# Collaborative Neural Painting

Nicola Dall'Asen[a,b,**], Willi Menapace[a], Elia Peruzzo[a], Enver Sangineto[c], Yiming Wang[d], Elisa Ricci[a,d]

[a]*University of Trento, Italy*
[b]*University of Pisa, Italy*
[c]*University of Modena and Reggio Emilia, Italy*
[d]*Fondazione Bruno Kessler, Trento, Italy*

## ABSTRACT

The process of painting fosters creativity and rational planning. However, existing generative AI mostly focuses on producing visually pleasant artworks, without emphasizing the painting process. We introduce a novel task, *Collaborative Neural Painting (CNP)*, to facilitate collaborative art painting generation between users and agents. Given any number of user-input *brushstrokes* as the context or just the desired *object class*, CNP should produce a sequence of strokes supporting the completion of a coherent painting. Importantly, the process can be gradual and iterative, so allowing users' modifications at any phase until the completion. Moreover, we propose to solve this task using a painting representation based on a sequence of parametrized strokes, which makes it easy both editing and composition operations. These parametrized strokes are processed by a Transformer-based architecture with a novel attention mechanism to model the relationship between the input strokes and the strokes to complete. We also propose a new masking scheme to reflect the interactive nature of CNP and adopt diffusion models as the basic learning process for its effectiveness and diversity in the generative field. Finally, to develop and validate methods on the novel task, we introduce a new dataset of painted objects and an evaluation protocol to benchmark CNP both quantitatively and qualitatively. We demonstrate the effectiveness of our approach and the potential of the CNP task as a promising avenue for future research. Project page and code: this https URL.

## 1. Introduction

In recent years we have experienced an explosive growth of *AI Art*, empowering users with the possibility to generate images and other media content given various forms of conditioning, such as text (Saharia et al., 2022; Rombach et al., 2022) or semantic maps (Zeng et al., 2022; Avrahami et al., 2023; Zhang and Agrawala, 2023). Generative AI art in the visual domain has led to astonishing results for image synthesis tasks (Nichol et al., 2022; Ramesh et al., 2022; Saharia et al., 2022; Rombach et al., 2022). It typically operates on the pixel space, both for content generation and editing, where users can modify the image by indicating the editing areas or describing the desired modifications (Nichol et al., 2022; Meng et al., 2022).

However, this process of art generation and editing is fundamentally different from how humans create art, of which painting is a primary example (Nakano, 2019). When painting, humans reason on individual *brushstrokes* rather than *pixels*, and creativity is fostered throughout the stroke design and their compositional planning. While several works have been proposed which employ a brushstroke formulation, they mainly focus on producing the entire painting given a reference image or artwork (Liu et al., 2021b; Zou et al., 2021; Kotovenko et al., 2021), without involving humans in the generation process. This limits the degree to which such methods promote creativity and empathy (Gerry, 2017; Pelowski et al., 2017) that are particularly important for enhancing educational and pedagogical development in children (Beh-Pajooh et al., 2018) or rehabilitation (Zhang et al., 2021). More recently, (Peruzzo et al., 2023) proposes an interactive formulation of neural painting in which the model suggests the next strokes to paint according to a given *reference* image. The reference-guided interaction

---
[**]Corresponding author:
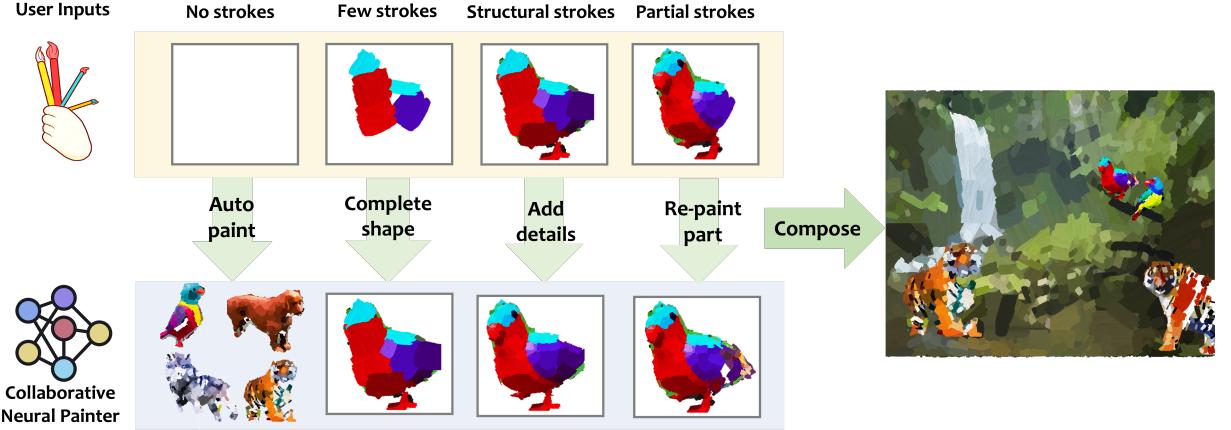*e-mail:* `nicola.dallasen@unitn.it` (Nicola Dall'Asen)

**Fig. 1. The proposed Collaborative Neural Painting task envisions a collaborative procedure in which users produce and compose artworks iteratively interacting with a Neural Painter. This interaction includes auto-painting without user input, and assistive painting/editing at any granularity level. We solve this generation task using a vectorized stroke parametrization whose joint distribution is learned using diffusion models.**

might be useful for the learning purpose, however, we argue that such setup may limit the fostering of creativity during painting, as the generation is bounded to the given reference picture.

In this work, we propose the novel task of **Collaborative Neural Painting** (CNP) which aims to enable a collaborative creation process of high-quality paintings that promotes active engagement of human users for creative generation. We design the interactive process based on brushstrokes that the human users can provide as context at any phase of their painting. This is different from previous works (Peruzzo et al., 2023; Zou et al., 2021; Liu et al., 2021b), as we do not require any reference image and the user can specify any arbitrary number of strokes. The agent, *i.e.*, the collaborative neural painter, subsequently reasons on the context strokes and produces new strokes in the painting to match the user's painting intentions *e.g.* completing the rough shape of the painting, adding details or regenerating a missing part, as shown in Figure 1.

CNP is a non-trivial task. It requires understanding the geometric properties of the strokes and how they translate to the final painting to generate rich and diverse outputs. It also requires to deal with an interactive generation process, where the user can intervene at any time during the painting creation. In this paper, we also propose to solve CNP using a novel Transformer architecture which models the relationships between the user's and the generated strokes. Specifically, we introduce a novel Position-aware Attention Bias (PAAB) mechanism in order to encourage neighboring strokes to share similar semantics. Furthermore, we propose to model user-agent interactions using an Interaction-aware Masking (IAM) procedure which simulates interactions with the user at training time. To capture the complex conditional distribution of the strokes given the user's input, we adopt the diffusion framework (Ho et al., 2020) which has demonstrated remarkable results in conditional generation tasks and naturally supports the generation of diverse outputs.

Moreover, to facilitate this study and to promote future research in the proposed task, we present a CNP benchmark, composed of a novel dataset of painted objects and a new evaluation protocol which takes into account both the quantitative

and qualitative aspects of the CNP task. Creating such a dataset is challenging, since objects in natural images are frequently occluded, incomplete, or have a low resolution, also depending on their relative position within the scene. To overcome these limitations, we employ a state-of-the-art model to create (Rombach et al., 2022) and segment (Liu et al., 2023; Zhao et al., 2023) objects from generated images, and convert the resulting segmented objects into a stroke representation (Zou et al., 2021). Our method achieves the best performance among all metrics when compared with other state-of-the-art models that work with sequential data.

The contributions of our paper are summarized as follows:

- We propose a novel task, *Collaborative Neural Painting*, to encourage art generation process with active human engagement.
- We create and release a novel curated dataset to facilitate research on the task.
- We design a benchmark for the task with an evaluation protocol with both objective measures and subjective evaluation.
- We propose a diffusion model framework to address Collaborative Neural Painting, with a novel attention mechanism and masking scheme to model the user-agent interaction. Our proposed method outperforms recent baselines on our benchmark. We also showcased its effectiveness in real demonstration for collaborative painting with human users.

## 2. Related Works

In this section, we thoroughly discuss research works that are related to Neural Painting, Controllable Image Generation and Editing and Masked Data Modelling in generative tasks.

### 2.1. Neural Painting

Neural Painting (NP) refers to the task of decomposing a natural image into a set of parameterized strokes. The resulting strokes can be rendered on the canvas, obtaining an artistic version of the original image which mimics the result of an actual painting. The seminal works of (Haeberli, 1990; Litwinowicz, 1997; Hertzmann, 1998) target this task for the first time,

proposing heuristic-based methods to decompose a given image into a set of strokes. A parallel line of works focuses on developing rendering techniques that could faithfully represent, on digital media, the effect of different physical brushes (Wang et al., 2014; Sochorová and Jamriška, 2021; Strassmann, 1986; Curtis et al., 1997). With the progress of deep learning, these methods have been replaced by learning-based methods. In particular, the task is formulated as a reinforcement learning (RL) problem, where the agent is the painter, the action space is represented by the different strokes the agent can displace in the canvas, and the reward is given by the similarity between the painted canvas and the reference image (Huang et al., 2019; Singh et al., 2022a; Schaldenbrand and Oh, 2021). Differently, Zou *et. al.* (Zou et al., 2021) formalizes the task as an optimization problem, where the stroke parameters are directly optimized to approximate the reference image leveraging a differentiable renderer. To overcome the time burden of the optimization process, (Liu et al., 2021b) proposes to adopt a Transformer-based architecture for regressing the stroke parameters and design a synthetic data generation pipeline to train the model.

Despite achieving good qualitative results, previous methods do not consider the dimension of user interaction in the generation process. Intelli-Paint (Singh et al., 2022a) introduces this idea in a reinforcement learning framework, rewarding the agent for painting in a layered and localized way, more like a human would do. Another line of work is represented by Interactive Neural Painting (Peruzzo et al., 2023) in which the model *interacts* with the user by predicting the *next* strokes based on the user input. The final goal of the model is to paint a reference image provided by the user.

Our work shares the same motivation with (Singh et al., 2022a; Peruzzo et al., 2023), but we further push the boundaries of interactivity by explicitly targeting a collaborative scenario with the user without a reference image and without limiting the number of strokes predicted by the model. The goal of our model is to coherently complete an *object* based on the interactions with the user. Moreover, we address the more realistic problem of generating a painting without the reference image, which is a task a human can easily accomplish but none of the previous NP methods can do.

