# Cloud2Curve: Generation and Vectorization of Parametric Sketches

Ayan Das[1,2]    Yongxin Yang[1,2]    Timothy Hospedales[1,3]    Tao Xiang[1,2]    Yi-Zhe Song[1,2]

[1] SketchX, CVSSP, University of Surrey, United Kingdom
[2] iFlyTek-Surrey Joint Research Centre on Artificial Intelligence
[3] University of Edinburgh, United Kingdom
{a.das, yongxin.yang, t.xiang, y.song}@surrey.ac.uk, t.hospedales@ed.ac.uk

## Abstract

*Analysis of human sketches in deep learning has advanced immensely through the use of waypoint-sequences rather than raster-graphic representations. We further aim to model sketches as a sequence of low-dimensional parametric curves. To this end, we propose an inverse graphics framework capable of approximating a raster or waypoint based stroke encoded as a point-cloud with a variable-degree Bézier curve. Building on this module, we present Cloud2Curve, a generative model for scalable high-resolution vector sketches that can be trained end-to-end using point-cloud data alone. As a consequence, our model is also capable of deterministic vectorization which can map novel raster or waypoint based sketches to their corresponding high-resolution scalable Bézier equivalent. We evaluate the generation and vectorization capabilities of our model on Quick, Draw! and K-MNIST datasets.*

## 1. Introduction

The analysis of free-hand sketches using deep learning [40] has flourished over the past few years, with sketches now being well analysed from classification [43, 42] and retrieval [27, 12, 4] perspectives. Sketches for digital analysis have always been acquired in two primary modalities - *raster* (pixel grids) and *vector* (line segments). Raster sketches have mostly been the modality of choice for sketch recognition and retrieval [43, 27]. However, generative sketch models began to advance rapidly [16] after focusing on vector representations and generating sketches as sequences [7, 37] of waypoints/line segments, similarly to how humans sketch. As a happy byproduct, this paradigm leads to clean and blur-free image generation as opposed to direct raster-graphic generations [30]. Recent works have studied creativity in sketch generation [16], learning to sketch raster photo input images [36], learning efficient hu-
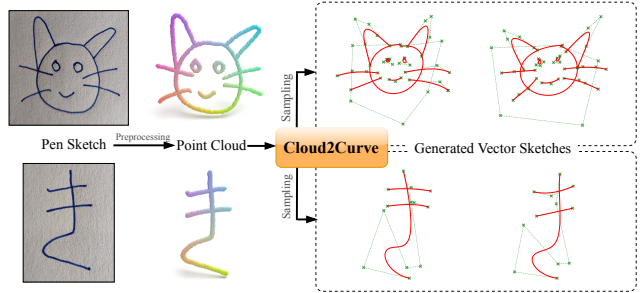


Figure 1. Cloud2Curve capability teaser. Top: Our trained model can vectorize a pen-on-paper sketch using scalable parametric curves. Bottom: Cloud2Curve can be trained on raster datasets such as KMNIST, where existing generative models cannot.

man sketching strategies [3], exploiting sketch generation for photo representation learning [39], and the interaction between sketch generation and language [18].

We present Cloud2Curve, a framework that advances the generative sketch modeling paradigm on two major axes: (i) by generating *parametric sketches*, i.e., a compositions of its constituent strokes as scalable parametric curves; and (ii) providing the ability to sample such sketches given point-cloud data, which is trivial to obtain either from raster-graphic or waypoint sequences. Altogether, our framework uniquely provides the ability to generate or deterministically vectorize scalable parametric sketches based on point clouds, as illustrated in Figure 1 and explained below.

First, we note that although existing frameworks like SketchRNN [16] and derivatives generate vector sketches, they do so via generation of a dense sequence of short straight line segments. Consequently, (i) the output sketches are not *spatially scalable* as required, e.g., for digital art applications, (ii) they struggle to generate *long* sketches due to the use of recurrent generators [28], and (iii) the generative process suffers from low interpretability compared to human sketching, where the human mental model relies more on composing sketches with smooth strokes. While there

have been some initial attempts at scalable vector sketch generative models [10], they were limited by the need for a two step training process, and to generate parametric Bézier curves of a fixed complexity, which poorly represent the diverse kinds of strokes in human sketches. In contrast, we introduce a machinery to dynamically determine the optimal complexity of each generated curve by introducing a continuous *degree* parameter; and our model can further be trained end-to-end.

Second, existing generative sketch frameworks require purely sequential data in order to train their sequence-to-sequence encoder-decoder modules. This necessitates the usage of specially collected sequential datasets like *Quick, Draw!* and not general purpose raster sketch datasets. In contrast, we introduce a framework that encodes point-cloud training data and decodes a sequence of parametric curves. We achieve this through an inverse-graphics [22, 33] approach of reconstructing training images/clouds by rendering its constituent strokes individually through a white-box Bézier decoder over several time-steps. To train this framework we compare reconstructed curves with the original segmented point cloud strokes with an Optimal Transport (OT) based objective. We show that such objectives, when coupled with appropriate regularizers, can lead to *controllable abstraction* of sketches.

To summarize our contributions: **1)** We introduce a novel formulation of Bézier curves with a continuous degree parameter that is automatically inferred for each individual stroke of a sketch. **2)** We develop Cloud2Curve, a generative model capable of training and inference on point cloud data to produce spatially scalable *parametric sketches*. **(3)** We demonstrate scalable parametric curve generation from point-clouds, using *Quick, Draw!* and a subset of K-MNIST [9] datasets.

## 2. Related Works

**Generative Models**    Generative models have been widely studied in machine learning literature [5] as a way of capturing complex data distributions. Probabilistic Graphical Models [20] were hard to scale with variational methods [6]. With the advent of deep learning, neural approaches based on Variational Autoencoder (VAE) [19] and Generative Adversarial Networks (GAN) [14] were more scalable and able to generate high-quality images. The general formulation of VAE has been adapted to numerous problem settings. [37] proposed the first model to encode sequential video data to a smooth latent space. Later, [7] showed an effective way to train sequential VAEs for generating natural language. SketchRNN [16] followed a sequence model very similar to [15] and learned a smooth latent space for generating novel sketches. [13] proposed a generative agent that learns to draw by exploring in the space of programs. Its

sample inefficiency was ameliorated by [44] through creating an environment model. More recently, the Transformer [38] has been used to model sketches [32] due to the permutation invariant nature of strokes.

