. Schuhmann, R. Vencu, R. Beaumont, R. Kaczmarczyk, C. Mullis, A. Katta, T. Coombes, J. Jitsev, and A. Komatsuaki. LAION-400M: open dataset of clip-filtered 400 million image-text pairs. *CoRR*, abs/2111.02114, 2021.
- [21] J. Sohl-Dickstein, E. A. Weiss, N. Maheswaranathan, and S. Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. *CoRR*, abs/1503.03585, 2015.
- [22] D. Soysaner. Attention mechanism in neural networks: where it comes and where it goes. *Neural Computing and Applications*, 34(16):13371–13385, may 2022.
- [23] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin. Attention is all you need. *CoRR*, abs/1706.03762, 2017.
- [24] B. Yu, P. Lu, R. Wang, W. Hu, X. Tan, W. Ye, S. Zhang, T. Qin, and T.-Y. Liu. Museformer: Transformer with fine- and coarse-grained attention for music generation, 2022.