### 2.2. Controllable Image Generation and Editing

There exists a substantial corpus of literature for the purposes of controllable image generation and editing. Numerous GAN-based approaches have been developed, which involve the inversion and manipulation of images within the latent space. Other methods condition the generation on additional signals, which are intuitive for the user to control and modify, like segmentation maps, sketches, and text (Park et al., 2019; Zhu et al., 2020; Ghosh et al., 2019; Liu et al., 2021a; Ling et al., 2021). Notably, (Singh et al., 2022b) proposes to use parameterized strokes, similar to the ones adopted by NP methods, as an intuitive manner of sketching an image, and train a StyleGAN-like model to generate realistic images from incomplete drawings. Recently, diffusion-based approaches emerged as state-of-the-art for controllable image generation. Most of these methods

are powered by large pre-trained diffusion models (Rombach et al., 2022; Nichol et al., 2022; Balaji et al., 2022; Saharia et al., 2022), which are then finetuned or adapted for the specific task at hand. One line of works explores generating image variations given a set of representative images, either by finetuning the whole model (Ruiz et al., 2022) or by optimizing a text embedding (Gal et al., 2022; Mokady et al., 2022). Exploiting the recent advancement in language and vision understanding, other works propose to use text as an intuitive way of modifying the generated image (Brooks et al., 2023; Hertz et al., 2022; Parmar et al., 2023; Tumanyan et al., 2022). To improve control over the spatial layout of the final output, segmentation masks are introduced as conditioning in (Zeng et al., 2022; Avrahami et al., 2023), while sketches are used in (Voynov et al., 2022). ControlNet (Zhang and Agrawala, 2023) proposes an effective condition mechanism by training a hypernetwork on a dataset of paired examples consisting of images and condition signals, and conditioning the main pre-trained network using skip-connections.

Our work differs from these methods because we do not operate directly in the pixel space, but adopt a vectorial representation of the image. By representing the image as a set of parameterized strokes, we can achieve editing capabilities by design. Our method offers the flexibility to make adjustments at any level of detail, without relying on external models. Full control is given to the user, who can change the layout of the image by modifying the position of the strokes, resizing them, and rendering at arbitrary output resolution.

As summarized in Tab. 1, we position our work with respect to the existing literature in terms of (i) the model's capability to synthesize images from scratch, (ii) the ability to support interactive generation, (iii) the dependency on large-scale datasets or external models (*e.g.* pretrained segmentation networks), and (iv) the vectorial representation of images, which provides built-in editing capabilities.

### 2.3. Masked Data Modeling in Generative tasks

Masked Data Modeling refers to the task of reconstructing partially corrupted data, and recently gained popularity both for representation learning and generative purposes. Notably, BERT (Devlin et al., 2019) proposes a masked reconstruction strategy as a strong pretraining objective for NLP tasks. Different from autoregressive formulation, masking enables bidirectional context, increasing the model's expressiveness and performance. In the visual domain, MaskGIT (Chang et al., 2022) successfully applies a similar strategy to image generation, introducing an iterative sampling policy at inference. Recently, denoising diffusion models (Sohl-Dickstein et al., 2015; Ho et al., 2020; Song and Ermon, 2019) have emerged as a class of generative models showing state-of-the-art performances in many tasks. Masked Data Modeling can be integrated into the training of Diffusion Models, treating the unmasked regions as conditioning information. MAGVIT (Yu et al., 2023) extends MaskGIT to video domains, using a similar masking strategy. Wei *et al.* (Wei et al., 2023) propose to condition diffusion models on masked input to formulate them as masked autoencoders, which enables image inpainting at different levels of

**Table 1. Positioning of our work w.r.t. Neural Painting and Image Editing models.**

| Task | Neural Painting | Image Editing | Ours |
|---|---|---|---|
| (i) Synthesis from scratch | ✗ | ✓ | ✓ |
| (ii) Interactive generation | ✗ | ✓ | ✓ |
| (iii) No dependency on external models | ✓ | ✗ | ✓ |
| (iv) Vectorial representation | ✓ | ✗ | ✓ |

detail. Following this line of work, Tashiro *et al.* (Tashiro et al., 2021) propose a score-based diffusion model for the imputation of missing values in time series, simulating the missing data at training time through a carefully designed masking procedure. We build on these emerging trends and introduce a diffusion-based method for controllable neural painting exploiting masked data modeling as a methodology to simulate users' interactions with the model. Differently from previous works, our masking strategies are not *random* but *interaction-aware*, making it suitable for interactive tasks.

## 3. Preliminaries

We first introduce some knowledge that serves as background for our method, in terms of diffusion models and the conditional generation with classifier-free guidance.

**Diffusion Models.** Gaussian diffusion models rely on a forward noising process that gradually applies noise to real data $x_0$: $q(x_t|x_0) = \mathcal{N}(x_t; \sqrt{\bar{\alpha}_t}x_0, (1-\bar{\alpha}_t)\mathbf{I})$, where constants $\bar{\alpha}_t$ are hyperparameters. The model then learns the reverse process to invert forward corruption process: $p_\theta(x_{t-1}|x_t) = \mathcal{N}(\mu_\theta(x_t), \Sigma_\theta(x_t))$, using neural networks to approximate the intractable distribution $q(x_{t-1}|x_t)$. Using the $\epsilon$-parametrization (Ho et al., 2020), the model is trained with the mean-squared error between the predicted noise $\epsilon_\theta(x_t)$ and the ground truth sampled Gaussian noise $\epsilon_t$:

$$\mathcal{L}_{simple}(\theta) = \|\epsilon_\theta(x_t) - \epsilon_t\|_2^2 \tag{1}$$

At inference time, new data can be sampled by initializing $x_T \sim \mathcal{N}(0, \mathbf{I})$, and sampling $x_{t-1} \sim p_\theta(x_{t-1}|x_t)$. The reverse process can be expressed as:

$$p_\theta(x_{t-1}|x_t) = \frac{1}{\sqrt{\alpha_t}}\left(x_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}}\epsilon_\theta(x_t)\right) + \sigma_t\mathbf{z} \tag{2}$$

We adapt the DDPM formulation (Ho et al., 2020) and modify it to work in the masked setting of interactive neural painting, as we aim at denoising only a part of the stroke sequence.

**Classifier-free guidance.** Conditional diffusion models take additional information, such as a class label $c$, as input. In this scenario, the reverse process becomes $p_\theta(x_{t-1}|x_t, c)$. To guide the probability mass towards data where the implicit classifier $p_\theta(c|x_t)$ has a high likelihood, *classifier-free guidance* can be employed (Ho and Salimans, 2022) and can produce considerably improved samples over generic sampling methods (Nichol et al., 2022; Ramesh et al., 2022; Peebles and Xie, 2022). This requires training the model in both conditional and unconditional cases and merging the predicted scores. During training, the evaluation of the diffusion model with $c = \emptyset$ is accomplished by randomly dropping out $c$ and replacing it with a learned "null" embedding $\emptyset$. At inference time, given a guidance scale $s > 1$, the modified score becomes:

$$\hat{\epsilon}_\theta(x_t, c) = \epsilon_\theta(x_t, \emptyset) + s \cdot (\epsilon_\theta(x_t, c) - \epsilon_\theta(x_t, \emptyset)) \tag{3}$$

## 4. Task formulation

In this section, we define the Collaborative Neural Painting (CNP) task, where a painting is iteratively generated according to different conditioning signals provided by the user. We represent a painting using a sequence of parametrized strokes $\mathbf{s} = (\mathbf{s}_1, ..., \mathbf{s}_L) \in \mathbb{R}^{L \times 8}$, where $L$ is the number of strokes. Following previous literature (Zou et al., 2021), each stroke $\mathbf{s}_i \in \mathbb{R}^8$ is defined as a set of 8 parameters which describe: the position $(x, y)$, the size $(w, h)$, the rotation $(\theta)$ and the color $(r, g, b)$. The stroke sequence can be used to render the painting in the pixel space using a *parameter-free renderer* (Liu et al., 2021b). Given a primitive brushstroke, we apply a set of affine transformations dictated by the stroke parameter, obtaining the foreground and the alpha matte of the given stroke. The final painting is obtained by sequentially composing the individual strokes on the canvas. We refer to (Liu et al., 2021b) for additional details.

To mimic how human paints, we define and organize the stroke sequence into different **granularity levels** (Zou et al., 2021): from coarse strokes that define the object's shape in the first level towards smaller and finer details in the following levels. Specifically, we progressively divide the input image into overlapping blocks of size $m \times m$, with $m$ representing the granularity level. For each block in the current level, we initialize $N = 12$ strokes and optimize the parameters following (Zou et al., 2021) (See Fig. 3). In this work, we set the maximum number of levels $m$ to 4, resulting in 400 strokes per image.

To enable collaborative painting, we also introduce two complementary conditioning signals which can be used by the user to interact with the painting process: a class label $c$ and a sequence of strokes $\mathbf{s}^{ctx} \in \mathbb{R}^{L' \times 8}$ representing a partially incomplete painting. More formally, the Collaborative Neural Painting process consists in a $\mathcal{G}$, a.k.a. the paining agent, that produces the sequence of strokes $\mathbf{s}$ given the user-provided conditioning signals:

$$\mathbf{s} = \mathcal{G}(c, \mathbf{s}^{ctx}). \tag{4}$$

This formulation enables the *collaborative* generation process with great freedom. In its simplest form, the CNP framework can be used for the completion of paintings without any context, *i.e.* by setting $\mathbf{s}^{ctx}$ to an empty sequence $\emptyset$, and specifying only the desired object class $c$. Other functionalities are also possible by varying the conditioning signals. For instance: i) a coarse representation of the painting can be predicted with few strokes from the user, ii) a detailed painting can be produced from a coarse drawing by representing the coarse painting as the conditioning sequence $\mathbf{s}^{ctx}$ with or without class supervision, iii) an entire part of a generated subject can be altered by erasing it and letting the generator complete the painting with or without class supervision. The generation can be iterative according to user's demand, by simply updating the sequence
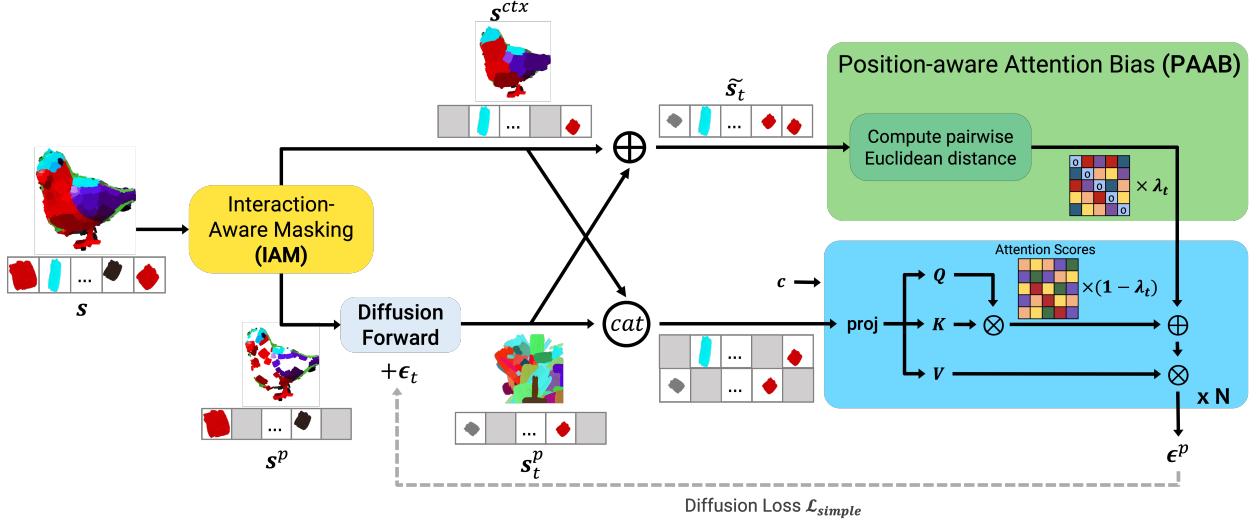
**Fig. 2. Given an input sequence s, the Interaction-Aware Masking block divides the sequence in two, a conditioning sequence $\mathbf{s}^{ctx}$ which acts as context to denoise the missing strokes $\mathbf{s}^p$ that, using a diffusion framework, are noised to $\mathbf{s}^p_t$. Our Position-aware Attention Bias modifies the attention scores of our Transformer based on the Euclidean distance between the conditioning and noised strokes.**
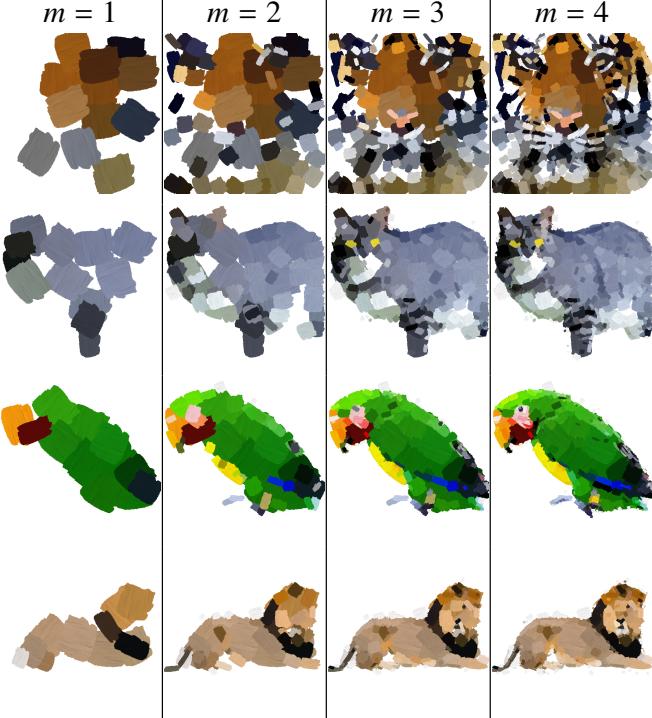


**Fig. 3. Examples of granularity levels. From left to right shows the granularity from the lowest ($m = 1$) to the highest ($m = 4$).**

$\mathbf{s}^{ctx}$ in Eq. 4, until the desired painting is obtained (see Fig. 1 for visual examples of these interactions).

## 5. Proposed Method

Our method is designed to mimic a collaborative scenario between the user and the painting agent $\mathcal{G}$. Given a sequence of strokes $\mathbf{s}$ describing the painting of an object of class $c$, we divide it into two subsequences: the conditioning sequence $\mathbf{s}^{ctx}$ and the target sequence $\mathbf{s}^p$. We propose an *Interaction-aware masking (IAM)* procedure to capture the different types of interactions between the user and $\mathcal{G}$ and use it to split the sequence $\mathbf{s}$ into context and target respectively.

Moreover, to effectively model the relationship between the target and conditioning stroke sequences, we introduce a *Masked Diffusion Transformer (MDT)* based on (Peebles and Xie, 2022). However, we find that the standard attention mechanism does not capture the relationship between strokes satisfactorily. To address this problem, we devise a *Position-aware Attention Bias (PAAB)* that encourages higher attention scores between spatially neighboring strokes.

We train our model using the $\epsilon$-parametrization of the diffusion framework (see Eq. 1). Random noise is added to the target sequence $\mathbf{s}^p$, and the model is trained to reconstruct the noise conditioned on the class information $c$ and the clean context sequence $\mathbf{s}^{ctx}$. At inference time, the user provides the context strokes (if any) and the class information. We then initialize the predicted strokes $\mathbf{s}^p_T$ with Gaussian noise and denoise them for $T$ steps conditioned on the user's inputs. In this stage, we leverage classifier-free guidance to increase the faithfulness of the predictions to the conditioning signals and extend it to a *Multiple Conditions Classifier-free guidance* to treat each of them separately.

In the next sections, we describe each part of our method in detail; Sec. 5.1 describes the Interaction-aware masking procedure, Sec. 5.2 describes our Transformer-based architecture, Sec. 5.3 describes our Position-aware Attention Bias, and finally, Sec. 5.4 describes the employed Multiple Conditions Classifier-free guidance procedure.

### 5.1. Interaction-aware Masking

The collaborative scenario described in Sec. 4 covers diverse ways the user and the painting agent can interact. They include modifying, adding, or deleting individual strokes, re-generating

a localized area of the painting, completing a coarse sketch provided by the user, and even generating an entire painting from scratch.

To mimic these types of interactions at training time, we introduce Interaction-aware masking (IAM) which, differently from the random masking strategies employed in previous works (Devlin et al., 2019), simulates at training time the interaction patterns that are expected at inference. In practice, we design different masks **m** to cover various use cases:

- **Granularity level**: simulates the task of generating fine details from coarse sketches of an object. We sample a random level of detail in the input sequence (see Sec. 4), and mask all the strokes belonging to successive levels with finer details.
- **Random**: simulates retouching operations on a given painting, involving the generation of strokes in different spatial locations and conveying different levels of details. We use random masking to mimic this use case.
- **Square**: simulates the generation of a local area of the painting. We remove multiple strokes in the sequence, based on the distance from a pivot stroke chosen at random, effectively creating a square mask in the rendered image. Note that spatially neighbor strokes are not necessarily consecutive in the stroke sequence.
- **Block**: simulates the undoing of the last $N$ consecutive strokes suggested by the model. We achieve it by masking adjacent strokes in the sequence which can span across granularity levels and be not spatially related to each other.
- **No context**: simulates the case where no context stroke $\mathbf{s}^{ctx}$ is provided by the user. In this case, we mask out the whole conditioning sequence and condition the model only on the desired object class $c$.

During training, for each sample, we randomly select one of these strategies and produce $\mathbf{s}^{ctx} = \mathbf{s} \odot \mathbf{m}$, where $\mathbf{s}$ is the stroke sequence sampled from the dataset.

## 5.2. Masked Diffusion Transformer

In this section, we describe the Masked Diffusion Transformer on which the collaborative painting generation process is based. We adopt the DDPM diffusion framework (Ho et al., 2020) and introduce a model $\mathcal{E}$ parametrized to predict the noise in the stroke sequence (see Fig. 2). The model is conditioned on the class $c$, the context strokes $\mathbf{s}^{ctx}$, and the diffusion timestep $t$:

$$\boldsymbol{\epsilon}^p = \mathcal{E}(\mathbf{s}_t^p | c, \mathbf{s}^{ctx}, t), \qquad (5)$$

with $\mathbf{s}_t^p$ being the noisy sequence of target stroke parameters, and $\boldsymbol{\epsilon}^p$ the predicted noise. Note that no noise is applied to the conditioning information $\mathbf{s}^{ctx}$ (Tashiro et al., 2021). We denote with $\tilde{\mathbf{s}}_t = \mathbf{s}^{ctx} \oplus \mathbf{s}_t^p$ the combination of conditioning and noisy strokes. At inference time, we initialize the predicted strokes $\mathbf{s}_T^p \sim \mathcal{N}(0, 1)$, and denoise them for $T$ steps, with $\tilde{\mathbf{s}}_0$ corresponding to the combination of the final predicted sequence and the context information.

Modeling an effective architecture for $\mathcal{E}$ is a challenging task, due to the sequential structure of strokes, and the necessity to model their relationship with the *variable-length* conditioning sequence provided by the user. We propose an approach based on a Transformer encoder architecture (Vaswani et al., 2017),

which can handle long input sequences and model pairwise stroke relationships thanks to self-attention. To accommodate our problem formulation, the network accepts two sequences of length $L$ as inputs. The first sequence is the noisy sequence of strokes $\mathbf{s}_t^p$, while the second is the conditioning strokes sequence $\mathbf{s}^{ctx}$, which provides the context for the denoising process. Each sequence is embedded using separate linear layers to form embedding sequences $\mathbf{e}^p, \mathbf{e}^{ctx} \in \mathbb{R}^{L \times F}$, with $F$ being the feature dimension of our Transformer. We consider strokes in corresponding positions to be mutually exclusive, *i.e.* a stroke is either a conditioning stroke, or a stroke to be predicted, and use a binary mask $\mathbf{m} \in \{0, 1\}^L$ which indicates, for each position in the sequence, if the stroke represents a conditioning signal (1) or is to be predicted (0). Finally, the input to the Transformer model is $\mathbf{e} = \mathbf{e}^p \odot (1 - \mathbf{m}) + \mathbf{e}^{ctx} \odot \mathbf{m}$.

Additionally, we condition the model on the class $c$, which is embedded into a learnable vector of size $F$. We combine this vector with the Fourier mapping of the logSNR corresponding to the diffusion time-step $t$ (Hoogeboom et al., 2023), and condition the model with the adaLN-Zero block (Peebles and Xie, 2022).