**Parametric representation**    Although used heavily in computer graphics [34], parametric curves like Bézier, B-Splines, Hermite Splines have not been used much in mainstream Deep Learning. An early application of splines to model handwritten digits [31] used a density model around b-splines and learns the parameters from a point-cloud using log-likelihood. B-Splines have been used as stroke-segments while representing handwritten characters with probabilistic programs [23]. SPIRAL [13] is a generative agent that produces program primitives including cubic Bézier curve. The font generation model in [25] and more recently DeepSVG Icon generator [8] treats fonts/icons as a sequence of SVG primitives. However, this requires the ground-truth SVG primitives. In contrast, we take an inverse graphics approach that learns to render point-clouds using parametric curves – without any parametric curve ground-truth in the training pipeline. Stroke-wise embeddings are studied in [1], but this produces non-interpretable representations of each stroke, and still requires sequence data to train, unlike our inverse graphics approach. In summary, none of these methods can apply to raster data such as K-MNIST which we demonstrate here.

**Learning parametric curves**    The field of computer graphics [34] has seen tremendous use of parametric curves [11] in synthesizing graphics objects. However, parametric curves are still not popular in mainstream deep learning due to their usage of an extra latent parameter which is difficult to incorporate into a standard optimization setting. The majority of algorithms for fitting Bézier [35, 26] or B-Splines [24, 29, 45] are based on alternatively switching between optimizing control points and latent $t$ values which is computationally expensive, requires careful initialization and not suitable for pluging into larger computational graphs trained by backpropagation. Recently, BézierEncoder [10] was proposed as a fitting method for Bézier curves by means of inference on any arbitrary deep recurrent model. We go beyond this to also infer curve degree, and fundamentally generalize it for training on point-cloud/raster data.

## 3. Methodology

The mainstream sketch representation popularized by *Quick, Draw!*, encodes a sketch as an ordered list of $L$ waypoints $[(\mathbf{x}_i, \mathbf{q}_i)]_{i=1}^L$ where $\mathbf{x}_i \in \mathbb{R}^2$ is the $i^{th}$ waypoint and $\mathbf{q}_i = (q_i^{stroke}, q_i^{sketch})$ is a tuple containing two binary variables denoting stroke and sketch termination. State-of-the-art sketch generation models such as SketchRNN [16] directly use this data structure for modelling the probability of a given sketch as the product of probabilities of individ-
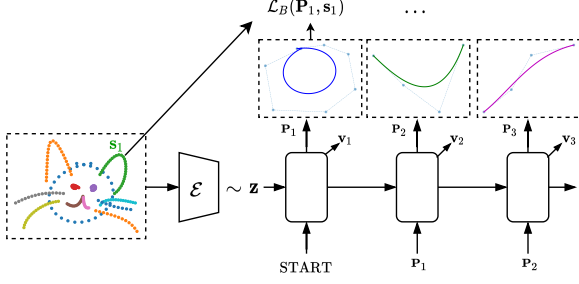
Figure 2. Overall diagram of the Cloud2Curve generative model.

ual waypoints given the previous waypoints

$$p_{sketchrnn}(\cdot) = \prod_{i=1}^{L} p(\mathbf{x}_i, \mathbf{q}_i | \mathbf{x}_{<i}, \mathbf{q}_{<i}; \theta) \qquad (1)$$

where the subscript $< i$ denotes the set of all indices before $i$. Our method treats sketches as a sequence of parametric curves, so we re-organize the data structure in terms of its strokes. We define a sketch $\mathcal{S} \triangleq [\mathbf{S}_k]_{k=1}^{K}$ as a sequence of its $K$ (may vary for each sketch) constituent strokes. Some recent methods substitute the original $\mathbf{S}_k$ with separately learned non-interpretable transformer-based embedding [1] or directly interpretable learned Bézier curves [10] with fixed degree. However, our method stands out due to the specific choice of variable degree Bézier curves as stroke embedding which is *end-to-end trainable* along with the generative model.

Furthermore, we simplify $\mathbf{S}_k$ by restructuring it into $(\mathbf{s}_k, \mathbf{v}_k)$ where $\mathbf{s}_k$ denotes a position-independent stroke and $\mathbf{v}_k$ is its position relative to an arbitrary but fixed point. This restructuring not only emphasizes the compositional relationship [1] within the strokes but also helps the generative model learn position invariant parametric curves.

## 3.1. Generative Model

Before defining our new variable-degree parametric stroke model in Section 3.2 and the training objective in Section 3.3, we introduce the overall generative model.

**Decoder/Generator** Unlike SketchRNN [16], but similar to BézierSketch [10] and CoSE [1], we model a sketch $\mathcal{S}$ autoregressively as a sequence of parametric strokes $\mathbf{P}_k$,

$$p(\mathcal{S}) = \prod_{k=1}^{K} p(\mathbf{P}_k, \mathbf{v}_k, \mathbf{q}_k | \mathbf{P}_{<k}, \mathbf{v}_{<k}, \mathbf{q}_{<k}; \theta) \qquad (2)$$

We use the binary random variable $\mathbf{q}_k$ to denote end-of-sketch as the usual way of terminating inference. Our model differs primarily in the fact that we model the density of the parametric stroke representation $\mathbf{P}_k$ at each step $k$. Since we do not use any pre-learned embedding as supervisory signal (unlike [10, 1]), we do not have ground-truth for $\mathbf{P}_k$

and hence can not train it by directly maximizing the likelihood in Eq. 2. We instead minimize an approximate version of the negative log-likelihood:

$$
\begin{aligned}
\mathcal{L}(\mathcal{S}) \approx - \sum_{k=1}^{K} & \Big[ \mathcal{L}_B(\widehat{\mathbf{P}}_k, \mathbf{s}_k) + \log p(\mathbf{v}_k | \widehat{\mathbf{P}}_{<k}, \mathbf{v}_{<k}) \\
& + \log p(\mathbf{q}_k | \widehat{\mathbf{P}}_{<k}, \mathbf{v}_{<k}) \Big], \\
& \text{with } \widehat{\mathbf{P}}_k \sim p(\mathbf{P}_k | \widehat{\mathbf{P}}_{<k}, \mathbf{v}_{<k})
\end{aligned}
\qquad (3)
$$

Instead of directly computing the log-likelihood of $\mathbf{P}_k$, we sample (re-parameterized) from the density and compute a downstream loss function $\mathcal{L}_B$ to act as a proxy. We describe the exact form of $\mathbf{P}_k$ and $\mathcal{L}_B(\cdot)$ in Sections 3.2 & 3.3.

**Encoder** We condition the generation with a global latent vector [16]. This is produced by a VAE-style [19] latent distribution whose parameters are computed using an encoder $\mathcal{E}_\theta$. A latent vector $\mathbf{z}$ is sampled as

$$\mathbf{z} | \mathcal{S} \sim \mathcal{N}(\mu, \boldsymbol{\Sigma}) \text{ with } [\mu, \boldsymbol{\Sigma}] = \mathcal{E}_\theta(\mathcal{S}) \qquad (4)$$

by means of the reparameterization trick [19]. $\mathbf{z}$ is then used to generate the mean and co-variance parameters at each step that define the distributions $p(\mathbf{P}_k | \cdot)$, $p(\mathbf{v}_k | \cdot)$ and $p(\mathbf{q}_k | \cdot)$ in Eq. 3. A high level diagram of the full architecture is shown in Fig. 2.

**Training** Given our encoder and decoder, training is conducted by optimising the following objective

$$\sum_{\mathcal{S} \sim \mathcal{D}} \mathcal{L}(\mathcal{S} | \mathbf{z}) + w_{KL} \cdot \text{KL} \left[ p_\theta(\mathbf{z} | \mathcal{S}) || \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}) \right] \qquad (5)$$

Depending on the nature of $\mathcal{E}_\theta$, we can have very different kind of generative models. One can use the usual SketchRNN-style [16] encoder. But we use a Transformer based set encoder $\mathcal{E}_\theta(\mathcal{S})$ that can parse a given sketch as a cloud of ink-points and produce a concise latent vector representation. In order to support learning parametric curve generation on point-cloud data, the remaining required components are a parametric curve model for $\mathbf{P}$ (Section 3.2) and the loss $\mathcal{L}_B(\mathbf{P}_k, \mathbf{s}_k)$ describing how well a parametric stroke $\mathbf{P}_k$ fits the relevant point subset $\mathbf{s}_k$, described in Section 3.3

## 3.2. Representing variable degree Bézier curves

The decoder generates the parameters of a Bézier curve representing a stroke at each step $k$. In this section we describe our new interpretable parametric curve representation $\mathbf{P}_k$. We formulate a flexible representation of $\mathbf{P}_k$ by defining it as a 2-tuple $(\mathcal{P}_k, r_k)$ comprised of: **1.** The parameters of a Bézier curve $\mathcal{P}_k \in \mathbb{R}^{(N+1) \times 2}$ where $N$ is the *maximum allowable* degree, **2.** A continuous variable $r_k \in [0, 1]$ that will be used to determine the *effective* degree of the Bézier

3

curve. We may drop the $k$ subscript when denoting an arbitrary time-step.

**Bézier curves** [34] are smooth and finite parametric curves used extensively in computer graphics. A degree $n$ Bézier curve is parameterized by $n + 1$ *control points*, and is usually modelled parametrically via an interpolation parameter $t \in [0, 1]$. A Bézier curve with control points $\mathcal{P}$ can be instantiated using a pre-specified set of $G$ interpolation points $[t_i]_{i=1}^{G}$ as

$$
C_{G \times 2} = \underbrace{\begin{bmatrix} \cdots \\ 1, t_i, t_i^2, \cdots, t_i^n \\ \cdots \end{bmatrix}}_{T_{G \times (n+1)}} \cdot M_{(n+1) \times (n+1)} \cdot \mathcal{P}_{(n+1) \times 2}
$$

(6)

where $G$ represents the granularity of rendering, which we treat as a hyperparameter. $T$ is the interpolation parameter matrix and $M$ is a matrix of Bernstein coefficients whose size and elements are dependent (only) on the degree $n$. For convenience, we will denote them as $M^{(n)}$. We next address how to use a continuous variable to induce an effective degree on $\mathcal{P}$.

**Soft Binning** [41] has been introduced as a *differentiable* way of binning a given real number into a predefined set of buckets. A real number $r$ needs to be tested against $n$ *cut points* in order to assign a one-hot $(n + 1)$-way categorical (or a continually relaxed) vector whose entries correspond to each of $n + 1$ buckets/bins. Please refer to [41] for the detailed formulation of *Soft Binning*. We interpret each bucket as an *effective degree* of a Bézier curve with its control points in $\mathcal{P}$. For our problem, we fix the cut points according to the maximum allowable degree $N$ as $\mathbf{U} = [i/N]_{i=0}^{N}$. This allows us to transform the continuous variable $r$ into a soft-categorical vector with $N + 2$ components. For practical benefit, we constraint the quantity $r$ to fall within the unit range of $[0, 1]$ by parameterizing it with an unconstrained variable $r' \in [-\infty, +\infty]$ as $r = \text{Sigmoid}(r')$. With this added constraint, $r$ can only fall into $N$ buckets (avoiding the first and last open buckets), each of which may denote a Bézier curve with $2, 3, \cdots, N + 1$ control points. Such design choice allows us to avoid representing a Bézier curve with 1 control point (which is invalid by definition) for any value of $r'$.

We define two quantities: **1)** A *degree selector* $\mathbf{R}(r)$ whose element $\mathbf{R}_i$ is 1 iff $r$ falls into the bin designated for degree $i$; **2)** A *control points selector* $\overline{\mathbf{R}}(r)$ defined as the *reversed cumulative summation* of $\mathbf{R}(r)$. Please refer to Fig. 3 for an illustration of a variable degree Bézier curve.

**Variable degree** We can now augment $\mathcal{P}$ with $r$ to create variable degree Bézier representation. We mask interpolation parameters $T$ and control points $\mathcal{P}$ using the *control*
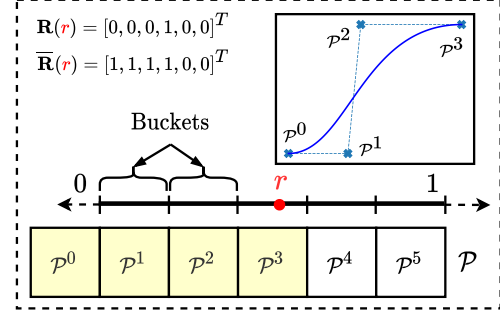


Figure 3. Visualization of the variable degree Bézier curve parameterized by its continuous degree parameter $r$.