The output of the Transformer blocks is projected by a linear layer to the predicted noise $\boldsymbol{\epsilon}^p \in \mathbb{R}^{L \times 8}$. We exclude the context from the loss computation by masking the predicted noise, and train the model with the loss in Eq. (1).

## 5.3. Position-aware Attention Bias

The default attention mechanism makes each token attend to any other independently from its semantics. We argue that this behavior is suboptimal for the painting generation task, where each stroke in the sequence bears semantic-rich information. In particular, we consider local proximity as a robust prior information source. It is likely that spatially neighboring strokes will have an impact on one another and present consistent features, such as color or orientation. These cues can come from either the conditioning sequence, which provides real grounding for missing parameters, or the spatial neighboring strokes in the noisy sequence.

Recalling that traditional attention scores in self-attention are defined as $\sigma_{\text{attn}}(Q, K) = \text{softmax}\left(QK^T / \sqrt{d_k}\right)$, we address this limitation by introducing a novel Position-aware Attention Bias, which biases attention based on the distance between strokes. We compute the distance on the sequence obtained by combining the conditioning and the noisy strokes:

$$\sigma_{\text{PAAB}}(\tilde{\mathbf{s}}_t) = \text{softmax}\left(-((\tilde{\mathbf{s}}_t^i|_x - \tilde{\mathbf{s}}_t^j|_x)^2 + (\tilde{\mathbf{s}}_t^i|_y - \tilde{\mathbf{s}}_t^j|_y)^2)\right), \quad (6)$$

where $i, j \in [1, L]$ denote the position of the stroke in the sequence. We then combine the two scores with a weighted sum:

$$\sigma(Q, K, \tilde{\mathbf{s}}_t) = \lambda_t \cdot \sigma_{\text{PAAB}}(\tilde{\mathbf{s}}_t) + (1 - \lambda_t) \cdot \sigma_{\text{attn}}(Q, K) \quad (7)$$

where $\lambda_t$ represents the weighting factor. In early experiments, we found that using constant weighting value leads to suboptimal performances (see Tab.4). In fact, the larger the corrupting noise, the less confident we can be that two strokes will be neighbors in the final generated sequence. To account for this

phenomenon, we propose to control the strength of the attention bias as a function of the logSNR of the diffusion process:

$$\lambda_t = \frac{\log \text{SNR}_t - \log \text{SNR}_{\min}}{\log \text{SNR}_{\max} - \log \text{SNR}_{\min}} \cdot (\lambda_{\max} - \lambda_{\min}) + \lambda_{\min} \quad (8)$$

### 5.4. Multiple Conditions Classifier-free guidance

Our setting considers two different types of conditioning: the class information $c$ and the conditioning strokes $\mathbf{s}^{ctx}$. Liu *et al.* (Liu et al., 2022) show that a conditional diffusion model can produce improved results by combining score estimates from various conditioning signals. We apply the same strategy to our model with two separate conditioning input types. Following Classifier-free guidance formulation (Ho and Salimans, 2022), we train a single model dropping the conditioning signals during training. We use a special learnable token $\emptyset$ to learn the case where the class is not provided, while our *IAM* strategy already captures the scenario where no context stroke is given. We introduce two guidance scales, $s_1$ and $s_2$, which control, respectively, the correspondence to the conditioning strokes $\mathbf{s}^{ctx}$ and the class $c$. The modified score estimate becomes:

$$
\begin{aligned}
\hat{\epsilon}(\mathbf{s}_t^p, c, \mathbf{s}^{ctx}) = \; & \underbrace{\hat{\epsilon}(\mathbf{s}_t^p, \emptyset, \emptyset)}_{\text{unconditional}} \\
& + s_1 \cdot (\hat{\epsilon}(\mathbf{s}_t^p, \emptyset, \mathbf{s}^{ctx}) - \hat{\epsilon}(\mathbf{s}_t^p, \emptyset, \emptyset)) \\
& + s_2 \cdot ( \underbrace{\hat{\epsilon}(\mathbf{s}_t^p, c, \mathbf{s}^{ctx})}_{\text{full conditional}} - \underbrace{\hat{\epsilon}(\mathbf{s}_t^p, \emptyset, \mathbf{s}^{ctx})}_{\text{no class}} )
\end{aligned} \quad (9)
$$

Setting $s_1$ and $s_2$ to 1 would leave only the **full conditional** part. Note that alternative formulations to combine the scores are possible, *e.g.* by switching the position of $c$ and $\mathbf{s}^{ctx}$, but this formulation naturally accommodates our CNP task.

## 6. Evaluation Benchmark

In this section, we introduce a novel benchmark for the Collaborative Neural Painting task, which is designed to evaluate the performance of different methods in a collaborative scenario. We curate a novel dataset (Sec. 6.1) that covers a wide range of possible interactions between humans and the agent in the context of painting (Sec. 6.2). Additionally, we devise a set of metrics to quantitatively evaluate the performance of these models with this dataset (Sec. 6.2).

### 6.1. Dataset

Given the novel nature of the proposed task, we could not find any ready-to-use public dataset for it. At the same time, existing large-scale datasets of natural images are not suitable as they do not specifically feature objects, with occlusion possibly compromising the stroke representation. Therefore, we design a data engine to produce a large curated dataset to enable the study on the CNP task. The data generation pipeline is divided into three steps and we showcase it in Fig. 6:

1. We feed the prompts *"a photo of a single full-size, full-body [obj], whole figure"*, *"a photo of a full-body [obj]"*, *"a high-resolution photo of a full-body [obj]"*, *"a DSLR photo of a full-body [obj]"* to Stable Diffusion (Rombach et al., 2022), with [obj] indicating the desired class. We use Stable Diffusion v2.1 and generate images at resolution $512 \times 512$ using the DPMSolver (Lu et al., 2022) scheduler.

2. We remove the background from the generated image, as we are interested in the object itself. To do so, we first use GroundingDINO (Liu et al., 2023) conditioned on the known class [obj] to obtain a bounding box of the object and Fast-SAM (Zhao et al., 2023) to obtain the corresponding segmentation mask. Then, with the bounding box and segmentation mask, we isolate the object from the background and we center crop the image to have the object in the middle of it. Since we know the desired class beforehand, we exploit the rich semantics from text to segment the animals. This also allows us to build an automatic pipeline from [obj] to the final stroke representation without manual intervention.

3. We represent the foreground image as a sequence of strokes, using Stylized Neural Painting (Zou et al., 2021) for its flexibility. We set the parameters of SNP to have four layers of details and limit the number of strokes describing each image to 400.

Note that our pipeline follows a modular design, where each module can be flexibly updated. Our current implementation is composed of off-the-shelf components that represents the state of the art in order to achieve the best quality dataset. Moreover, our data-generation pipeline is fully automatic and class agnostic. Thus, it can readily include more object categories or increase the dataset size and diversity. In this work, we choose to generate data of 10 animal classes: *cat, dog, eagle, elephant, lion, parrot, rabbit, squirrel, tiger* and *wolf*. They are easily recognizable with great inter and intra-class diversity, thus allowing for objective and subjective evaluation with less ambiguity to kick off the research. We generate at least 10K images for each class, with 500 left out for testing. There are a total of 101.052 decomposed sequences for the whole dataset. We show some examples in Fig. 4 while we show examples of generated and segmented images in Fig. 6.

### 6.2. Evaluation Protocol

Existing benchmarks for image synthesis are inadequate to capture the interactive nature of the CNP task. In Sec. 5.1 we discussed several masking strategies, namely *block, level, random, no context*, and *square*, designed to capture multiple interactive scenarios. We use those masking strategies as tasks to mimic, in an automatic way, the interaction between the user and the painting agent, and use them for evaluation.

**Metrics.** We introduce a set of metrics to quantitatively evaluate the proposed methods. First, we render the predicted and context strokes $\tilde{\mathbf{s}}_0$, and compare it with the rendered test set. Following the image inpainting literature (Zheng et al., 2022; Liu et al., 2019), we compute the Fréchet Inception Distance (**FID**) (Heusel et al., 2017), and a sample-wise $\mathcal{L}_2$ distance. Second, we design a metric specific for the stroke-based formulation and evaluate the similarity between predicted parameters

**Table 2. Quantitative results of our models compared to the baselines. Stroke $\mathcal{L}_1$ is reported in $\times 10^{-3}$, Image $\mathcal{L}_2$ is reported in $\times 10^{-4}$. All the metrics the lower the better.**

| Method | Block | | | Level | | | Random | | | Square | | | No ctx |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | FID | Stroke $\mathcal{L}_1$ | Image $\mathcal{L}_2$ | FID | Stroke $\mathcal{L}_1$ | Image $\mathcal{L}_2$ | FID | Stroke $\mathcal{L}_1$ | Image $\mathcal{L}_2$ | FID | Stroke $\mathcal{L}_1$ | Image $\mathcal{L}_2$ | FID |
| Continuous Transformer | 35.87 | 144 | 188 | 81.12 | 145 | 269 | 101.54 | 127 | 665 | 25.54 | **135** | 288 | 443.86 |
| BERT | 142.89 | 198 | 328 | 250.54 | 196 | 406 | 247.07 | 196 | 117 | 33.89 | 197 | 195 | 320.43 |
| MaskGIT | 149.54 | 205 | 265 | 261.27 | 206 | 398 | 250.75 | 200 | 111 | 35.16 | 211 | 193 | 336.18 |
| Ours | **6.20** | **127** | **94** | **7.29** | **125** | **100** | **12.69** | 126 | 103 | **5.53** | 142 | **154** | **30.12** |



**Fig. 4. Examples of rendered sequences from the dataset.**

and ground truth strokes computing a $\mathcal{L}_1$ distance. Although ordering is crucial for determining the rendering priority of overlapping strokes and subsequent levels, non-overlapping strokes in the same granularity level can be swapped without affecting the result. For this reason, before computing the distance, we perform a Hungarian Matching (Kuhn, 1955).

We compute the metrics for each task separately to gauge their respective difficulty. In practice, we sample ~5K sequences from the test set and create the conditioning sequence $\mathbf{s}^{ctx}$ by applying the masking corresponding to the given task. Note that for the *no context* case, we cannot compute sample-wise metrics, thus we rely only on the FID.