*point selector* as

$$
\widehat{T}(r) = \underbrace{\begin{bmatrix} \cdots \\ 1, t_i, t_i^2, \cdots, t_i^N \\ \cdots \end{bmatrix}}_{T_{G \times (N+1)}} \odot \begin{bmatrix} \overline{\mathbf{R}}^T(r) \\ \vdots \\ \overline{\mathbf{R}}^T(r) \end{bmatrix}_{G \times (N+1)}
$$

$$
\widehat{\mathcal{P}}(r) = \mathcal{P} \odot \left[ \overline{\mathbf{R}}(r), \overline{\mathbf{R}}(r) \right]
$$

For $M$, we need to select the correct one from the set $\{M^{(n)}\}_{n=0}^{N}$ according to the value of $r$. We accomplish this by first defining a 3D tensor $\mathcal{M} \in \mathbb{R}^{(N+1) \times (N+1) \times (N+1)}$ where

$$
\mathcal{M}[i, \cdots] = \begin{bmatrix} M^{(i)} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}_{(N+1) \times (N+1)}
$$

The $\mathbf{0}$s denote appropriately sized zero matrices used as fillers. We can then compute $\widehat{M}$ using the *degree selector* as

$$
\widehat{M}(r) = \mathbf{R}^T(r) \cdot \mathcal{M}
$$

With all three components augmented with masks, we can write the variable degree version of Eq. 6 as a function of the degree parameter $r$ as

$$
\widehat{C}(r)_{G \times 2} = \widehat{T}(r) \cdot \widehat{M}(r) \cdot \widehat{\mathcal{P}}(r)
$$

(7)

### 3.3. Learning and Inference

Given the generative model described in Section 3.1 and our new curve representation in Section 3.2, we now describe how to train on point-cloud data, and how to use our trained model to vectorize new point cloud inputs.

**Bézier Loss: Cloud** The final component required by our generative model (Eq. 3) is a loss $\mathcal{L}_B(\mathbf{P}, \mathbf{s})$ to measure the similarity between a point cloud based stroke $\mathbf{s}$ and the curve $\mathbf{P}$. We discard the sequential information in the set $\widehat{C}$ and can compute any Optimal Transport (OT) based loss like EMD (Earth mover's distance) or Wasserstein Distance

[2]. The specific distance we used in our experiments is the Sliced Wasserstein Distance (SWD) [21]:

$$\mathcal{L}_B(\mathbf{P}, \mathbf{s}) = \text{SWD}(\mathbf{P}, \mathbf{s}) \tag{8}$$

Since OT-based losses are theoretically designed to measure the difference between two distributions, it is necessary to ensure the cardinality of the sets (either $\mathbf{P}$ or $\mathbf{s}$) are sufficiently high. We use a large enough $G$ for instantiating the Bézier curve and also densely resample $\mathbf{s}$ to the same granularity.

**Bézier Loss: Sequence** If *optional* sequence information for $\mathbf{s}$ is available, we can compute a point-to-point MSE loss:

$$\mathcal{L}_B(\mathbf{P}, \mathbf{s}) = \sum_{g=1}^{G} ||\widehat{C}^g(r) - \mathbf{s}^g||_2.$$

between each point $s^g$ on the curve and each interpolation point $\widehat{C}^g$ on the rendered curve $\mathbf{P}$ (Eq. 7).

**Regularization** The purpose of introducing a variable degree Bézier curve formulation is to provide the model with flexibility to encode strokes with perfect fit. To avoid overfitting strokes by using a complex curve to fit a simple stroke, the learning phase should be provided with incentive to reduce the degree whenever possible. A simple regularizer on our degree variable $r$ could achieve this:

$$\widehat{\mathcal{L}}_B(\mathbf{P}, \mathbf{s}) = \mathcal{L}_B(\mathbf{P}, \mathbf{s}) + \lambda_d \cdot r \tag{9}$$

Apart from this, we also added another regularizer to reduce the level of overfitting given a degree. Overfitting in learning Bézier curve can occur when control points can move anywhere during the optimization. Following [10], we add another term to the loss function to penalize the consecutive control points moving away from each other

$$\widehat{\mathcal{L}}_B(\mathbf{P}, \mathbf{s}) = \mathcal{L}_B(\mathbf{P}, \mathbf{s}) + \lambda_d \cdot r \\ + \lambda_c \cdot \left( \sum_{i=0}^{N} ||\mathcal{P}^{i+1} - \mathcal{P}^i||_2 \right) \odot \overline{\mathbf{R}}(r) \tag{10}$$

where we have masked out control points in $\mathcal{P}$ that are not meaningful given the degree value $r$.

**Implementation Details** Our training objective (Eq 5) encodes whole images but fits one parametric curve at a time to the set of points corresponding to a stroke (Eq 10). In principle all the strokes could be emitted by the generator, and then compute the loss between the full parametric sketch and the full point cloud. However for stability of optimisation, and limiting the cost of OT computation between curves and cloud, we proceed stroke-wise. To do this for genuine raster data, we pre-process the input pointcloud with 2D clustering to segment into strokes, and then iterate over the strokes in random order to train the model. Note that this priveleged information is only required during training, after training we can vectorize an unsegmented raster image into parametric curves, as shown in Figure 1.

**Inference: Generation & Vectorization** Given our trained model, we can use it for conditional generation. Given a sketch $\mathcal{S}$ as pointcloud, we simply sample a latent vector $\widehat{\mathbf{z}} \sim \mathcal{N}(\mathcal{E}_{\theta*}(\mathcal{S}))$ following Eq. 4 and use it to construct the parameters of the distributions $p(\mathbf{P}_k|\cdot)$, $p(\mathbf{v}_k|\cdot)$ and $p(\mathbf{q}_k|\cdot)$. We further sample

$$\widehat{\mathbf{P}}_k, \widehat{\mathbf{v}}_k, \widehat{\mathbf{q}}_k \sim p(\mathbf{P}_k|\cdot) \cdot p(\mathbf{v}_k|\cdot) \cdot p(\mathbf{q}_k|\cdot)$$

iteratively at each time-step and stop only when $\widehat{\mathbf{q}}_k$ is in *end-of-sequence* state. To visualize, we simply render all the $(\widehat{\mathbf{P}}_k, \widehat{\mathbf{v}}_k)$ pair on a canvas.

We can *vectorize* a given sketch $\mathcal{S}$ deterministically by following a similar procedure as generation but with discarding the source of stochasticity while sampling. We can simply assign all the co-variance parameters of $p(\mathbf{P}_k|\cdot)$, $p(\mathbf{v}_k|\cdot)$ and $p(\mathbf{q}_k|\cdot)$ to zero.

## 4. Experiments

**Datasets** *Quick, Draw!* [16] is the largest free-hand sketch dataset available till date. *Quick, Draw!* is created by collecting drawings from a fixed set of categories drawn under a game played by millions all over the world. Although *Quick, Draw!* is collected on a vast array of digital devices like smartphone, tablets etc., the data acquisition technique is kept uniform. The sketches are collected as a series of 2D waypoints along the trajectory of ink flow. To demonstrate our model, we use *Quick, Draw!* but discard it's waypoint sequence information, using the waypoints as pointcloud data. We use *Quick, Draw!* stroke-level segmentation (given by pen-up and pen-down indicators) as priveleged information during training. To demonstrate our model's ability to train on pure raster data, for which neither point-sequence nor stroke-sequence information is available, we also validate our model on a few classes $(0, 1, 8, 9)$ of K-MNIST [9]. We extract a point-cloud representation from K-MNIST by simple binarization and thinning. For training we segment into strokes by Spectral Clustering.

### 4.1. Variable-degree Bézier curve

We first validate the ability of our new curve representation to fit variable-degree Bézier curves to isolated strokes represented as point clouds.

**Setup** We collected few strokes from the sketches of *Quick, Draw!* [16] dataset. Since each stroke may have dif-
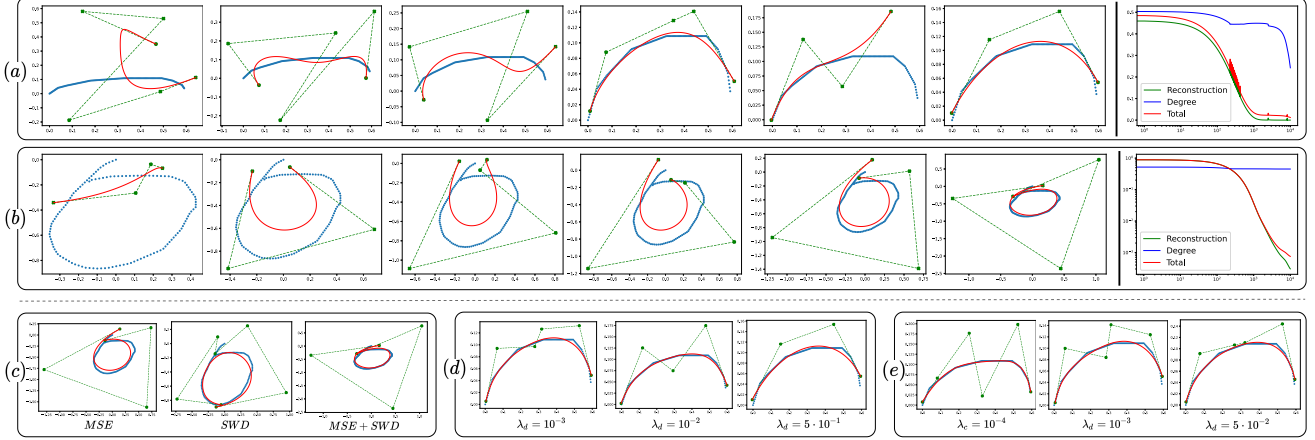
Figure 4. Visualization of fitting a variable degree Bézier curve (red) to individual strokes represented as point-clouds (blue). (a & b) Two examples with corresponding training loss components on the right. The training iteration increases from left to right. (c) Different choices of loss (section 3.3). (d) Different choices of degree regularizer $\lambda_d$. (e) Different choices of control point regularizer $\lambda_c$.

ferent number of $2D$ waypoints, we densely resample them uniformly with a fixed granularity of $G_0$. Hence, given a waypoint-based stroke $\mathbf{s} \in \mathbb{R}^{G_0 \times 2}$, we directly optimize the Bézier curve parameters by means of Eq. 10

$$\mathcal{P}^*, r'^* = \underset{\mathcal{P}, r'}{\arg\min} \widehat{\mathcal{L}}_B(\mathcal{P}, r, \mathbf{s}), \text{ with}$$
$$r = \text{Sigmoid}(r'), \text{ and } (\mathcal{P}, r) \triangleq \mathbf{P} \quad (11)$$

The gradients of the loss w.r.t parameters, i.e. $\left( \frac{\partial \mathcal{L}_B}{\partial \mathcal{P}}, \frac{\partial \mathcal{L}_B}{\partial r'} \right)$ are computed simply by backpropagation and updated using any SGD-based algorithm.

Additionally, we experimented with point-to-point MSE loss as described in Section 3.3 in cases where sequential information is available. Usage of privileged information helps learning the Bézier parameters quickly. However, the usage of *only* MSE loss degrades the fitting quality. We noticed that the dense uniform resampling of strokes $\mathbf{s}$ do not have a proper point-to-point alignment with an instantiated Bézier curve with uniform set of $t$-values of same granularity $G$. An instantiated Bézier curve tends to have dense distribution of points in places of bending and sparse otherwise, while the point cloud data derived from a waypoint sequence does not. Using the Sliced Wasserstein Distance (SWD) [21] with MSE loss proved most effective with the sequential information in the MSE loss helping speed of convergence, and the SWD loss helping quality of fit.

**Results** We show qualitative results of learning variable degree Bézier curves from point-cloud based strokes $\mathbf{s}$. Figure 4(a, b) shows two examples of learning the Bézier control points $\mathcal{P}$ and degree parameter $r'$. To interpret the learned value of $r$, we run binning on it using the predefined cut-points $\mathbf{U}$ in Section 3.2 and retrieve the degree $n$. In Fig. 4(a), we see the degree reduces from $n = 5$

to $n = 4$ and then $n = 3$, since the curve is simple. Similarly, Fig. 4(b) shows an increase in degree due to the higher complexity stroke. The last figure in each example shows each loss component over iterations. We used $\lambda_d = 10^{-3}$, $\lambda_c = 5 \cdot 10^{-2}$, $G_0 = 128$ and maximum degree $N = 6$ for this experiment. Fig. 4(c) shows the qualitative difference between learning with SWD, MSE alone and both combined. We see that SWD tends to ignore small details but preserves the overall structure. Fig. 4(d, e) also shows the effect of both regularizer strengths ($\lambda_c$ and $\lambda_d$). It is evident from the figure that $\lambda_d$ is responsible for pressuring the reduction of degree in fitting, whereas $\lambda_c$ reduces overfitting by keeping the control points close to each other.

### 4.2. Generation and & Vectorization Model

**Setup** We pre-process the sketches in *Quick, Draw!* and K-MNIST into a sequence of position-independent strokes and their starting points $\mathcal{S} = \left[ (\mathbf{s}_k, \mathbf{v}_k) \right]_{k=1}^{K}$. To stabilize the training dynamics, we also center the whole sketch and scale it down numerically to fit inside a unit circle. As part of data augmentation, we performed the following: (i) We made sure the strokes do not have too sharp bends. We split a stroke into multiple strokes from a point with high curvature or when it is too long. Such transformation alleviates the problem of learning overly smooth representations up to a great extent. (ii) We added standard Gaussian noise to each $2D$ point in the stroke.

**Implementation Details** Although our decoder model in Eq. 2 is very generic, we chose to use a standard RNN. Alternatives such as Transformer [32] could also be used.