# 7. Experiments

To the best of our knowledge, no existing model is able to operate interactively on painting strokes for image generation. Hence, for comparative analysis, we used state-of-the-art methods for sequence modeling and adapted them to fit our task. We first provide the implementation details of our method in Sec 7.1. We then present comparison against the baselines and discuss the results in Sec. 7.2. Besides automatic metrics, we have conducted a user study to assess human preference among different methods and report the results in Sec. 7.3. Moreover, we evaluate the effect of our proposed components and training

strategies with an ablation study in Sec. 7.4. Lastly, we report the real-time analysis of our method in Sec. 7.5.

## 7.1. Implementation Details

**Architecture.** Our model is based on the Transformer (Vaswani et al., 2017) architecture. Given an input sequence of strokes of dimension $L \times 8$, we split it into context and target using our IAM masking strategy (see Fig. 2 of main paper). In practice, the sequence is multiplied by two binary masks, with zero to mask out strokes. We then concatenate the context and target sequences on the feature dimension, resulting in a sequence of size $L \times 16$. Next, we use a linear layer to project the input sequence of strokes to the feature size of the Transformer model $L \times F$. Since the two sequences are mutually exclusive with value 0 in the *empty* positions, the linear projection does not act on these values, effectively acting as two separate projections on the two sequences. Lastly, we add sinusoidal positional embeddings (Vaswani et al., 2017) to the sequence and feed it to the Transformer model. GELU non-linearities (Hendrycks and Gimpel, 2016) (approximated with tanh) are used in the core Transformer.

To inject the diffusion timestep information $t$ and control the diffusion process with the class information $c$, we employ the adaLN-Zero block (Peebles and Xie, 2022), due to its demonstrated performance. To embed the input time steps we use a 256-dimensional frequency embedding (Dhariwal and Nichol, 2021), followed by a two-layer MLP with SiLU activations with a output dimension equal to the Transformer's hidden size $F$. Next, we sum the time-step and the class embedding and feed it to the adaLN-Zero layer. In addition to regressing the values of the scale $\gamma$ and the shift $\beta$ as in AdaIN (Huang and Belongie, 2017), adaLN-Zero also regresses dimension-wise scaling parameters $\delta$, which are applied before any residual connection within the blocks. We regress the scaling parameters using a linear projection, which is initialized to output a zero-vector for all the $\delta$ values. For each operation $op$ (attention or feed-forward) the equation becomes: out $= x + \delta \cdot op(\gamma \cdot x + \beta)$.

**Masking.** The *random* masking technique is implemented by randomly selecting a masking ratio in the range $[0.1, 0.9]$ and dropping strokes accordingly in the sequence. For the *level* masking, we randomly select a granularity level between 1 and 3, and mask all the strokes that belong to subsequent levels, effectively training the model to add details to an existing rough representation of the given object. The *block* masking technique involves randomly selecting a starting point in the sequence and masking a block of contiguous strokes with a length in the range $[10, 0.75 \cdot L]$, with $L$ the sequence length. Lastly, in *square* case we select a stroke in the sequence at random as the center and build a square of side length 0.5, masking all
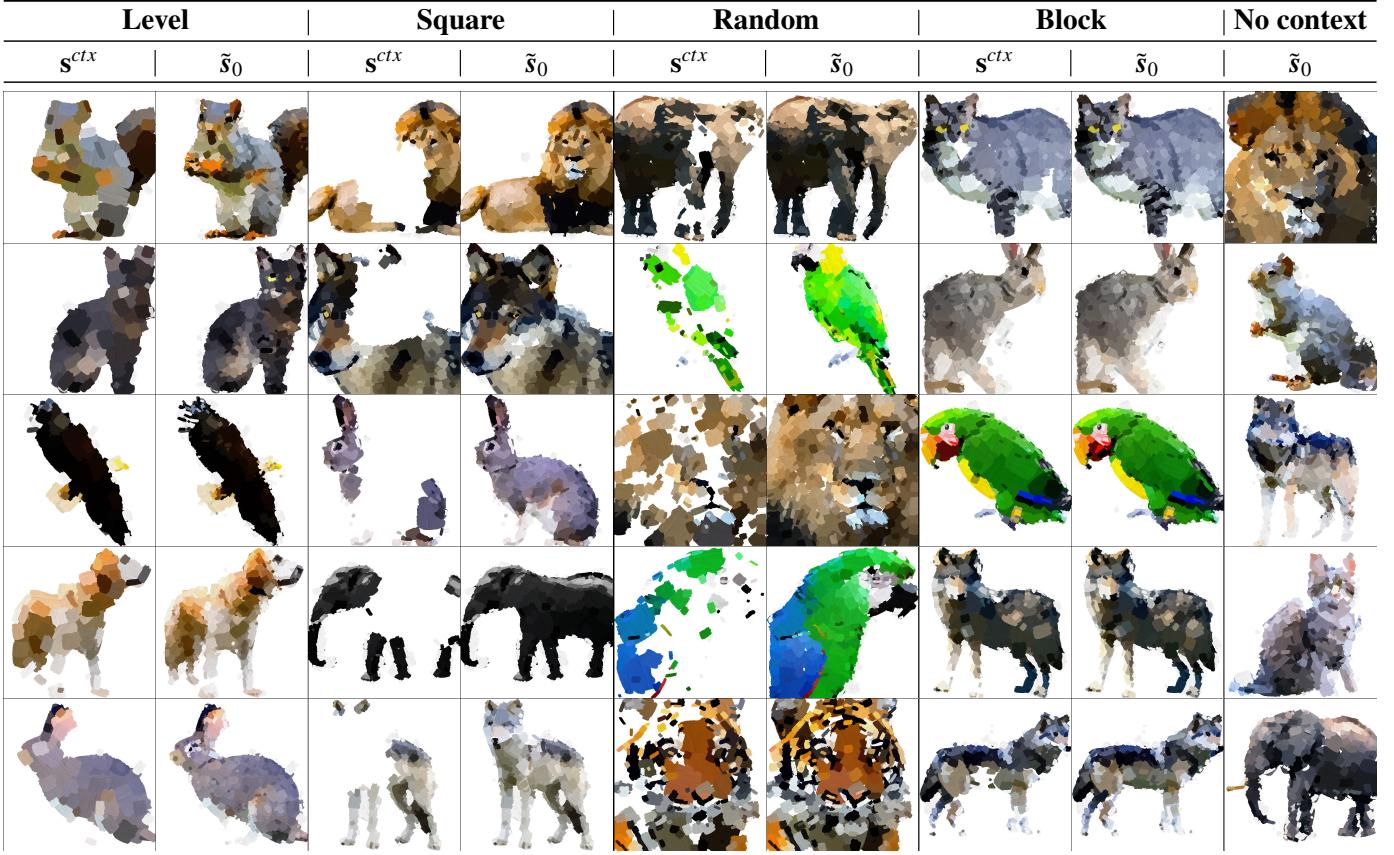
| Level | | Square | | Random | | Block | | No context |
|---|---|---|---|---|---|---|---|---|
| $\mathbf{s}^{ctx}$ | $\tilde{s}_0$ | $\mathbf{s}^{ctx}$ | $\tilde{s}_0$ | $\mathbf{s}^{ctx}$ | $\tilde{s}_0$ | $\mathbf{s}^{ctx}$ | $\tilde{s}_0$ | $\tilde{s}_0$ |



**Fig. 5. Qualitative completion results of our method. Images are organized in blocks of 2. On the left the conditioning sequence $\mathbf{s}^{ctx}$, on the right the conditioning with the completion added,** *i.e.* $\tilde{s}_0$.

the strokes spatially contiguous to the center. During training, we select one of the above strategies with equal probability and apply it to split the sequence into context and target strokes.

**Diffusion Model.** We followed the formulation of "simple diffusion" (Hoogeboom et al., 2023), defining the noising schedule in terms of SNR $= \alpha_t^2/\sigma_t^2$ and keeping the default values min log SNR $= -15$ and max log SNR $= 15$. We do not employ any weighting on the losses and train the model with the $\mathcal{L}_{simple}$ formulation (Ho et al., 2020).

**Training hyper-parameters.** We train all models with Adam (Kingma and Ba, 2014) and a cosine annealing schedule for the learning rate, where the maximum value is $1 \times 10^{-4}$. We train with a global batch size of 256 divided on 4 GPUs. Following the common practice in the generative modeling literature, we maintain an exponential moving average (EMA) of the weights over training with a decay of 0.9999. All reported results use the EMA model.

**Sampling.** For the diffusion reverse process we use the standard DDPM (Ho et al., 2020) sampler with 1000 steps.

**MC-CFG.** We choose the hyper-parameters of our Classifier-free guidance, $s_1$ and $s_2$ to be 1.5 in order to weigh the importance of the conditioning sequence and the class information.

**Inference-time stroke position.** At inference time, *e.g.* the demo, we start without strokes and we do not use the *IAM* module. Since Stylized Neural Painting fixes the maximum number of strokes in each cell, when the user draws strokes on the canvas at inference time, the size and position of the strokes allow

finding the corresponding cell. Given a level $l$ and the grid described above, a stroke in that level cannot have a size greater than $1/l$, therefore we can pick the correct level by choosing the first level whose maximum stroke size is the closest to the input. After the level, the correct cell can be identified by the position of the strokes. Then, the stroke gets added to the part of the sequence corresponding to the given cell in the first available slot, and we update the occupant mask $\mathbf{m}$ accordingly.

### 7.2. Comparison

**Baselines.** Inspired by the success of Transformer-based models in NLP tasks, where the input belongs to a fixed-sized vocabulary, we compare our method with a baseline operating with discrete input. In practice, we discretize the stroke parameters independently into a codebook of size 256. We train the model in a BERT-like fashion (Devlin et al., 2019), masking the strokes and predicting the original tokens with cross-entropy loss. At inference time, we operate this baseline in two ways: (i) predicting the missing strokes in one step (referred to as *BERT*), or (ii) following (Chang et al., 2022), iteratively sampling based on the network confidence on the generated tokens, leading to a multi-step sampling process similar to diffusion-based models (referred to as *Mask-GIT*). Additionally, we compare with a baseline working with continuous stroke parameters, named the *Continuous Transformer* baseline. In this case, the model is trained to regress the masked stroke parameters with an MSE loss. We employ a Transformer architecture with

similar configurations, FLOPS, and parameter counts, and apply our IAM to all the models.

As shown in Tab. 2, the *Continuous Transformer* baseline fails to capture the distribution of target strokes and to use the contextual information provided by the conditioning sequence. Due to the MSE loss used during training, it converges to a mean representation for each class, with failure cases more pronounced in the unconditional case (see Fig. 7). On the other hand, the baselines working with discrete strokes fail completely at modeling the relationship between the conditioning strokes and the masked ones. We impute this to the independent conversion of each stroke parameter to a different token, leading to an input sequence of length $L \times 8$, and making the computation of relationships between strokes challenging. Our method consistently outperforms other approaches and achieves strong performances across all the different tasks.

**Qualitative Results.** Qualitative results are presented in Fig. 5 and 8, where the images are displayed in the format of a conditioning sequence $s^{ctx}$ and predicted completion $\tilde{s}_0$. We notice how our method predicts completion consistent with the context information in various tasks. For example, in the first column (*Level*), given a rough sketch of the desired image as context, the model adds fine-grained details in a consistent manner. At the same time, when a large portion of a detailed painting is missing, the model generates both coarse and detailed strokes that harmonize with the context and complete it in a plausible manner (e.g. second column, *Square*). Finally, our method can generate realistic paintings without any conditioning but the class information (last column, *No context*), showcasing another useful use case.

*Robustness to random masking.* We investigate the robustness of our method to varying levels of random masking and complement the quantitative results of Tab. 3 of the main paper with the qualitative results in Fig. 9. This experiment showcases the effectiveness of the conditioning sequence as a strong prior for generation, even in high masking regimes. The results suggest that our model is capable of producing accurate predictions with limited input strokes.

*Diverse suggestions from the same context.* In Fig. 10 we exhibit the strength of diffusion models to produce diverse results given the same conditioning sequence. This feature is crucial for the Collaborative Neural Painting task, providing the user with diverse but coherent completion given the same context. In practice, we input the same context strokes $s^{ctx}$ with the class $c$, and sample different initial random noise $s_T \sim \mathcal{N}(0, 1)$. We can observe how all the competitions suggested by the model are coherent with the context while providing some variety and diversity among which the users can make their final decision.

*Automatic inference without providing class.* We probe the role of the class conditioning $c$ by dropping it and providing only the context strokes to the model. We show the results in Fig. 11. We can notice that the competition suggested by the model reflects the expected class, suggesting that during training the model learns to rely on the context information and implicitly learns to associate it with a specific class. This phenomenon is further encouraged by the CFG training procedure employed, in which the class is not always used to train the model.

We also refer the reader to the Supp. Mat. for a demonstration video with human interaction.

## 7.3. User Study

We complement our quantitative analysis with a user study to compare the users' preferences among the different methods. We ask the users to rank the completion from different methods from worst to best, given the rendered image of the conditioning sequence $s^{ctx}$. We compare the predictions of our method with the ones obtained with *Continuous Transformer*, *MaskGIT*, and the *ground truth*. The study has been conducted on 20 different users, collecting a total of 894 votes. We report the preferences in Tab. 3, where we exhibit competitive results with the ground truth as the best completion, and outperform the other methods as the second-best choice.

**Table 3. User study results, we compare our model against baselines and ground truth completion.**

|  | Ground truth | Ours | Continuous | MaskGIT |
|---|---|---|---|---|
| 1st choice | 81% | 17% | 1% | 1% |
| 2nd choice | 18% | 80% | 1% | 1% |

## 7.4. Ablation Study

In this section, we evaluate the effectiveness of the components introduced in Sec. 5, and present the results in Tab. 4.

**IAM.** We start our analysis by training a model without any of our components, with simple random masking similar to BERT (first row). We then introduce IAM and, as expected, its role is crucial to perform well at inference time on the synthetic tasks, especially in the *no context* case (second row).

**PAAB.** We then incorporate the Position-aware Attention bias (PAAB), using a constant weighting of $\lambda = 0.5$ across the diffusion time-steps $t$. We observe an performance degradation in all the tasks, suggesting this fixed bias is not optimal to guide the generative process.

**Ada-$\lambda$.** To counter this effect, we introduce adaptive weighting $\lambda$ which modifies the strength of the bias as a function of log SNR of the diffusion model. The results in the fourth row prove the effectiveness of this choice, reducing the FID in all the tasks, and particularly in the *no ctx* case.

**MC-CFG.** In the last row, we explore the role of MC-CFG at inference time, which leads to further improvement in the performances.

**Robustness.** Moreover, we study the robustness of our model to different masking ratios. We test the model on *random* task, varying the percentage of masked strokes (see Tab. 5). Our model exhibits good performance even in a high masking regime, *i.e.* 80%, showing that even a few strokes are sufficient to provide context to the network.

**Scaling.** Lastly, we investigate the scaling laws governing our model and design three configurations: MDT-S(mall), MDT-B(ase), and MDT-L(arge), as described in Tab. 6. The results reported in the table correspond to the ablation with all components included except MC-CFG. We observed a consistent performance improvement by scaling up the model, suggesting that performance could be further improved by increasing the capacity.

**Table 4. Ablation study of our proposed components. IAM, PAAB, and MC-CFG stand for Interaction-aware Masking, Position-aware Attention Bias, and Multiple Conditions Classifier-free guidance, respectively.**

| IAM | PAAB | Ada-$\lambda$ | MC-CFG | FID ↓ | | | | |
|-----|------|---------------|--------|-------|-------|--------|--------|--------|
| | | | | Block | Level | Random | Square | No ctx |
| ✗ | ✗ | ✗ | ✗ | 8.78 | 11.66 | 18.41 | 7.51 | 233.21 |
| ✓ | ✗ | ✗ | ✗ | 9.23 | 11.65 | 21.92 | 7.17 | 36.96 |
| ✓ | ✓ | ✗ | ✗ | 9.51 | 11.85 | 22.61 | 7.34 | 37.94 |
| ✓ | ✓ | ✓ | ✗ | 8.43 | 10.34 | 19.25 | 6.63 | 30.46 |
| ✓ | ✓ | ✓ | ✓ | **6.20** | **7.29** | **12.69** | **5.53** | **30.12** |

**Table 5. Performance of our model w.r.t. masking percentage in the *random* setting.**

| Masking % | 20% | 40% | 60% | 80% | 100% |
|-----------|-----|-----|-----|-----|------|
| FID ↓ | 6.57 | 10.42 | 13.67 | 19.05 | 30.12 |

**Table 6. Details of MDT models. Model configurations for the Small (S), Base (B), and Large (L) variants and corresponding performance on the proposed task.**

| Model | Layers N | F | Heads | GFLOPS | FID ↓ | | | | |
|-------|----------|---|-------|--------|-------|-------|--------|--------|---------------|
| | | | | | Block | Level | Random | Square | Unconditional |
| MDT-S | 6 | 576 | 6 | 9.09 | 8.67 | 11.24 | 22.69 | 7.34 | 55.25 |
| MDT-B | 8 | 768 | 12 | 21.27 | 6.81 | 8.39 | 14.96 | 6.02 | 36.29 |
| MDT-L | 12 | 768 | 12 | 31.89 | **6.20** | **7.29** | **12.69** | **5.53** | **30.12** |

*7.5. Real-time analysis*

**Table 7. Model performance w.r.t. number of sampling steps.**

| Sampling steps | FID ↓ | | | | |
|----------------|-------|-------|--------|--------|---------------|
| | Block | Level | Random | Square | Unconditional |
| 35 | 6.42 | 7.83 | 12.81 | 5.40 | 36.74 |
| 50 | 6.10 | 7.59 | 11.97 | 5.34 | 31.64 |
| 70 | 5.90 | 7.12 | 11.71 | 5.28 | 28.30 |
| 125 | 5.95 | 7.01 | 12.06 | 5.24 | 27.95 |
| 250 | 5.98 | 7.19 | 12.12 | 5.42 | 28.05 |
| 500 | 6.08 | 7.09 | 12.37 | 5.47 | 29.73 |
| 1000 | 6.20 | 7.29 | 12.69 | 5.53 | 30.12 |

We test our model on a single NVIDIA A100. Performing 1000 denoising steps requires ∼ 8 seconds, meaning ∼ 8 ms per step. In Table 7, we present the effect of using a reduced number of steps for sampling, and we find that the model performs consistently well after ∼ 70 steps. A number of steps under 70 leads to degraded performance. Thus, we opt to use 70 steps for demo purposes, resulting to a sampling time of ∼ 560 ms, a reasonable time for real-time interactions.

## 8. Conclusions

We introduced Collaborative Neural Painting, a novel task to facilitate collaborative art generation. We exploit a parametrized vector formulation for strokes to achieve editability and composability in painting generation. We proposed a novel Transformer-based architecture to solve the Collaborative Neural Painting task and designed a novel attention mechanism and masking scheme to tackle the challenges brought by the stroke formulation and the arbitrary user interaction, exhibiting state-of-the-art performance both quantitatively and qualitatively. Our proposed method is class-agnostic thus being potentially applicable to datasets of higher semantic cardinality.

Future study involves extending the dataset for further investigation and incorporating language models to enable open-set semantics. Moreover, we plan to investigate more efficient sampling and diffusion strategies in order to enable a more reactive user experience with the model. Finally, it is also interesting to study how to better blend the generated objects with the background.

## References

Avrahami, O., Hayes, T., Gafni, O., Gupta, S., Taigman, Y., Parikh, D., Lischinski, D., Fried, O., Yin, X., 2023. Spatext: Spatio-textual representation for controllable image generation, in: Proceedings of IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR).

Balaji, Y., Nah, S., Huang, X., Vahdat, A., Song, J., Kreis, K., Aittala, M., Aila, T., Laine, S., Catanzaro, B., et al., 2022. ediff-i: Text-to-image diffusion models with an ensemble of expert denoisers. arXiv preprint arXiv:2211.01324 .

Beh-Pajooh, A., Abdollahi, A., Hosseinian, S., 2018. The effectiveness of painting therapy program for the treatment of externalizing behaviors in children with intellectual disability,‖ vulnerable child. Vulnerable Children and Youth Studies .

Brooks, T., Holynski, A., Efros, A.A., 2023. Instructpix2pix: Learning to follow image editing instructions, in: Proceedings of IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR).

Chang, H., Zhang, H., Jiang, L., Liu, C., Freeman, W.T., 2022. Maskgit: Masked generative image transformer, in: Proceedings of IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR).

Curtis, C.J., Anderson, S.E., Seims, J.E., Fleischer, K.W., Salesin, D.H., 1997. Computer-generated watercolor, in: Proceedings of the Conference on Computer Graphics and Interactive Techniques (SIGGRAPH).

Devlin, J., Chang, M.W., Lee, K., Toutanova, K., 2019. BERT: Pre-training of deep bidirectional transformers for language understanding, in: Proceedings of the Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL-HLT).

Dhariwal, P., Nichol, A., 2021. Diffusion models beat gans on image synthesis.