The consequence of using position independent strokes is that the predicted parametric curve must also be position independent. A trick to guarantee that $\mathcal{P}^0 = [0, 0]^T$ is to
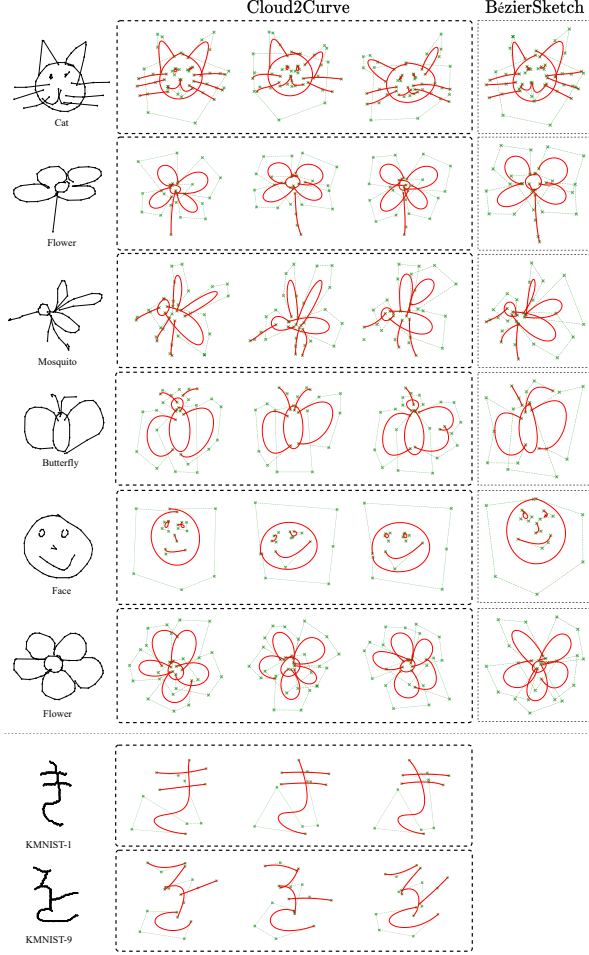
6

Figure 5. Qualitative results for conditional generation on *Quick, Draw!* and a limited subset of KMNIST. Three vector sketches are generated by means of sampling from Cloud2Curve. For a qualitative comparison, we also provide one sample from BézierSketch [10] for the same input sketch. Due to unavailability of purely sequential information, BézierSketch cannot be trained on KMNIST.
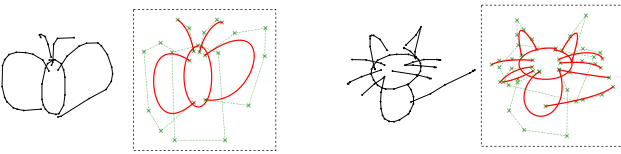


Figure 6. Deterministically vectorized sketches from the model trained with $\lambda_d = 10^{-3}$.

only predict $N$ control points for $i = 1 \rightarrow N$ and explicitly fix $\mathcal{P}^0 := [0, 0]^T$.

We set the family of distributions $p(\mathbf{P}_k | \cdot; \theta)$ and $p(\mathbf{v}_k | \cdot)$ at each time-step as isotropic Gaussian with their mean and std vector predicted by our encoder $\mathcal{E}$. The reason we chose to not use GMM here is because of its difficulty to reparameterize (as required in Eq. 3). However, the loss of expressiveness is compensated by increasing model capacity.

Initiation of inference requires a special input state. Since $(\mathcal{P}, r, \mathbf{v}) = (\mathbf{0}, 0, \mathbf{0})$ is a semantically invalid state, it can be used as a special token to kickstart the generation and also use in teacher-forcing while training.

As mentioned earlier, we use a set transformer as the encoder $\mathcal{E}_\theta$. Specifically, we use the strategy of [32] to compute a compact latent distribution by simply applying a learnable self-attention layer on the encoder features. Note that, unlike the decoding side, we never use stroke level segmentation for the encoder. Thus we can perform inference without stroke-segmented sketches.

For the decoder, we used a 2-layer LSTM with hidden vectors of size $1024D$. We fixed the value of maximum allowable degree to $N = 9$, therefore the Bézier curves can have at max 10 control points which is more than enough to represent fairly complex geometry. For computing the Bézier loss $\mathcal{L}_B$, we used a granularity of $G = 128$. For the transformer-based encoder model, we used a $512D$ transformer with 8 heads and 4 layers. For stable training, we follow the usual KL annealing trick [7, 16] by varying $w_{KL}$ from 0 to 1 over epochs.

**Qualitative Results of Sketch Generation and Vectorization** We trained one model for each class from both *Quick, Draw!* and K-MNIST datasets. Since K-MNIST dataset is fairly small, we do not train a model from scratch. Instead, we finetune on a model pre-trained on *Quick, Draw!*. Fig. 5 shows qualitative results in terms of conditional generation from our generative sketch model. We also show qualitative results for deterministic raster-to-curve vectorizations in Fig. 6. Both results are computed following the procedures described in 3.3. Trained with two categories, the model can also interpolate between two given sketches [16] by interpolating on the latent space (refer to Fig. 7 for an example).

**Cloud2Curve Generates Compact Sketches** One important technical benefit of generating sketches using Bézier curves is that the generated sketches are shorter in terms of number of points to be stored. Moreover, short representation of sketches lead to less complex downstream models like classification and retrieval. The length histograms both at stroke-level (points per stroke) and sketch-
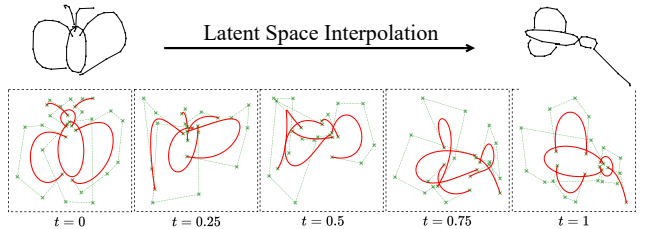