Gal, R., Alaluf, Y., Atzmon, Y., Patashnik, O., Bermano, A.H., Chechik, G., Cohen-Or, D., 2022. An image is worth one word: Personalizing text-to-image generation using textual inversion. arXiv preprint arXiv:2208.01618 .

Gerry, L.J., 2017. Paint with me: stimulating creativity and empathy while painting with a painter in virtual reality. IEEE transactions on visualization and computer graphics .

Ghosh, A., Zhang, R., Dokania, P.K., Wang, O., Efros, A.A., Torr, P.H.S., Shechtman, E., 2019. Interactive sketch & fill: Multiclass sketch-to-image translation, in: Proceedings of IEEE/CVF International Conference on Computer Vision (ICCV).

Haeberli, P., 1990. Paint by numbers: Abstract image representations, in: Proceedings of the Conference on Computer Graphics and Interactive Techniques (SIGGRAPH).

Hendrycks, D., Gimpel, K., 2016. Gaussian error linear units (gelus). arXiv preprint arXiv:1606.08415 .

Hertz, A., Mokady, R., Tenenbaum, J., Aberman, K., Pritch, Y., Cohen-Or, D., 2022. Prompt-to-prompt image editing with cross attention control. arXiv preprint arXiv:2208.01626 .

Hertzmann, A., 1998. Painterly rendering with curved brush strokes of multiple sizes, in: Proceedings of the Conference on Computer Graphics and Interactive Techniques (SIGGRAPH).

Heusel, M., Ramsauer, H., Unterthiner, T., Nessler, B., Hochreiter, S., 2017. Gans trained by a two time-scale update rule converge to a local nash equilibrium.

Ho, J., Jain, A., Abbeel, P., 2020. Denoising diffusion probabilistic models, in: Advances in Neural Information Processing Systems (NeurIPS).

Ho, J., Salimans, T., 2022. Classifier-free diffusion guidance. arXiv preprint arXiv:2207.12598 .

Hoogeboom, E., Heek, J., Salimans, T., 2023. simple diffusion: End-to-end diffusion for high resolution images. arXiv preprint arXiv:2301.11093 .

Huang, X., Belongie, S., 2017. Arbitrary style transfer in real-time with adap-

tive instance normalization, in: Proceedings of IEEE/CVF International Conference on Computer Vision (ICCV).

Huang, Z., Zhou, S., Heng, W., 2019. Learning to paint with model-based deep reinforcement learning, in: Proceedings of IEEE/CVF International Conference on Computer Vision (ICCV).

Kingma, D.P., Ba, J., 2014. Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980 .

Kotovenko, D., Wright, M., Heimbrecht, A., Ommer, B., 2021. Rethinking style transfer: From pixels to parameterized brushstrokes, in: Proceedings of IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR).

Kuhn, H.W., 1955. The hungarian method for the assignment problem. Naval research logistics quarterly .

Ling, H., Kreis, K., Li, D., Kim, S.W., Torralba, A., Fidler, S., 2021. Editgan: High-precision semantic image editing, in: Advances in Neural Information Processing Systems (NeurIPS).

Litwinowicz, P., 1997. Processing images and video for an impressionist effect, in: Proceedings of the Conference on Computer Graphics and Interactive Techniques (SIGGRAPH).

Liu, H., Jiang, B., Xiao, Y., Yang, C., 2019. Coherent semantic attention for image inpainting, in: Proceedings of IEEE/CVF International Conference on Computer Vision (ICCV).

Liu, H., Wan, Z., Huang, W., Song, Y., Han, X., Liao, J., Jiang, B., Liu, W., 2021a. Deflocnet: Deep image editing via flexible low-level controls, in: Proceedings of IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR).

Liu, N., Li, S., Du, Y., Torralba, A., Tenenbaum, J.B., 2022. Compositional visual generation with composable diffusion models, in: Proceedings of IEEE/CVF European Conference on Computer Vision (ECCV).

Liu, S., Lin, T., He, D., Li, F., Deng, R., Li, X., Ding, E., Wang, H., 2021b. Paint transformer: Feed forward neural painting with stroke prediction, in: Proceedings of IEEE/CVF International Conference on Computer Vision (ICCV).

Liu, S., Zeng, Z., Ren, T., Li, F., Zhang, H., Yang, J., Li, C., Yang, J., Su, H., Zhu, J., et al., 2023. Grounding dino: Marrying dino with grounded pre-training for open-set object detection. arXiv preprint arXiv:2303.05499 .

Lu, C., Zhou, Y., Bao, F., Chen, J., Li, C., Zhu, J., 2022. Dpm-solver: A fast ode solver for diffusion probabilistic model sampling in around 10 steps. Advances in Neural Information Processing Systems 35, 5775–5787.

Lüddecke, T., Ecker, A., 2022. Image segmentation using text and image prompts, in: Proceedings of IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR).

Meng, C., He, Y., Song, Y., Song, J., Wu, J., Zhu, J., Ermon, S., 2022. Sdedit: Guided image synthesis and editing with stochastic differential equations, in: Proceedings of International Conference on Learning Representations (ICLR).

Mokady, R., Hertz, A., Aberman, K., Pritch, Y., Cohen-Or, D., 2022. Null-text inversion for editing real images using guided diffusion models. arXiv preprint arXiv:2211.09794 .

Nakano, R., 2019. Neural painters: A learned differentiable constraint for generating brushstroke paintings. arXiv preprint arXiv:1904.08410 .

Nichol, A.Q., Dhariwal, P., Ramesh, A., Shyam, P., Mishkin, P., McGrew, B., Sutskever, I., Chen, M., 2022. GLIDE: towards photorealistic image generation and editing with text-guided diffusion models, in: Proceedings of International Conference on Machine Learning (ICML).

Park, T., Liu, M.Y., Wang, T.C., Zhu, J.Y., 2019. GauGAN: Semantic image synthesis with spatially adaptive normalization, in: Proceedings of the Conference on Computer Graphics and Interactive Techniques (SIGGRAPH).

Parmar, G., Singh, K.K., Zhang, R., Li, Y., Lu, J., Zhu, J.Y., 2023. Zero-shot image-to-image translation. arXiv preprint arXiv:2302.03027 .

Peebles, W., Xie, S., 2022. Scalable diffusion models with transformers. arXiv preprint arXiv:2212.09748 .

Pelowski, M., Leder, H., Tinio, P.P., 2017. Creativity in the visual arts. The Cambridge handbook of creativity across domains .

Peruzzo, E., Menapace, W., Goel, V., Arrigoni, F., Tang, H., Xu, X., Chopikyan, A., Orlov, N., Hu, Y., Shi, H., Sebe, N., Ricci, E., 2023. Interactive neural painting. Computer Vision and Image Understanding .

Ramesh, A., Dhariwal, P., Nichol, A., Chu, C., Chen, M., 2022. Hierarchical text-conditional image generation with clip latents. arXiv preprint arXiv:2204.06125 .

Rombach, R., Blattmann, A., Lorenz, D., Esser, P., Ommer, B., 2022. High-resolution image synthesis with latent diffusion models, in: Proceedings of IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR).

Ruiz, N., Li, Y., Jampani, V., Pritch, Y., Rubinstein, M., Aberman, K., 2022. Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. arXiv preprint arxiv:2208.12242 .

Saharia, C., Chan, W., Saxena, S., Li, L., Whang, J., Denton, E.L., Ghasemipour, K., Gontijo Lopes, R., Karagol Ayan, B., Salimans, T., et al., 2022. Photorealistic text-to-image diffusion models with deep language understanding, in: Advances in Neural Information Processing Systems (NeurIPS).

Schaldenbrand, P., Oh, J., 2021. Content masked loss: Human-like brush stroke planning in a reinforcement learning painting agent, in: Proceedings of the AAAI conference on artificial intelligence.

Singh, J., Smith, C., Echevarria, J., Zheng, L., 2022a. Intelli-paint: Towards developing more human-intelligible painting agents, in: Proceedings of IEEE/CVF European Conference on Computer Vision (ECCV).

Singh, J., Zheng, L., Smith, C., Echevarria, J., 2022b. Paint2pix: interactive painting based progressive image synthesis and editing, in: Proceedings of IEEE/CVF European Conference on Computer Vision (ECCV).

Sochorová, v., Jamriška, O., 2021. Practical pigment mixing for digital painting. ACM Transactions on Graphics .

Sohl-Dickstein, J., Weiss, E.A., Maheswaranathan, N., Ganguli, S., 2015. Deep unsupervised learning using nonequilibrium thermodynamics, in: Proceedings of International Conference on Machine Learning (ICML).

Song, Y., Ermon, S., 2019. Generative modeling by estimating gradients of the data distribution, in: Advances in Neural Information Processing Systems (NeurIPS).

Strassmann, S., 1986. Hairy brushes, in: Proceedings of the Conference on Computer Graphics and Interactive Techniques (SIGGRAPH).

Tashiro, Y., Song, J., Song, Y., Ermon, S., 2021. CSDI: conditional score-based diffusion models for probabilistic time series imputation, in: Advances in Neural Information Processing Systems (NeurIPS).

Tumanyan, N., Geyer, M., Bagon, S., Dekel, T., 2022. Plug-and-play diffusion features for text-driven image-to-image translation. arXiv preprint arXiv:2211.12572 .

Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L., Polosukhin, I., 2017. Attention is all you need, in: Advances in Neural Information Processing Systems (NeurIPS).

Voynov, A., Abernan, K., Cohen-Or, D., 2022. Sketch-guided text-to-image diffusion models. arXiv preprint arXiv:2211.13752 .

Wang, M., Wang, B., Fei, Y., Qian, K., Wang, W., Chen, J., Yong, J.H., 2014. Towards photo watercolorization with artistic verisimilitude. IEEE Transactions on Visualization and Computer Graphics .

Wei, C., Mangalam, K., Huang, P.Y., Li, Y., Fan, H., Xu, H., Wang, H., Xie, C., Yuille, A., Feichtenhofer, C., 2023. Diffusion models as masked autoencoders. arXiv preprint arXiv:2304.03283 .

Yu, L., Cheng, Y., Sohn, K., Lezama, J., Zhang, H., Chang, H., Hauptmann, A.G., Yang, M.H., Hao, Y., Essa, I., et al., 2023. Magvit: Masked generative video transformer.

Zeng, Y., Lin, Z., Zhang, J., Liu, Q., Collomosse, J., Kuen, J., Patel, V.M., 2022. Scenecomposer: Any-level semantic image synthesis. arXiv preprint arXiv:2211.11742 .

Zhang, L., Agrawala, M., 2023. Adding conditional control to text-to-image diffusion models. arXiv preprint arXiv:2302.05543 .

Zhang, Y., Wang, H., Shi, B.E., 2021. Gaze-controlled robot-assisted painting in virtual reality for upper-limb rehabilitation, in: International Conference of the IEEE Engineering in Medicine & Biology Society (EMBC).

Zhao, X., Ding, W., An, Y., Du, Y., Yu, T., Li, M., Tang, M., Wang, J., 2023. Fast segment anything. arXiv preprint arXiv:2306.12156 .

Zheng, H., Lin, Z., Lu, J., Cohen, S., Shechtman, E., Barnes, C., Zhang, J., Xu, N., Amirghodsi, S., Luo, J., 2022. Image inpainting with cascaded modulation gan and object-aware training, in: Proceedings of IEEE/CVF European Conference on Computer Vision (ECCV).

Zhu, P., Abdal, R., Qin, Y., Wonka, P., 2020. SEAN: image synthesis with semantic region-adaptive normalization, in: Proceedings of IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR).

Zou, Z., Shi, T., Qiu, S., Yuan, Y., Shi, Z., 2021. Stylized neural painting, in: Proceedings of IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR).

| Generated image | Segmented animal | Stroke representation | Generated image | Segmented animal | Stroke representation |
|---|---|---|---|---|---|



Fig. 6. Dataset generation pipeline. We show the images generated with Stable Diffusion Rombach et al. (2022), the segmentation obtained with Lüddecke and Ecker (2022), and the stroke representation Zou et al. (2021).

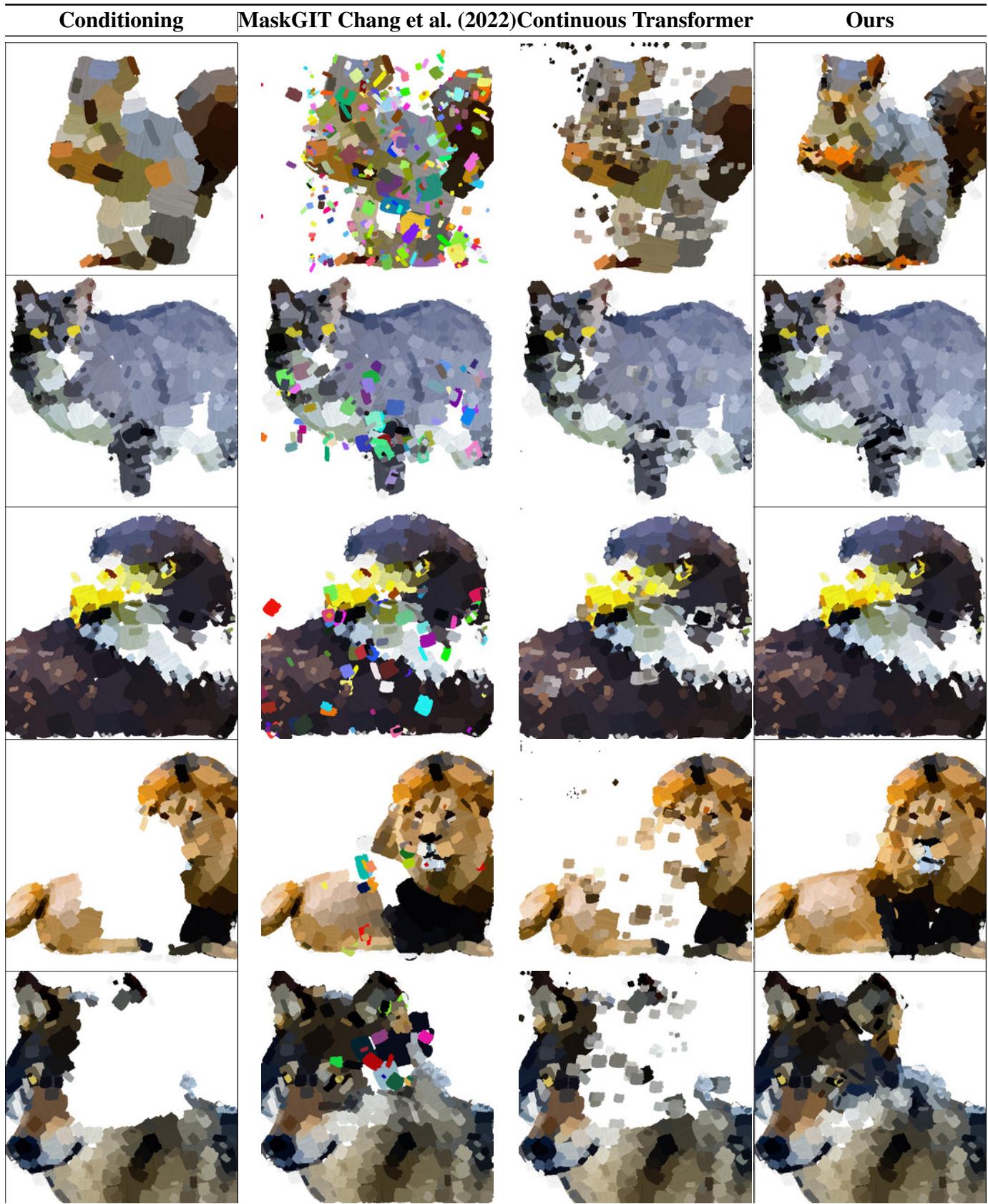| Conditioning | MaskGIT Chang et al. (2022) | Continuous Transformer | Ours |
| --- | --- | --- | --- |



Fig. 7. Comparison with baselines. Given the same conditioning sequence, we show the completion obtained with MaskGIT Chang et al. (2022), Continuous Transformer and our method.

**Fig. 8. Qualitative completion results of our method. Images are organized in blocks of 2. On the left the conditioning sequence s^{ctx}, on the right the conditioning with the completion added, *i.e.* $\tilde{s}_0$.**
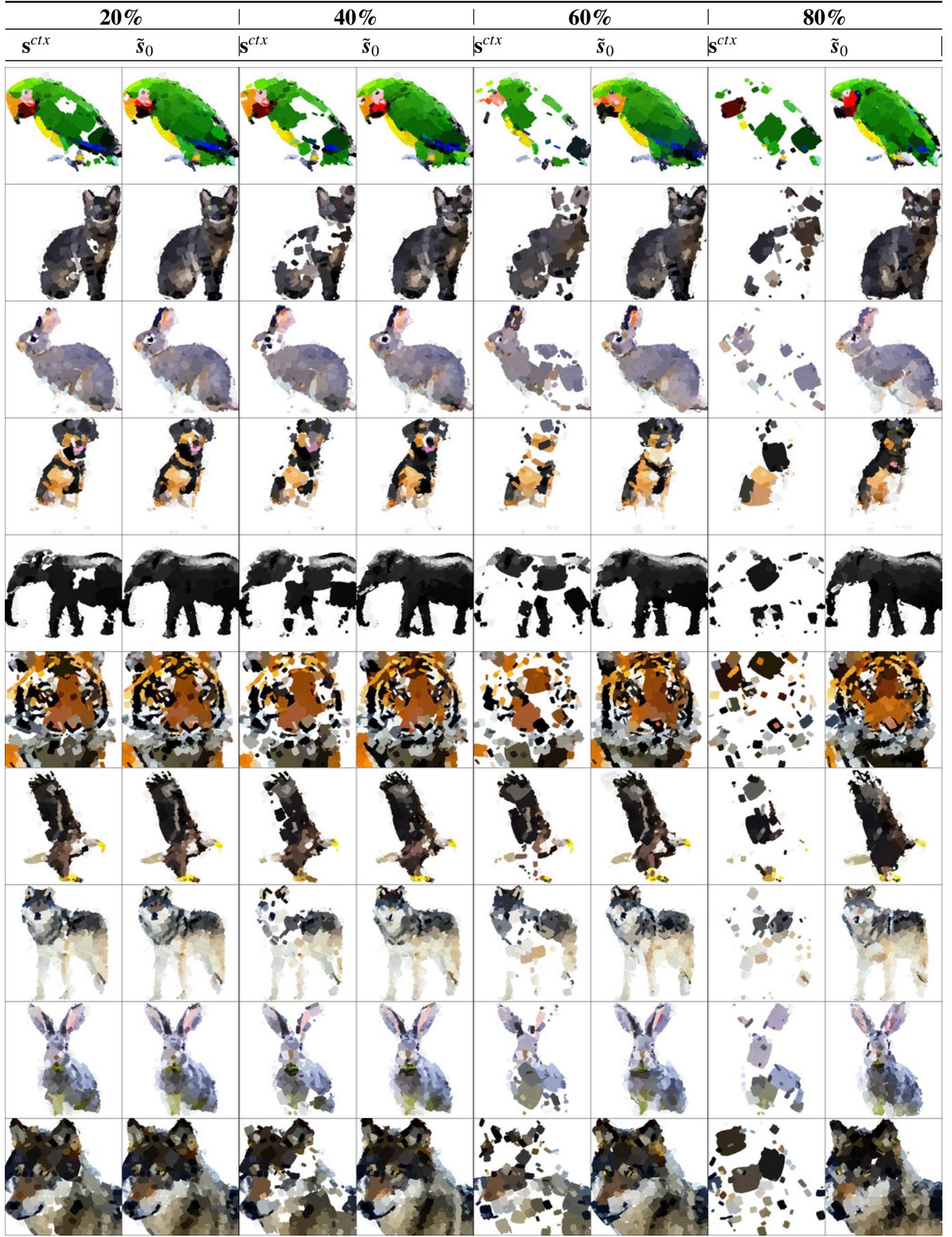
**Fig. 9. Robustness to random masking. Images are organized in blocks of 2. On the left the conditioning sequence $s^{ctx}$, on the right the conditioning with the completion added, *i.e.* $\tilde{s}_0$. The images are presented with an increasing percentage of strokes masked as conditioning from the left to the right.**

| Conditioning | Completion 1 | Completion 2 | Completion 3 | Completion 4 | Completion 5 |
|---|---|---|---|---|---|



Fig. 10. Diverse suggestions from the same context. On the left is the conditioning sequence $s^{ctx}$, and on the right are the diverse completions.

Fig. 11. Automatic inference without providing class. Images are organized in blocks of 2. On the left the conditioning sequence $\mathbf{s}^{ctx}$, on the right the conditioning with the completion added but no class provided.