Figure 7. Interpolating between two sketches by walking on the latent space as: $\mathbf{z} = \mathbf{z}_{butterfly} \cdot (1 - t) + \mathbf{z}_{mosquito} \cdot t$
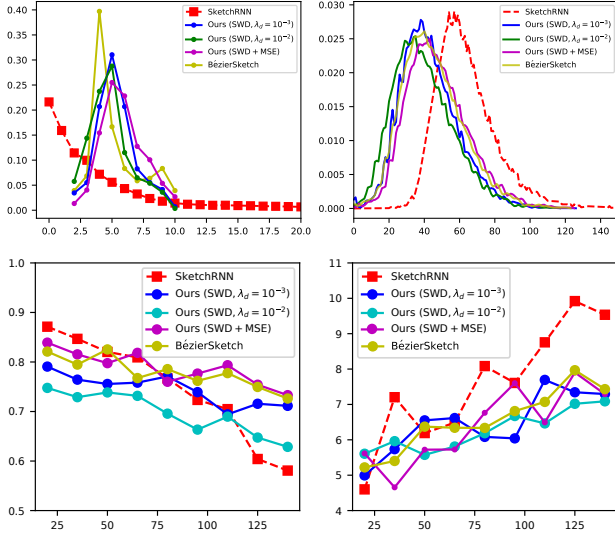
Figure 8. Comparison between Cloud2Curve and other generative models in terms of (Top left, a) stroke-histogram, (Top right, b) sketch-histogram, (Bottom left, c) Classification accuracy-vs-length and (Bottom right, d) FID-vs-length for the generated samples. Cloud2Curve generates more compact sketches and scales better to higher lengths in terms of FID and recognition accuracy.
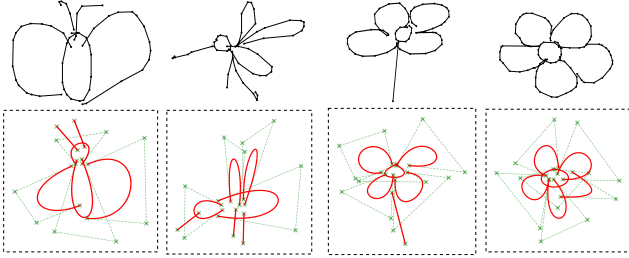


Figure 9. Input sketches and conditionally generated samples from a model trained with $\lambda_d = 10^{-2}$.

level (points-per-sketch) of generated sketches are shown in Figure. 8(a,b) demonstrate this point for Cloud2Curve vs SketchRNN and BézierSketch. We notice that the number of points per stroke is heavily concentrated around 5 (and bounded by $2 \rightarrow 10$) for all our models but SketchRNN has a much larger spread. As a direct consequence, the average number of points in an entire sketch is also smaller (by a margin of ~20). We also show that increasing $\lambda_d$ reduces the lengths even further, thereby achieving more *abstract* sketches (Refer to Fig. 9 for qualitative examples). Using MSE loss has an effect of increasing the length because it tries to fit the data more precisely considering small artifact.

**Quantitative Analysis on Sketch Generation Quality**
To confirm that our generative model samples are semantically plausible, we rendered them and classified them using a CNN classifier [42] trained on real sketches from *Quick, Draw!*. The results in Figure 8(c) show that our parametric

curve sketches are indeed accurately recognizable by state-of-the-art image classifiers even when longer in length.

To further assess generation quality, the FID score [17] is also computed. We compute a modified FID score using rasterized samples after projecting them down to the activations of penultimate layer of a pre-trained Sketch-a-Net [42] classifier. Figure. 8(d) shows that our model can generate long sketches better than sequential models like SketchRNN in terms of FID (lower is better).

**Quantitative Analysis on Sketch Vectorization Quality**
We finally validate our model as a raster-to-curve vectorizer. We vectorize testing sketches and calculate the test loss between vectorized sketch and the available stroke-segmented ground-truth to evaluate the quality of sketch curve fitting. Furthermore, to evaluate the vectorization quality from a semantic perspective, we calculate classification accuracy of the generated vector sketches by first rasterizing and then classifying them with a pre-trained Sketch-a-Net [42]. The results are shown in Table 1. We notice that although the usage of sequential information (MSE loss) helps reducing the SWD error a lot, the perceptual quality (classification score) remains fairly similar. We also evaluated BézierSketch [10] on *Quick, Draw!*, with and without using stroke sequence information, which correspond to an upper bound and a baseline respectively to our contribution of cloud-based sketch-rendering. The results show that Cloud2Curve almost matches BézierSketch with full sequence supervision, and is significantly better than BézierSketch using raw cloud data (and random sequence assignment).

## 5. Conclusion

In this paper, we introduced a model capable of generating scaleable vector-graphic sketches using parametric curves – and crucially it is able to do so by training on point cloud data, thus being widely applicable to general raster image datasets such as K-MNIST. Our framework provides accurate and flexible fitting due to the ability to chose curve complexity independently for each stroke. Once trained, our architecture can also be used to vectorize raster sketches into flexible parametric curve representations. In future work we will generalize our parametric stroke model from overly smooth Bézier curves to more general parametric curves such as B-splines or Hermite splines.

|  | QuickDraw | | KMNIST | |
|---|---|---|---|---|
|  | Classif. acc. | Test Loss | Classif. acc. | Test Loss |
| SWD, $\lambda_d = 10^{-3}$ | 0.80 | 0.00123 | 0.82 | $8.4 \cdot 10^{-3}$ |
| SWD, $\lambda_d = 10^{-2}$ | 0.69 | 0.02850 | 0.73 | $7.6 \cdot 10^{-2}$ |
| SWD + MSE | 0.84 | 0.00034 | 0.85 | $2.1 \cdot 10^{-4}$ |
| BézierSketch (Seq.) [10] | 0.86 | 0.00012 | NA | NA |
| BézierSketch (Cloud) [10] | 0.41 | 0.2572 | NA | NA |
| Real data | 0.92 | NA | 0.89 | NA |

Table 1. Quantitative validation of the vectorization model.

# References

[1] Emre Aksan, Thomas Deselaers, Andrea Tagliasacchi, and Otmar Hilliges. Cose: Compositional stroke embeddings. *NeurIPS*, 2020. 2, 3

[2] Martin Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein generative adversarial networks. In *PMLR*, pages 214–223, 2017. 5

[3] Ayan Kumar Bhunia, Ayan Das, Umar Riaz Muhammad, Yongxin Yang, Timothy M. Hospedales, Tao Xiang, Yulia Gryaditskaya, and Yi-Zhe Song. Pixelor: A competitive sketching ai agent. so you think you can beat me? *SIGGRAPH Asia*, 2020. 1

[4] Ayan Kumar Bhunia, Yongxin Yang, Timothy M. Hospedales, Tao Xiang, and Yi-Zhe Song. Sketch less for more: On-the-fly fine-grained sketch based image retrieval. In *CVPR*, 2020. 1

[5] Christopher M Bishop. *Pattern recognition and machine learning*. springer, 2006. 2

[6] David M Blei, Alp Kucukelbir, and Jon D McAuliffe. Variational inference: A review for statisticians. *Journal of the American statistical Association*, 112(518):859–877, 2017. 2

[7] Samuel R. Bowman, Luke Vilnis, Oriol Vinyals, Andrew Dai, Rafal Jozefowicz, and Samy Bengio. Generating sentences from a continuous space. In *CoNLL*, 2016. 1, 2, 7

[8] Alexandre Carlier, Martin Danelljan, Alexandre Alahi, and Radu Timofte. Deepsvg: A hierarchical generative network for vector graphics animation. *NeurIPS*, 2020. 2

[9] Tarin Clanuwat, Mikel Bober-Irizar, Asanobu Kitamoto, Alex Lamb, Kazuaki Yamamoto, and David Ha. Deep learning for classical japanese literature. *CoRR*, abs/1812.01718, 2018. 2, 5

[10] Ayan Das, Yongxin Yang, Timothy Hospedales, Tao Xiang, and Yi-Zhe Song. Béziersketch: A generative model for scalable vector sketches. In *ECCV*, 2020. 2, 3, 5, 7, 8

[11] Carl De Boor, Carl De Boor, Etats-Unis Mathématicien, Carl De Boor, and Carl De Boor. *A practical guide to splines*, volume 27. Springer-Verlag New York, 1978. 2

[12] Sounak Dey, Pau Riba, Anjan Dutta, Josep Llados, and Yi-Zhe Song. Doodle to search: Practical zero-shot sketch-based image retrieval. In *CVPR*, 2019. 1

[13] Yaroslav Ganin, Tejas Kulkarni, Igor Babuschkin, S. M. Ali Eslami, and Oriol Vinyals. Synthesizing programs for images using reinforced adversarial learning. In *ICML*, 2018. 2

[14] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *NIPS*, 2014. 2

[15] Alex Graves. Generating sequences with recurrent neural networks. *arXiv preprint arXiv:1308.0850*, 2013. 2

[16] David Ha and Douglas Eck. A neural representation of sketch drawings. In *ICLR*, 2018. 1, 2, 3, 5, 7

[17] 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 *NIPS*, 2017. 8

[18] Forrest Huang, Eldon Schoop, David Ha, and John F. Canny. Scones: towards conversational authoring of sketches. In *IUI*, 2020. 1

[19] Diederik P. Kingma and Max Welling. Auto-encoding variational bayes. *ICLR*, 2014. 2, 3

[20] Daphne Koller and Nir Friedman. *Probabilistic graphical models: principles and techniques*. MIT press, 2009. 2

[21] Soheil Kolouri, Phillip E Pope, Charles E Martin, and Gustavo K Rohde. Sliced wasserstein auto-encoders. In *ICLR*, 2018. 5, 6

[22] Tejas D. Kulkarni, Will Whitney, Pushmeet Kohli, and Joshua B. Tenenbaum. Deep convolutional inverse graphics network. *NIPS*, 2015. 2

[23] Brenden M. Lake, Ruslan Salakhutdinov, and Joshua B. Tenenbaum. Human-level concept learning through probabilistic program induction. *Science*, 350(6266):1332–1338, 2015. 2

[24] Yang Liu and Wenping Wang. A revisit to least squares orthogonal distance fitting of parametric curves and surfaces. In *GMP*, 2008. 2

[25] Raphael Gontijo Lopes, David Ha, Douglas Eck, and Jonathon Shlens. A learned representation for scalable vector graphics. In *ICCV*, 2019. 2

[26] Asif Masood and Sidra Ejaz. An efficient algorithm for robust curve fitting using cubic bezier curves. In *ICIC*, 2010. 2

[27] Kaiyue Pang, Ke Li, Yongxin Yang, Honggang Zhang, Timothy M. Hospedales, Tao Xiang, and Yi-Zhe Song. Generalising fine-grained sketch-based image retrieval. In *CVPR*, 2019. 1

[28] Razvan Pascanu, Tomas Mikolov, and Yoshua Bengio. On the difficulty of training recurrent neural networks. In *ICML*, 2013. 1

[29] Michael Plass and Maureen Stone. Curve-fitting with piecewise parametric cubics. In *SIGGRAPH*, 1983. 2

[30] Alec Radford, Luke Metz, and Soumith Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. In *ICLR*, 2016. 1

[31] Michael Revow, Christopher K. I. Williams, and Geoffrey E. Hinton. Using generative models for handwritten digit recognition. *IEEE T-PAMI*, 18(6):592–606, 1996. 2

[32] Leo Sampaio Ferraz Ribeiro, Tu Bui, John Collomosse, and Moacir Ponti. Sketchformer: Transformer-based representation for sketched structure. In *CVPR*, 2020. 2, 6, 7

[33] L. Romaszko, C. K. I. Williams, P. Moreno, and P. Kohli. Vision-as-inverse-graphics: Obtaining a rich 3d explanation of a scene from a single image. In *ICCV Workshops*, 2017. 2

[34] David Salomon. *Curves and surfaces for computer graphics*. Springer Science & Business Media, 2007. 2, 4

[35] Lejun Shao and Hao Zhou. Curve fitting with bezier cubics. *Graphical models and image processing*, 58(3):223–232, 1996. 2

[36] Jifei Song, Kaiyue Pang, Yi-Zhe Song, Tao Xiang, and Timothy M. Hospedales. Learning to sketch with shortcut cycle consistency. In *CVPR*, 2018. 1

[37] Nitish Srivastava, Elman Mansimov, and Ruslan Salakhudinov. Unsupervised learning of video representations using lstms. In *ICML*, 2015. 1, 2

[38] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszko-reit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NIPS*, 2017. 2

[39] Alexander Wang, Mengye Ren, and Richard Zemel. Sketchembednet: Learning novel concepts by imitating drawings. *arXiv preprint arXiv:2009.04806*, 2020. 1

[40] Peng Xu. Deep learning for free-hand sketch: A survey. *arXiv preprint arXiv:2001.02600*, 2020. 1

[41] Yongxin Yang, Irene Garcia Morillo, and Timothy M Hospedales. Deep neural decision trees. *arXiv preprint arXiv:1806.06988*, 2018. 4

[42] Qian Yu, Yongxin Yang, Feng Liu, Yi-Zhe Song, Tao Xiang, and Timothy Hospedales. Sketch-a-net: A deep neural network that beats humans. *IJCV*, 122:411–425, 2017. 1, 8

[43] Qian Yu, Yongxin Yang, Yi-Zhe Song, Tao Xiang, and Timothy Hospedales. Sketch-a-net that beats humans. In *BMVC*, 2015. 1

[44] Ningyuan Zheng, Yifan Jiang, and Dingjiang Huang. Strokenet: A neural painting environment. In *ICLR*, 2019. 2

[45] Wenni Zheng, Pengbo Bo, Yang Liu, and Wenping Wang. Fast b-spline curve fitting by l-bfgs. *CAGD*, 2012. 2